SESSION VISUALIZATION, TOOLS AND TECHNIQUES

Chair(s)

TBA

Visualizing Camera Pose for Augmented Reality Applications

Andrew Risse, Dino Schweitzer, and Scott Teel

Department of Computer Science, United States Air Force Academy, Colorado, USA

Abstract - Augmented reality applications are becoming increasingly popular with technological advancements in mobile technology. Such applications rely on available sensor information to determine current position in the world, and direction in which the camera is pointing. This is known as the camera pose problem. The goal of this research was to experiment with the accuracy of the available sensors in common mobile devices and attempt to devise means to increase the accuracy. To visually understand the accuracy, a visualization tool, CameraPose, was developed to compare photographs taken with a mobile device to a drawn model of the world based on the estimated camera pose from sensor data. This paper describes the project and tool, demonstrates how the tool is used to determine accuracy, and proposes alternative approaches.

Keywords: Camera Pose, Augmented Reality, Visualization

1 Background

Augmented reality (AR) is the process by which computer stored information is actively displayed over real time images. In mobile applications, this is typically seen as information overlaid on the view from the camera. This information can be text, image, video, or sound data that is superimposed over the real world image. The two most common types of augmented reality applications today are AR browsers and AR applications that render virtual objects [1]. AR browsers use geopositioning data to display relevant location based information to the user. AR applications that render virtual objects are frequently used for gaming and entertainment purposes. One example would be the AR gaming experience delivered by Nintendo's 3DS. Augmented reality is commonly used for navigation, entertainment, and sightseeing, but its potential applications are many.

Today, AR technology is growing quickly in popularity as it becomes more interactive and applications are being developed in diverse areas. Augmented reality has recently been ported to mobile platforms. There are several open source systems that enable mobile AR development such as AndAR and Mixare [2,3]. Closed source options include Qualcomm's AR SDK "Vuforia" or Layar's "Reality Browser" and "Layar Player," [4,5]. On mobile platforms, augmented reality applications often rely on the accuracy of the on-board sensors. In order for a device to accurately display data integrated with a real-time image, the device must know where it is in the world and where it is facing. In order for the AR application to make it seem as if the augmented data actually exists in the real world, an accurate orientation of the device is crucial [6]. This is known as the "camera pose" problem and there are currently three common ways to solve it.

The first method of solving the "camera pose" problem is to use the on-board compass, accelerometer, gyroscope, and GPS sensors. The compass, accelerometer, and gyroscope determine the orientation of the device in the X, Y, Z axes. Many devices do not have a gyroscope and must default to using only the accelerometer, giving a less accurate reading. GPS data is utilized to determine where the device is physically located in the world. Combining orientation and device location give a camera pose which can then be used to augment the experience with relevant computer-generated This approach relies on the accuracy of the information. sensors which varies greatly between devices. GPS availability on mobile phones has been found to be near 100% outdoors, however, its accuracy has been found to be somewhat unreliable with positions up to 100 meters off [7]. Accelerometer sensors offer more reliability for orientation purposes, but magnetic compasses can be very inaccurate for certain devices.

Another common way of solving the camera pose problem is to utilize *fiducials*. Fiducials are markers placed in an image that are used as a point of reference or as a trigger to display some predefined data. Fiducials can look like QR codes and are often placed on a card or some sort of flat surface [9]. Layar's "Layar Vision" scans images and creates a "fingerprint" which is then recognized and used to deliver the corresponding AR experience. This "fingerprint" is essentially a fiducial. With fiducials, GPS data is irrelevant as visual recognition of the marker is all that is needed to deliver the corresponding information. Fiducials are commonly used to display virtual objects and allow the user to move around the virtual object in the real world while the object maintains its virtual position. Nintendo's 3DS uses fiducials to create a video game that is seen in the real world, for example, on the living room carpet or kitchen table [10]. In order for a system using fiducials to work, the designer must have some degree of control over the environment in which the application is running. The designer must be able to place or generate fiducials in the desired location. This method is not suitable for unprepared environments [6].

A third method to solve the camera pose problem is to perform natural feature recognition with the device's built in camera. Natural feature recognition is computationally expensive and has been found difficult to implement on mobile devices due to their limited computational power [11]. It is also possible to determine camera pose by combining natural feature recognition with on-board sensor data. Researchers at the Christian Doppler Laboratory used this method by first creating a 3D model of an environment. They then analyzed the mobile device's camera feed and performed natural feature tracking. By comparing what was seen by the camera and what was known from on-board device sensors, the device was able to calculate a full 6DOF camera pose [12].

2 The project

The purpose of our research was to investigate the accuracy of Android device on-board sensors when used to determine camera pose. We also wanted to devise and experiment with ways to enhance the accuracy of the pose estimation. The devices used were an HTC Incredible and a Samsung Galaxy Tablet 10.1.

Our approach was to first create a simple 3D model of the environment. For simplicity, we chose to model our local area - the cadet terrazzo and buildings at the United States Air Force Academy. We modeled the five prominent buildings in the area: Vandenberg Hall, the Chapel, Sijan Hall, Mitchell Hall, and Fairchild Hall. Using Google Earth, we obtained latitude and longitude coordinates for each corner of the buildings. Each building had a total of eight coordinates: four at ground level and four at roof level. An arbitrary origin point was chosen (a local statue on the cadet terrazo), and X, Y coordinates for the lat/lon value were determined by computing distances to the origin point. Building heights were obtained from building schematics.

In addition to the 3D model, we needed real world data captured by the device. To do this, we developed an Android OS application that brings up the camera view and dynamically displays the real-time data for device orientation and GPS position. GPS data is captured using the GPS_PROVIDER, an on-board GPS chip as opposed to network based GPS. GPS information from the GPS chip is slower, but more accurate than the data received from the NETWORK_PROVIDER. The orientation data is labeled as yaw, pitch, and roll. Yaw corresponds to the compass heading of the device where the compass vector is pointing out of the back side camera and away from the user. To account for magnetic declination, a 9 degree offset was used for our geographic location. Pitch describes the elevation of the device where 0° occurs when the device is lying face up on a

flat surface and 90° when it is upright in portrait configuration. Roll is measured from the bottom left hand corner of the screen and changes as the device rotates around the axis coming out of the camera. 0° roll is when the device is right side up and 180° would be upsidedown.

The app has a "Take Picture" button that, when pressed, writes the current orientation data and GPS location to a text file and the picture taken is saved. To compare the collected data to the 3D model, we wrote a visualization tool called "CameraPose".

2.1 CameraPose

CameraPose is a Java-based application to compare camera pose information captured from a mobile device to a computer-generated display of what the scene "should" look like. Figure 1 shows a sample screen shot. The tool consists of the following four panes:

- Menu pane (upper left) for manually setting map scale, camera orientation, and camera FOV
- 2D map (lower left) showing buildings, camera location, and camera direction
- Camera picture (lower right) showing the captured image from the mobile device
- Model view (upper right) of the 3D scene based on the captured camera orientation

CameraPose reads in the model of the world and displays it in the 2D map pane which is scalable and can be panned. The camera can be dragged to any location. The orientation of the camera is set by changing the Yaw, Pitch, Roll, and Height sliders in the menu pane. The camera's field of view (FOV) can be adjusted in both the horizontal and vertical direction. The 3D wireframe scene of the model as viewed from the camera's location and orientation is shown in the model view pane.

To compare the model view with the captured data, CameraPose reads in a text file of captured data points with their associated camera image. The user selects an image, and the application adjusts the camera location and orientation to the values in the data file, regenerating the model view. In theory, if the data is accurate, the wireframe model view should match the camera picture.



Figure 1. CameraPose screen shot.



Figure 2. Overlaying camera picture with model view.



Figure 3. Adjusting parameters to match images.

2.2 Determining accuracy

To ensure the model view was a true representation of what the camera "should" see, the FOV for tested devices were physically determined by measuring widths and heights of images at known distances. The FOV slider values are set appropriately for the device and used in the perspective transformation. Care is taken to make sure the correct aspect ratio is used when displaying the picture in the smaller camera picture pane.

Overlaying the model view and picture allows a more accurate comparison of the two images (Figure 2). In this example, the top of the building matches fairly accurately, while the bottom and right front corner of the building are obviously off. The camera parameter sliders and camera position can be adjusted in the tool until the two images match exactly. At this point, the difference between the captured data values and adjusted parameters represents the error term in the captured data. Figure 3 shows an example of an original overlay with the same overlay after adjustments. In this example, yaw, pitch, and camera location were adjusted. Note the left building (Mitchell Hall) is modeled from the overhang at the top.

3 Results

The Samsung tablet was significantly more accurate than the HTC Incredible. Table 1 shows the results of applying the above image registration process to over 20 sample images for the HTC and 19 for the tablet. For orientation parameters (yaw, pitch, roll), the average error is shown in degrees. For camera position, the error is shown as distance in feet. GPS coordinates were not recorded for the Samsung tablet. Rather, pictures were taken from the same locations as the HTC data.

Table 1. Average error terms for devices.

| Device | Yaw | Pitch | Roll | Position |
|--------|-------|-------|------|----------|
| HTC | 19.99 | 2.4 | 1.1 | 16.6 |
| Tablet | 11.0 | 1.3 | 3.3 | N/A |

The largest error factor is in the yaw parameter, or camera direction. For the HTC Incredible, two outliers, (a 63 and 88 degree error) resulted in the large 20 degree average error. Without those outliers, the average error was still 14.7 degrees which is significant for an augmented reality application. If you were looking at a point, such as a building, 100 feet away, this error term represents 25 feet on either side of the point. This is a result of the inherent inaccuracy of the mobile magnetic compass.

The pitch and roll angles, which rely on the built-in accelerometers were much more accurate than the compass. The GPS accuracy of 16 feet is consistent with reported measurements in the literature.

3.1 Attempts to improve accuracy

The unreliability of compass readings in Android devices is a known problem [1]. In order to get accurate readings, the device must frequently be calibrated. This is done by either moving the device in a figure eight pattern, or by rotating the device repeatedly over all three axes. Even when calibrated, if the device is subject to sudden extreme or jerky movements the data jumps wildly and provides inaccurate information. Close proximity of metal objects also distorts the compass data.

Because of this inaccuracy, we wanted to find a better way of determining the yaw vector. Our first approach was to analyze recent GPS positions to try and generate a "direction moving" vector. The idea is to capture recent GPS coordinates as the user moves and use them to identify the line of travel. Assuming the user was moving in the desired direction of the taken image, this line of travel would then be considered the yaw vector. We experimented with this approach to determine feasibility, how many GPS coordinates were required to establish a stable "look" vector, and how accurate this approach was compared to using the compass.

The Android data capture tool was modified to keep a circular buffer of GPS coordinates that was constantly GPS coordinates were captured whenever they updated. changed. When the user selected the "Take Picture" button, the previous 100 GPS coordinate were dumped to the text file along with the other sensor data and image. The CameraPose application was modified to display the "trail" of GPS points on the 2D map leading up to the picture (Figure 4). The vector represents the captured compass direction. The true direction was true north (up). Unfortunately, as evident in the picture, the distance between captured points were too separated and scattered to form an accurate direction vector, unless the user was moving in the exact direction of the desired image for several feet. In addition, the relative error in each measurement did not provide sufficient consistency to get an accurate direction. This approach was deemed impractical for estimating an accurate direction vector.



Figure 4. Trail of GPS coordinates prior to picture.

Our second attempt to improve accuracy was to see if taking an average of the compass readings would yield a better result. After plotting 100 sequential readings from the devices' compasses in each of the cardinal directions, we found that readings were generally centralized around a common value. This value however, was still significantly displaced from the true heading of the device. While taking an average might help reduce some of the noise within the compass readings, the end result would still be an unacceptably inaccurate heading. The graphs below (Figures 5 and 6) show the compass reading plotted over time. It is interesting to note that the HTC Incredible's compass would pause on certain values then begin to change, then pause again. The Samsung tablet's readings were constantly changing.



Figure 5. 100 sequential HTC Incredible compass readings, true heading 180°, not accounting for declination.



Figure 6. 100 Sequential Samsung tablet compass readings, true heading 180°, not accounting for declination

4 Future research areas

Future research for this project could include looking at the application of Kalman filters for achieving a more accurate compass reading. When the user selects the button to take a picture, a Kalman is applied to the last n compass readings to account for the error term and attempt to increase the accuracy. Kalman filter approaches have been used successfully in robotics and autonomous vehicle navigation. Another avenue would be to examine the accuracy when combining the gyroscope and compass sensors to achieve a faster and steadier reading, as has been done in applications such as "Steady Compass." A third approach might be to experiment with methods to reduce electromagnetic interference, such as leaving the device in "airplane mode".

5 Conclusions

Augmented reality is an increasingly popular area for mobile applications. The effectiveness of an AR approach relies heavily on the accuracy of the location and orientation of the mobile device. Current sensor technology limits the ability to rely on internal sensors, and applications that require a high degree of accuracy must use image processing techniques with external registration. Approaches such as image registration with a 3D model or known fiducials in a scene provide greater accuracy, but can only be used in known environments.

The CameraPose application provides a visual approach to "seeing" the accuracy of estimated camera pose in a known environment. It can be used to compare different estimation techniques without requiring high quality sensors. It also provides a manual method for measuring the accuracy of the estimation.

6 References

- Sarmento, A., Amor, M., Padron, E, and Regueiro, C. "An Analysis of Android Smartphones as a Platform for Augmented Reality Games"; Proceedings of UBICOMM 2011: The Fifth International Conference on Mobile Ubiquitous Computing, Systems, Services and Technologies (Nov 2011).
- [2] <u>http://code.google.com/p/andar/</u>
- [3] http://www.mixare.org/
- [4] http://www.qualcomm.com/solutions/augmented-reality
- [5] http://www.layar.com/
- [6] Azuma, R., et. al. "Recent Advances in Augmented Reality". IEEE Computer Graphics and Applications 21, 6, pp. 24-47. (2001)
- [7] Zandbergen, P. and Barbeau, S. "Positional Accuracy of Assisted GPS Data from High-Sensitivity GPS-enabled Mobile Phones". The Journal of Navigation. Vol. 64 No. 3. pp. 381-399, (2011).
- [8] Perey, C. "Print and publishing and the future of Augmented Reality". Journal of Information Services and Use. Vol 31, No. 1-2. pp 31-38, (2011).
- [9] http://www.layar.com/documentation/browser/howtos/ layar-vision-doc/
- [10] http://www.nintendo.com/3ds/built-in-software/#/4
- [11] Wagner, D., et al « Pose Tracking from Natural Features on Mobile Phones ». In Proceedings of the 7th IEEE/ACM International Symposium on Mixed and Augmented Reality (ISMAR -08)I. pp. 125-134, (2008).
- [12]http://handheldar.icg.tugraz.at/localization.php

Volumetric Intravascular Ultrasound Visualization Using Shape-based Nonlinear Interpolation

Y. Rim¹, D. D. McPherson¹, and H. Kim¹

¹Department of Internal Medicine, The University of Texas Health Science Center at Houston, Houston, Texas, USA

Abstract - Intravascular ultrasound (IVUS) has been utilized primarily to evaluate plaque formation of the coronary and carotid arteries. Several studies have reported three-dimensional (3D) IVUS image reconstruction algorithms to better provide anatomical information of the arterial wall structure. However, most of previous 3D IVUS studies provided unrealistic morphology and acoustic information in the arterial wall due to the low quality of input image data and the linear interpolation algorithm. In the present study, we proposed an improved algorithm to generate intermediary slices using a shape-based nonlinear interpolation method, and created volumetric 3D IVUS images using radio frequency (RF) raw IVUS signal data of a porcine femoral artery. Both arterial structure and acoustic intensity information of the intermediary slice were created by the cubic spline interpolation. This novel volumetric 3D IVUS visualization strategy has the potential to improve vascular ultrasound imaging for better determination of atheroma distribution in the arterial structure.

Keywords: Intravascular ultrasound (IVUS), Shape-based nonlinear interpolation, Volumetric visualization, Threedimensional ultrasound, Intermediary slice, Acoustic intensity

1 Introduction

Clinical procedure of the intravascular ultrasound (IVUS) has been utilized to identify the pathological alterations of atherosclerotic plaque by evaluating acoustic intensity information of the cross-sectional images of the artery [1-3].

Comparison of the relatively proximal vascular segment to distal segments in the standard clinical procedure of IVUS requires a repeated review of sequentially recorded cross-sectional two-dimensional (2D) IVUS images to determine the spatial relation of the regions of interest (ROI) [1]. Therefore, it is often difficult for clinicians to determine the type and morphology of atherosclerotic plaque or lesion within the arterial wall [4]. IVUS imaging can demonstrate a series of cross-sectional images of the arterial wall structure in the longitudinal direction using a pull-back device with IVUS catheter. This sequential imaging of 2D IVUS images can provide threedimensional (3D) IVUS visualization.

Two methodologies have been commonly utilized to create 3D IVUS imaging. A popular method is simply to generate a longitudinal cut-view of the artery showing acoustic intensity information of the arterial wall along the blood flow direction [5]. However, this method only creates another 2D image inside the arterial wall along the longitudinal direction, and thus cannot provide comprehensive information pertaining to spatial plaque distribution over the entire arterial structure. The other popular technique for 3D IVUS is to focus on creating anatomically realistic arterial border of the lumen and media-adventitia contour using a smooth 3D surface reconstruction algorithm [6]. But detailed information of acoustic intensity distribution through the arterial wall structure is discarded in this technique once surface rendering is performed.

Several studies have reported three-dimensional (3D) IVUS image reconstruction algorithms to better provide anatomical information of the arterial wall structure and overcome the limited 3D visualization capability of currently available commercialized IVUS systems [1, 6-8]. An early study demonstrated a computer-automated 3D reconstruction method to generate a tangible format with a series of 2D IVUS images using linear interpolation method, and compared the reconstructed 3D images to the sequential images obtained during IVUS examination [1]. Another study introduced a shape-based interpolation method of multi-dimensional grayscale images to create coarser and finer discretized images to combine images of the same ROI from two independent image modalities [9]. Although previous studies on 3D IVUS reconstruction presented 3D images of arterial structure, most of the 3D images had unrealistic morphology and acoustic information in the arterial wall due to the linear interpolation to create intermediary slices and the low quality of input image data from video tapes.

We have recently proposed a volumetric 3D IVUS visualization methodology for early and inflammatory arterial atheroma characterization [10]. Volumetric 3D images can provide both visualization of the 3D arterial structure and acoustic intensity distribution for plaque and atheroma detection. However, linear interpolation was utilized in the volumetric 3D IVUS visualization. In order to

quantitate the geometric characteristics of plaque volume and atheroma formation, an improved volumetric 3D IVUS visualization method is required to more distinctly demonstrate acoustic intensity information overlaid on a 3D structure image of the artery.

In the present study, we proposed an improved algorithm to generate intermediary slices between adjacent slices using the shape-based nonlinear interpolation method and create a volumetric 3D IVUS image using raw radio frequency (RF) IVUS signal data of a porcine femoral artery. Since the distance between neighboring pixels in a crosssectional image is usually smaller than the distance between the images leading to non-isotropic voxel dimensionality between neighboring pixels, this often yields deterioration in the quality of the 3D image [11]. Therefore, increasing the level of discretization between the acquired 2D IVUS images may provide improved quality of the consequent volumetric 3D IVUS image of the arterial structure. We compared the improved volumetric 3D IVUS images created by the proposed shape-based nonlinear interpolation method with those created by the linear interpolation in previous studies.

2 Materials and Methods

We developed a straightforward algorithm to convert a series of raw RF IVUS signal data into a fully volumetric 3D visualization (Fig. 1). The entire protocol including 2D image reconstruction from raw RF signal data, border tracing, segmentation, sequential alignment of 2D images, and intermediary slice generation was conducted in a single image processing platform. Volumetric visualization of 3D IVUS images was performed using ImageJ, an open-source Java-based image processing software provided by the National Institutes of Health (NIH).

2.1 Volumetric 3D IVUS Visualization

An atherosclerotic Yucatan miniswine atheroma model (20 kg, Sinclair Research Center Inc., Columbia, MO) was used. The animal protocol was approved by the Institutional Animal Care and Use Committee of The University of Texas Health Science Center at Houston. Following full anesthesia, the right femoral artery was exposed with groin incisions, and an arteriotomy was performed. A 5F sheath was inserted in the femoral artery. A high frequency (20 MHz, 3.5F) IVUS imaging catheter was utilized connected to a Volcano s5i IVUS Imaging System (Volcano Co, Rancho Cordova, CA). The IVUS catheter was inserted through the arterial segment past the region of interest. The IVUS catheter was withdrawn using an automatic pullback device at a constant speed of 0.5 mm/s, while IVUS images and raw RF signal data of the arterial segment were continuously recorded. A total of 256 scan lines with 1,024 sampling data per each scan line were recorded (dynamic range of 40-60 dB). Since there was no curvature in the arterial segment evaluated, it was assumed that the direction of pullback of the IVUS catheter was parallel to the longitudinal direction of the artery. Electrocardiogram (ECG) was synchronized with IVUS imaging in real time thus serving as a time reference for systolic and diastolic phases while recording the RF IVUS signal data. In order to construct 2D grayscale images from the raw RF signal in beam space, the pre-determined cutoff threshold (1,050 mV) was applied to the RF signal data to create 2D IVUS images comparable to that directly generated from the Volcano IVUS system. The enveloped amplitude (i.e., acoustic intensity) under the threshold value was utilized to reconstruct grayscale images in beam space. In this beam space, x- and y- axes refer to radial and circumferential directions, respectively.

The reconstructed image in beam space was transformed to the Cartesian coordinate system for standard vascular imaging. A graphical user interface (GUI)-based image processing system was developed for interactive tracing and segmenting procedure under MATLAB (Mathworks Inc., Natick, MA) platform.



Figure 1. Protocol of the volumetric 3D IVUS visualization

The endothelium/atheroma border and the outer edge of the dense adventitia in each image were manually segmented, and a series of segmented RF data set in the ROIs were placed in tomographic sequence for intermediary slice generation. The extracted RF signal data of the ROIs were utilized to generate intermediary slices using the shape-based nonlinear interpolation method.

2.2 Shape-based Nonlinear Interpolation

Fig. 2 demonstrates the algorithm to generate intermediary slices between 2D IVUS image slices (collected and segmented from original raw RF data) using the shape-based nonlinear interpolation. The basic concept in this shape-based nonlinear interpolation method is to interpolate the vascular structure geometry using three neighboring slices (Step 1, Fig. 2). We first applied the cubic spline interpolation method to obtain the segmented ROI of the arterial wall along the longitudinal direction using the traced 20 cross-sectional IVUS image data. The interpolation was performed on the 256 scan lines along the circumferential direction between original slices, and the boundary information was calculated along the same 256 scan lines in the intermediary slices. Next, nonlinear acoustic intensity interpolation was performed (Step 2, Fig. 2). The segmented ROI of the arterial wall in each slice contains acoustic intensity distribution profile. This information was utilized to interpolate acoustic intensity distribution within the ROI in the consequent intermediary slice using the cubic spline interpolation.

We collected RF data set of the segmented ROIs in 20 IVUS images with a distance of 0.5 mm between images to generate 5 intermediary slice images between adjacent images resulting in a total of 115 cross-sectional images along the longitudinal direction. Acoustic intensity information in the 2D IVUS image data was preserved in the RF data of the segmented ROIs in each slice. In order to interpolate RF data (i.e. acoustic intensity distribution) within the segmented ROIs between 20 slices in a 3D space, we utilized the cubic spline interpolation considering the nonlinearity of the vascular structure geometry and acoustic intensity in the arterial wall. The essential idea of the interpolation is to fit a piecewise function of the form

$$S(x) = \begin{cases} s_{1}(x) & \text{if} \quad x_{1} \le x \le x_{2} \\ s_{2}(x) & \text{if} \quad x_{2} \le x \le x_{3} \\ \vdots & \vdots \\ s_{n-1}(x) & \text{if} \quad x_{n-1} \le x \le x_{n} \end{cases}$$
(1)

where s_i is a third degree polynomial defined by

$$s_i(x) = a_i(x - x_i)^3 + b_i(x - x_i)^2 + c_i(x - x_i) + d_i \quad (2)$$

for $i = 1, 2, ..., n-1$.

The first and second derivatives of these n-1 equations are fundamental to this process, which are

$$s'_{i}(x) = 3a_{i}(x - x_{i})^{2} + 2b_{i}(x - x_{i}) + c_{i}$$
 (3)

$$s_i^{"}(x) = 6a_i(x - x_i) + 2b_i$$
 (4)

for *i* = 1, 2, ..., *n*-1.

Using the following four stipulations, we can determine the weights for n-1 equations.

- 1. The piecewise function S(x) will interpolate data points.
- 2. *S* (*x*) will be continuous on the interval $[x_1, x_n]$.
- 3. S(x) will be continuous on the interval $[x_1, x_n]$.
- 4. S'(x) will be continuous on the interval $[x_1, x_n]$.



Step1. Nonlinear geometry interpolation



Step 2. Nonlinear acoustic intensity interpolation



Shape-based nonlinear interpolation

Figure 2. Intermediary slice generation using the shapebased nonlinear interpolation

Detailed information of the shape-based nonlinear interpolation procedure is described in Fig. 3. The first step is to register the original 2D IVUS slices in the global coordinate, and align the centroids of the arterial structure on each slice to a centerline. Considering the distance between slices, the interpolated curvatures were created on each of the 256 scan lines by means of the cubic spline interpolation method. Radial distances $r(x_1)$, $r(x_2)$ and $r(x_3)$ on the slice #1, #2, and #3, respectively, were utilized to calculate r(x) on the interpolated intermediary slice. Both inner and outer boundary points of the arterial wall structure on each scan line were generated on the intermediary slices. Since there were a varying number of data points between these two boundary points of the arterial wall structure on each scan line, it is necessary to create the same number of acoustic intensity data points within the ROI along each scan line. Therefore, the acoustic intensity data within the ROI on each scan line were resampled with 100 data points allowing generation of the same number of discretized acoustic intensity data points in the corresponding ROI on the intermediary slices. Nonlinear interpolation of the acoustic intensity values at 100 data points on 256 scan lines was performed to complete the shape-based nonlinear interpolation for both vascular structure geometry and acoustic intensity information.



Figure 3. Schematic diagram of the shape-based nonlinear interpolation

3 Results

3.1 Intermediary Slice Generation Using the Shape-based Nonlinear Interpolation

Fig. 4 shows the original slices of the porcine femoral artery from the raw RF IVUS data and the consequent intermediary slices created by using the developed shape-based nonlinear (left) and linear interpolations (right). The

first, third and fourth row images indicate three original slices, and the second row images demonstrate the intermediary slices between slice #1 and #2 using two different interpolation methods.

The linear interpolation method demonstrated a blur intermediary image calculated by a simple averaging function between the two adjacent slices for both segmentation and acoustic intensity information resulting in an unrealistic image of the arterial intermediary slice. The shape-based nonlinear interpolation method developed in this study, however, provided a clearly reconstructed intermediary slice image with accurately predicted segmentation and acoustic intensity information calculated based on the corresponding information of the three neighboring original image slices. The shadow effect in the linear interpolation method was not observed in the shapebased nonlinear interpolation method. Shape-based nonlinear interpolation method can, therefore, provide better interpolation outcome and generate more realistic geometry of the arterial segment with accurate acoustic intensity values.

| Slice #1 | | \bigcirc |
|-----------------------|-----------------|--------------------------|
| Intermediary slice | \bigcirc | |
| Slice #2 | | |
| Slice #3 | \bigcirc | \bigcirc |
| | (A) Shape-based | (B) Linear interpolation |

Figure 4. Intermediary slice generation using two different interpolation methods

3.2 Volumetric 3D IVUS Visualization

Volumetric 3D IVUS images reconstructed by the developed shape-based nonlinear interpolation and the conventional linear interpolation methodologies are demonstrated in Fig. 5. Both volumetric 3D reconstruction models were created utilizing the raw RF IVUS signal data from a total of 20 slices, and presented along the longitudinal direction to facilitate 3D visualization. It is clearly observed that the shape-based nonlinear interpolation method provided better demonstration of the 3D structure of the arterial segment and more realistic acoustic intensity distribution compared to the linear interpolation method.





Full volume

(A) Shape-based nonlinear interpolation



Half-cut



Full volume

(B) Linear interpolation

Figure 5. Volumetric 3D IVUS visualization using the shape-based nonlinear and linear interpolation methods

In particular, the half-cut view of the 3D IVUS image created by the shape-based nonlinear interpolation provided excellent information with smooth acoustic intensity distribution over the luminal surface of the arterial wall where plaques are usually observed (Fig. 5A top). The full volume view of the arterial segment well described the boundary surface of the outer edge of the dense adventitia with corresponding acoustic intensity distribution along the longitudinal direction (Fig. 5A bottom).

On the other hand, the linear interpolation method poorly demonstrated the vascular structure geometry with noticeable discontinuity between slices and striped acoustic intensity distribution in both half-cut and full volume images (Fig. 5B). Due to the blur images with gray shaded areas in the intermediary slices, the morphology of lesion along the longitudinal as well as circumferential directions were not correctly demonstrated leading to unrealistic volumetric visualization of the arterial structure.

4 Discussion

Conventional cross-sectional images of an arterial structure from most of the commercialized IVUS imaging systems can hardly provide comprehensive information pertaining to complex spatial distribution of lesions such as atherosclerosis in a 3D space. Volumetric 3D IVUS visualization of the arterial structure can provide a powerful tool to overcome this limited spatial demonstration issue with 2D IVUS alone. Visualization of both 3D morphology of the arterial structure and corresponding acoustic intensity distribution within the arterial segment can help better understand the extent and stage of plaque formation.

Previous studies have established 3D visualization methods of the vascular structure using 2D IVUS images captured from analog video tapes [1, 6-8]. However, the volumetric 3D reconstruction in these studies demonstrated images with a poorly low resolution as non-isotropic voxel dimensionality was used due to the smaller distance between neighboring pixels than the distance between slice images along the longitudinal direction (blood flow direction). In general, the longitudinal resolution of commercial IVUS imaging systems is much lower than that of the 2D crosssectional IVUS images. This often yields deterioration in image quality. In order to provide better visualization of volumetric 3D IVUS images, it is imperative to improve interpolation techniques for intermediary slice generation. In the present study, we developed a novel volumetric 3D IVUS visualization strategy to create intermediary slices between original IVUS image slices using the shape-base nonlinear interpolation method.

An important aspect of the volumetric 3D IVUS image reconstruction is to improve image quality of the 3D visualization of the arterial structure. Poor image quality may reduce the accuracy of 3D quantitation and visualization of plaque formation. In most of previous studies, low image resolution has hampered 3D IVUS visualization. In this study, we utilized raw RF IVUS signal data with high resolution to generate intermediary slices between the IVUS image data. The shape-based nonlinear interpolation method using the cubic spline algorithm was proposed to create vascular structure geometry and acoustic intensity distribution information in intermediary slices between neighboring 2D IVUS slice images. The intermediary slices clearly demonstrated accurately predicted segmentation and acoustic intensity information of the arterial segment. Volumetric 3D IVUS reconstruction with the shape-based nonlinear interpolation provided better visualization of morphology of lesion in the arterial wall with realistic acoustic intensity distribution. In particular, the longitudinal half-cut view of the arterial structure demonstrated excellent continuity between original IVUS slices with respect to both geometry and acoustic intensity information.

There are some limitations in the present study. Manual tracing was performed to segment the endothelium/atheroma border and outer edge of the dense adventitia on the cross-sectional 2D IVUS images along the artery. We are developing a semi-automated image segmentation algorithm for more accurate volumetric 3D visualization and quantitation. In addition, we were not able to perfectly adjust the rotational offset of IVUS catheter probe which can occur during pullback recording. Volumetric 3D IVUS image can be affected by this rotational offset effect.

In summary, we have successfully developed an improved 3D reconstruction algorithm using the shapebased nonlinear interpolation method, and performed volumetric 3D IVUS visualization of a porcine artery. A superiority of the shape-based nonlinear interpolation method over the conventional linear interpolation algorithm in terms of the quality of volumetric 3D IVUS visualization was clearly demonstrated. This novel volumetric 3D IVUS visualization strategy has the potential to improve vascular ultrasound imaging for better determination of atheroma distribution in the arterial structure. Moreover, precise volumetric 3D visualization with accurate acoustic intensity information may improve advanced molecular ultrasound imaging of atheroma components.

5 References

[1] K. Rosenfield, D. W. Losordo, K. Ramaswamy, J. O. Pastore, R. E. Langevin, S. Razvi, B. D. Kosowsky, and J. M. Isner, "Three-dimensional reconstruction of human coronary and peripheral arteries from images recorded during two-dimensional intravascular ultrasound examination," *Circulation*, vol.84(5), pp.1938-1956, Nov, 1991.

[2] S. E. Nissen, and P. Yock, "Intravascular ultrasound: novel pathophysiological insights and current clinical applications," *Circulation*, vol.103(4), pp.604-616, Jan 30, 2001.

[3] D. E. Goertz, M. E. Frijlink, N. de Jong, and A. F. van der Steen, "Nonlinear intravascular ultrasound contrast imaging," *Ultrasound Med Biol*, vol.32(4), pp.491-502, Apr, 2006.

[4] A. G. Bors, L. Kechagias, and I. Pitas, "Binary morphological shape-based interpolation applied to 3-D tooth reconstruction," *IEEE Trans Med Imaging*, vol.21(2), pp.100-108, Feb, 2002.

[5] F. Prati, E. Arbustini, A. Labellarte, L. Sommariva, T. Pawlowski, A. Manzoli, A. Pagano, M. Motolese, and A. Boccanelli, "Eccentric atherosclerotic plaques with positive remodelling have a pericardial distribution: a permissive role of epicardial fat? A three-dimensional intravascular ultrasound study of left anterior descending artery lesions," *Eur Heart J*, vol.24(4), pp.329-336, Feb, 2003.

[6] R. Sanz-Requena, D. Moratal, D. R. Garcia-Sanchez, V. Bodi, J. J. Rieta, and J. M. Sanchis, "Automatic segmentation and 3D reconstruction of intravascular ultrasound images for a fast preliminar evaluation of vessel pathologies," *Comput Med Imaging Graph*, vol.31(2), pp.71-80, Mar, 2007.

[7] J. G. van den Broek, C. H. Slump, C. J. Storm, A. C. van Benthem, and B. Buis, "Three-dimensional densitometric reconstruction and visualization of stenosed coronary artery segments," *Comput Med Imaging Graph*, vol.19(2), pp.207-217, Mar-Apr, 1995.

[8] J. L. Evans, K. H. Ng, S. G. Wiet, M. J. Vonesh, W. B. Burns, M. G. Radvany, B. J. Kane, C. J. Davidson, S. I. Roth, B. L. Kramer, S. N. Meyers, and D. D. McPherson, "Accurate three-dimensional reconstruction of intravascular ultrasound data. Spatially correct three-dimensional reconstructions," *Circulation*, vol.93(3), pp.567-576, Feb 1, 1996.

[9] G. J. Grevera, and J. K. Udupa, "Shape-based interpolation of multidimensional grey-level images," *IEEE Trans Med Imaging*, vol.15(6), pp.881-892, 1996.

[10] H. Kim, M. R. Moody, S. T. Laing, P. H. Kee, S. L. Huang, M. E. Klegerman, and D. D. McPherson, "In vivo volumetric intravascular ultrasound visualization of early/inflammatory arterial atheroma using targeted echogenic immunoliposomes," *Invest Radiol*, vol.45(10), pp.685-691, Oct, 2010.

[11] G. P. Penney, J. A. Schnabel, D. Rueckert, M. A. Viergever, and W. J. Niessen, "Registration-based interpolation," *IEEE Trans Med Imaging*, vol.23(7), pp.922-926, Jul, 2004.

Spatial Data Structures, Sorting and GPU Parallelism for Situated-agent Simulation and Visualisation

A.V. Husselmann and K.A. Hawick

Computer Science, Institute for Information and Mathematical Sciences, Massey University, North Shore 102-904, Auckland, New Zealand email: { a.v.husselmann, k.a.hawick }@massey.ac.nz Tel: +64 9 414 0800 Fax: +64 9 441 8181

March 2012

ABSTRACT

Spatial data partitioning techniques are important for obtaining fast and efficient simulations of N-Body particle and spatial agent based models where they considerably reduce redundant entity interaction computation times. Highly parallel techniques based on concurrent threading can be deployed to further speed up such simulations. We study the use of GPU accelerators and highly data parallel techniques which require more complex organisation of spatial datastructures and also sorting techniques to make best use of GPU capabilities. We report on a multiple-GPU (mGPU) solution to grid-boxing for accelerating interaction-based models. Our system is able to both simulate and also graphically render in excess of $10^5 - 10^6$ agents on desktop hardware in interactive-time.

KEY WORDS

grid-boxing; sorting; GPU; thread concurrency; data parallelism

1 Introduction

Simulation problems involving interacting spatially located entities such as N-Body particle models [1–4] or spatial agent-based models often require considerable computational power to support simulation of adequately large scale systems. Large size is important to expose and investigate complex and emergent phenomena that often only appear on logarithmic length and systems size scales. To reduce the $\mathcal{O}(N^2)$ computational complexity of such interacting systems, spatial partitioning methods are often used. Depending upon how effective the scheme used is, the computational complexity can be reduced down to $\mathcal{O}(N \log N)$. Fur-



Figure 1: A visualisation of a uniform grid datastructure as generated by GPU, using NVIDIA's CUDA.

thermore, a good scheme will also support - and not impede - the introduction of parallelism into the computation so that an appropriate parallel computer architecture [5] can be used to further reduce compute time per simulation step. As a vehicle, we use the Boids model originally by Reynolds [6,7].

Spatial partitioning [8–10] has been employed as far back as the earliest attempts of Andrew Appel [11] who produced the first multipole partitioning method, with many subsequent developments [12]. The N-Body particle simulation has served as a simple but effective testing platform for improving the performance of particle simulations and generally interaction-based models which rely on particle-particle interactions [13]. The vast majority of particle simulations find themselves suffering an inherent performance problem and a severe lack of system size scaling. This is mostly due to particle interaction with every other particle inducing an $\mathcal{O}(N^2)$ complexity. Even for a conceptually simple in-



(a) Screenshot of our Boids sim- (b) Screenshot of the same simulation. ulation state with visualisation of the grid-box algorithm superposed.

Figure 2: Visualisation of the same set of agents with and without a uniform grid.

teraction such as a nearest neighbour (NN) interaction, it is still necessary, in the simplest cases, to perform calculations (such as Euclidean distance) for every combination of two particles. Collision detection is one such problem.

This problem is mitigated by using methods such as Barnes-Hut [14, 15] treecodes, Fast Multipole methods [16], and tree organisational methods [17, 18] including Octrees [19-26], K-D trees [27, 28], Kautz trees [29], Adaptive Refined Trees [30], and more, all of which fall under the category of spatial partitioning. Not all of these algorithms would necessarily be suitable for any one particle simulation. Some of these algorithms are suitable for simulations which frequently have high density in one or several spatial locations. Some perform best with uniformly distributed particles, and some attempt to systematically treat several particles in the same light as a point mass to reduce computation in exchange for losing a small amount of precision. This technique is often well-suited to applications in astrophysics. There also exists algorithms which sacrifice precision for performance, with techniques as simple as updating only a certain number of agents per simulation time step, which dramatically increases performance but the loss of precision is so great it will no longer reflect the original simulation code. Choosing which one of these algorithms or methods to use in a simulation therefore ultimately depends on the purpose of the simulation.

This work contributes towards efficiency in the simulation of large multi-agent systems which do not require high precision in very small clusters. To accomplish such a task in a single-threaded environment, one would normally employ standard techniques such as Octrees, or K-d trees which are relatively simple to implement in single-threading. K-d trees offer some parallelisation opportunities with some authors reporting piece-wise GPU datastructure construction algorithms which are of $\mathcal{O}(N \log^2 N)$ and also $\mathcal{O}(NlogN)$ complexity [27,28].



(a) A visual example of a K-d tree datastructure used extensively on CPU implementations for spatial partitioning purposes.



(b) K-d tree datastructure, with particles clustered so densely that the datastructure loses its effectiveness.

Figure 3: Examples of datastructures used in other implementations for the same purpose.

The former is normally the case when constructing a Kd tree by performing an O(NlogN) sort for every node in the tree to obtain the median. The latter is the case when using a linear median-finding algorithm, which is not trivial to implement on GPU, nor well-suited on GPU either.

One space partitioning algorithm responds very well to parallelism: grid-boxing [8] (sometimes known as uniform grid space partitioning). The NVIDIA CUDA SDK comes with a particle simulator which makes use of this technique [31]. Our single-GPU implementation of grid-boxing is heavily based on this.

Graphical Processing Units (GPUs) [32] lend themselves particularly well to accelerating simulations especially after the advent of NVidia's CUDA. CUDA is a powerful addition to C++ which allows the programmer to write code in C or C++ containing sections of special syntax which NVidia's special compiler in the CUDA SDK compiles into GPU instructions before handing the rest of the program code to the system's built in C or C++ compiler.

Figures 3(a) and 3(b) show examples of K-d trees,

which are very popular for accelerating ray tracing programs [27, 28]. Since many aspects of ray tracing are inherently parallel, K-d trees have also received a large amount of research effort in order to offload as much as possible computation from the CPU to GPU. For this reason, GPU K-d trees are of great interest, and numerous implementations already exist [33–35]. Each of these implementations have different ratios of CPU code to GPU code. Significant parallelism can be achieved by building the tree level-by-level instead of node-by-node, but even with this advantage, it is difficult to retain most computation on the GPU.

Our article is structured as follows. In Section 2 we present the method by which we construct the gridboxing datastructure, and also the method by which we use the result to accelerate body-body interactions with multiple-GPUs (mGPU) and also single-GPU in a standard simulation of the Boids model as propounded by Reynolds. Following this, in Section 3 we discuss the results we obtain and the approximate scaling we observe from this algorithm, from low numbers of agents (16384) to excessive numbers (1,000,000+). In Section 4 we discuss our results and implementation, and how it compares to other techniques. Section 5 provides conclusions and future work that we may pursue in this area.

2 Method

We use NVidia's CUDA as the parallelisation platform. With this choice comes several restrictions, and perhaps the most prominent of these comes with the memory architecture. A CUDA device's memory architecture is divided into several kinds: constant, global, shared, and texture. Each of these have different scopes and access penalties. Global memory takes by far the longest at about 200 cycles, followed by faster constant memory and very fast shared memory. The scope of these range from application to block level, and also kernel level. A CUDA block is conceptually a 1, 2 or 3D grid of CUDA threads, whose dimensions are arbitrarily defined by the user. For optimal results, it is beneficial to adjust these parameters to suit the program. This is especially important to allow latency hiding when reading global memory (the slowest).

Due to this architecture, it is impractical to use pointers and other simple data structures which require them on GPU. Efficient practical algorithms for generating trees on GPUs generally make use of hash tables, or spacefilling curves. Morton ordering (also known as Z-ordering or Nordering) is a very popular space-filling curve for the purpose of ensuring spatial locality [36]. Treecodes on GPU make use of Morton codes, which serve the purpose of encoding tree nodes for storing in hashtables [20, 34]. We use this method to increase coalesced global memory reads. The CUDA architecture generally reads fixed-size blocks of memory, and caches this for a short time when a thread requires it. It is for this reason that scattered reads from many threads earns a very large time penalty.

Algorithm 1 Single-GPU implementation of grid-boxing in boids.

| Allocate and initialise a <i>vec4</i> array as velocity |
|---|
| Allocate and initialise a vec4 array as position |
| Allocate a <i>uint</i> array as hashes |
| Allocate a <i>uint</i> array as indices |
| Allocate a <i>uint</i> array as gridBoxStart |
| Allocate a <i>uint</i> array as gridBoxEnd |
| Copy simulation parameters to device |
| copyVectorsToDevice() |
| //For <i>n</i> frames. |
| for $i \leftarrow 0$ to n OR NOT exit_condition do |
| $i \leftarrow i + 1$ |
| //Calculate hashes for all boids. |
| for $j \leftarrow 0$ to NUM_BOIDS do |
| hashes[j] = calculate_hash(position[j]) |
| end for |
| Sort by hash key (hashes, indices) |
| Populate gridBoxStart and gridBoxEnd |
| Scatter write boids |
| Perform Boid kernel |
| copyVectorsFromDevice() |
| drawBoids() |
| swapDeviceBuffers() |
| end for |
| |

Algorithm 1 presents the single-GPU version of gridboxing. All actions in this algorithm are performed in parallel, except for the loop containing the exit condition. This loop is simply used to advance to the next frame.

The very first kernel launched in this algorithm is the hash calculation kernel. This kernel has a simple task, and that is to populate the hashes and indices arrays. Algorithm 2 contains a pseudocode version of this kernel. Once this kernel has completed, the next step is to sort by the hash key. This step is accomplished using Thrust [37]. Thrust makes use of a parallel merge sort for our purposes. In appropriate situations, it uses an aggressively optimised GPU Radix sort written by

Algorithm 2 Data-parallel CUDA kernel which calculates hashes for all boids.

```
interview of an oblast
func calc_hash_d(...)
i = blockIdx.x * blockDim.x
+ threadIdx.x
if i < numBoids then
    p = positions[i]
//Morton ordering for hashes
    hash = calc_grid_hash(p)
    hashes[index] = hash
    indices[index] = index
end if
end
```

Merrill and Grimshaw [38].

Following the sorting phase of the algorithm, a single kernel is executed to populate the grid box starting and ending indices, as well as to perform a scatter write of the boids into their sorted positions, so that the hashes will be effective when the boid kernel is executed. The boid kernel itself simply evaluates the boids within the 8 grid boxes surrounding the grid box that the current boid is in. During this process, the sorting phase will afford improved coalesced memory reads for the device, and this greatly speeds up the process. The Morton ordering used when the hashes were calculated also increase memory locality of the grid boxes surrounding the current boid.

The kernel which calculates the grid box starting and ending indices does so by having each thread keep track of one grid box hash, and comparing itself against the next thread's hash. The threads which do not match the hash of the next thread are the ones which mark the boundary points of the grid box, and these threads simply write their starting and ending indices into an array in global memory. To accelerate this process, shared memory is used. The threads then reuse themselves by copying the boid indicated by their indices into the new sorted positions indicated by the index array reordered by the Thrust sorting phase.

We modified this algorithm to execute across several GPUs in order to both reduce computation time, and also to reach into even larger numbers of agents. In order to achieve this, we only parallelise one part of the algorithm across more than one GPU - the agentagent interaction kernel. In larger systems, the computation time required by this kernel dwarfs that of the datastructure construction time and all other subsequent steps taken. Therefore we find it most suitable to extend the algorithm in this manner first.

The first modification we made was the addition of host

POSIX threads. One thread is created for every GPU, and these control the memory and execution of each. A POSIX thread barrier is used to synchronise between the threads. Following this change, we modified the kernel execution code to perform interaction to parts of the total number of agents. For every POSIX thread, the range r of agents to evaluate is r = t/g where t is the total number of agents, and g is the total number of GPUs available. The starting index s is simply s = it/g. Once the agent interactions have been evaluated, the POSIX threads are synchronised and the complete parts are copied from each GPU back to the host. This process is effectively a gather operation. Once a frame is drawn or finished, each POSIX thread copies the entire host memory into its GPU, and the process repeats. It is worth noting that one could either have one GPU calculate the initial hashes, sort, and reordering, then copy (effectively scatter) the results to the other GPUs via host memory, or simply have each GPU independently perform these steps. It is a large waste of computation, but copying between GPUs will likely cause a higher overhead. In future, we hope to parallelise the previous steps also.

Following from Algorithm 1, the modified algorithm for use with more than one GPU is presented in Alg. 3. This algorithm is a very simple parallelism of Alg. 1, but it affords an excellent (albeit single-factor) increase in performance, and it also brings into reach even higher numbers of agents, with reasonable computing time.

Firstly, all vectors are copied onto all devices (a scatter operation). Then, for n frames, or until the application quits, each GPU will calculate hashes for each boid, then sort by hash, then populate the grid box starting and ending indices, scatter the boids into their sorted positions, and finally, perform a *partial* agentagent interaction kernel. Once this is complete, the results are gathered back onto the host, where it is either drawn on the screen, and/or scattered back to the GPUs and the process repeats.

The simulation model that we use for benchmarking these techniques is the Boids model, originally by Craig Reynolds [6, 7]. In our simulation, we make use of a goal rule to direct boids toward the center of a fixed size box. Figures 4(a) and 4(b) show a visual indication of boid clustering near the start and end of a simulation, respectively. This shows more clearly that during simulation, boids will cluster more tightly in the middle of the box. We ensure that this is the case so that we can determine exactly how the effectiveness of grid-boxing is lost as boid clustering increases, and gridboxes contain more and more boids.

Algorithm 3 Multiple-GPU (mGPU) implementation of grid-boxing in boids.

Allocate arrays

Copy simulation parameters to constant memory on all devices HOST.copyVectorsToDevices() for $i \leftarrow 0$ to n OR NOT exit_condition do $i \leftarrow i + 1$ for each gpu g in parallel do for $j \leftarrow 0$ to NUM_BOIDS in parallel do hashes[j] = calculate_hash(position[j]) end for Sort by hash key (hashes, indices) Populate gridBoxStart and gridBoxEnd Scatter write boids Perform Boid kernel end for copyVectorsPartiallyFromDevices() drawBoids() copyVectorsToDevices()

end for



(a) A histogram of the spatial lo- (b) Another histogram of spatial cations of boids near the begin- locations of boids several hunning of the simulation. dred frames later.

Figure 4: Histogram plots of spatial locations of boids during the simulation. Red refers to a cell with 1 boid, green 2 boids, blue 3 boids, and yellow 4 boids.

3 Results

We tested the performance of our mGPU and single-GPU versions of the grid-boxing algorithm. The results for mGPU interaction and datastructure construction times are presented in Figures 5 and 6 in the form of a y-normalised plot for interaction computation, and a linear plot for the datastructure construction time. In both of these plots, the data points taken are 10-frame averages. For the mGPU implementation, we average the compute time across the four GPUs that we use. The testing GPUs we used are two NVidia GTX 590s. Each of these cards have two physical GPUs each.

Our simulation is simply a flock of boids moving to-



Figure 5: A y-normalised plot of the agent-agent interaction CUDA kernel.



Figure 6: A linear plot of the construction time needed for the datastructure.

ward the centre of the box they are in. Movement towards the centre causes a clustering effect in the centre of the box, which compresses and pulsates a number of times before reaching what can be described as equilibrium. This accounts for the initial increase in computation time, followed by slight harmonic fluctuations.

Figure 7 shows how the mGPU implementation scales against the single-GPU version of the algorithm. It shows an impressive performance gain over the single-GPU version, particularly as the system becomes larger. Smaller systems are simply more suited to being evaluated on one GPU to avoid the large overhead of memory copies across the PCI-E bus.



Figure 7: A linear plot of the mGPU and single-GPU interaction compute time vs the system size.

4 Discussion

Uniform grid-boxing might well be a more simplistic space partitioning method, but it is readily parallelisable, as opposed to other datastructures inspired from single-threaded solutions. These solutions often rely on pointers for the simplest implementation, as well as recursion. These two techniques are the bane of dataparallel programs.

The CUDA-GPU memory architecture does not allow for a heap, and branch diversion is an even larger problem. Branching is allowed in CUDA, but it comes at a performance loss. Each time threads in a warp diverge from one another, some threads simply execute NOOPs until they converge again. Thus, a CUDA kernel is most effective when its threads execute the same instructions.

Other space partitioning algorithms also exist which perform well, but they only slightly resemble the algorithms which inspire them. For example, it may well be possible to replace a K-d tree with a hashtable to make it suitable for execution on GPU, but this multi-threaded hashtable requires mutexes - another branch diversion issue. It seems the only way to perform space partitioning in reasonable time on GPU to still make the effort of writing a CUDA-based program is to utilise some kind of sorting algorithm to place particles in very deliberate places in memory, to gain both spatial locality and also some structured method of traversing this data.

Sorting on GPU has been studied extensively, and many excellent solutions already exist, most notably the Radix sort by Merrill and Grimshaw merged into the Thrust library, which is shipped with CUDA [38]. The uniform nature of this grid-boxing datastructure makes mGPU an attainable goal for multi-agent simulations, as it allows for distributing data with minimal copying between GPUs, while still preserving accuracy. This technique has seen a lot of use already.

Finally, it is worth noting also that grid-boxing is only well-suited to algorithms which have a fixed interaction distance. The smaller this distance is, the better the algorithm will perform, given that the grid is appropriately sized. For the results to remain fully accurate and representative of the original simulation code, the interaction distance must be smaller than the grid box size.

5 Conclusions

We have presented a multiple-GPU (mGPU) and single-GPU solution for grid-boxing in multi-spatial-agent simulations. Performance measurements were made, and the two algorithms compared, and an overwhelming clear advantage of mGPU was seen over the single-GPU implementation as expected. Even though memory sacrifices were made in the process, larger systems continued to receive an excellent increase in performance. These techniques are vital for facilitation of large system sizes of agents and the investigation of logarithmic scaled emergent phenomena in simulations.

In future we hope to parallelise the hash calculation and the sorting and reordering phases also. These phases are more difficult to parallelise over multiple GPUs, but this will greatly improve performance, and allow the ability to reach into much larger systems as well. These approaches will allow interactive-time visualisation of large scale spatial agent simulations and hence make it easier to search the parameter space of the models.

References

- Couchman, H.M.P.: Cosmological simulations using adaptive particle-mesh methods. Web (1994)
- [2] Couchman, H.M.P.: Gravitational n-body simulation of large-scale cosmic structure. Number 13. In: The Restless Universe Applications of Gravitational N-Body Dynamics to Planetary Stellar and Galactic Systems. Taylor and Francis (2001) 239–264 ISBN: 978-0-7503-0822-9.
- [3] Couchman, H.M.P.: Simulating the formation of largescale cosmic structure with particle-grid methods. J. Comp. and Applied Mathematics 109 (1999) 373–406
- [4] Katzenelson, J.: Computational structure of the n-body problem. Technical Report AI Memo 1042, MIT AI Lab (1988)

- [5] Liewer, P.C., Decyk, V., Dawson, J., Fox, G.C.: A universal concurrent algorithm for plasma particle-in-cell simulation codes. In: Proc. Third Hypercube Conference. Number C3P-362 (1988) 1101–1107
- [6] Reynolds, C.: Flocks, herds and schools: A distributed behavioral model. In: SIGRAPH '87: Proc. 14th Annual Conf. on Computer Graphics and Interactive Techniques, ACM (1987) 25–34 ISBN 0-89791-227-6.
- [7] Reynolds, C.: Boids background and update. http: //www.red3d.com/cwr/boids/ (2011)
- [8] Hawick, K.A., James, H.A., Scogings, C.J.: Gridboxing for spatial simulation performance optimisation. In T.Znati, ed.: Proc. 39th Annual Simulation Symposium, Huntsville, Alabama, USA (2006) The Society for Modeling and Simulation International, Pub. IEEE Computer Society.
- [9] Guo, D.: Spatial cluster ordering and encoding for highdimensional geographic knowledge discovery. Technical report, GeoVISTA Center and Department of Geography, Pennsylvania State University (2002)
- [10] Guo, D., Gehagen, M.: Spatial ordering and encoding for geographic data mining and visualization. J. Intell. Inf. Sys. 27 (2006) 243–266
- [11] Appel, A.W.: An efficient program for many-body simulation. SIAM J. Sci. Stat. Comput. 6 (1985) 85–103
- [12] Singer, J.K.: Parallel implementation of the fast multipole method with periodic boundary conditions. East-West J. Numer. Maths 3 (1995) 199–216
- [13] Gelato, S., Chernoff, D.F., Wasserman, I.: An adaptive hierarchical particle-mesh code with isolated boundary conditions. The Astrophysical Journal 480 (1997) 115– 131
- [14] Barnes, J., Hut, P.: A hierarchical o(n log n) forcecalculation algorithm. Nature 324 (1986) 446–449
- [15] Burton, A., Field, A.J., To, H.W.: A cell-cell barnes-hut algorithm for fast particle simulation. Australian Computer Science Communications 20 (1998) 267–278
- [16] Darve, E., Cecka, C., Takahashi, T.: The fast multipole method on parallel clusters, multicore processors, and graphics processing units. Compte Rendus Mecanique 339 (2011) 185–193
- [17] Warren, M.S., Salmon, J.K.: Astrophysical n-body simulations using hierarchical tree data structures. In: Proc. ACM IEEE Conf on Supercomputing (SC92). (1992)
- [18] Samet, H.: Data structures for quadtree approximation and compression. Communications of the ACM 28 (1985) 973–993
- [19] O'Connell, J., Sullivan, F., Libes, D., Orlandini, E., Tesi, M.C., Stella, A.L., Einstein, T.L.: Self-avoiding random surfaces: Monte carlo study using oct-tree datastructure. J. Phys. A: Math. and Gen. 24 (1991) 4619– 4631
- [20] Warren, M.S., Salmon, J.K.: A parallel hashed oct-tree n-body algorithm. In: Supercomputing. (1993) 12–21
- [21] Libes, D.: Modeling dynamic surfaces with octrees. Computers & Graphics 15 (1991) 383–387
- [22] Bedorf, J., Gaburov, E., Zwart, S.P.: A sparse octree

gravitational n-body code that runs entirely on the gpu processor. J. Comp. Phys. **Online December** (2011) 1–34

- [23] Dubinski, J.: A parallel tree code. Interface **1** (1996) 1–19
- [24] Gope, D., Jandhyala, V.: Oct-tree-based multilevel lowrank decomposition algorithm for rapid 3-d parasitic extraction. IEEE Trans. Computer-Aided Design of Integrated Circuits and Systems 23 (2004) 1575–1580
- [25] Jackins, C.L., Tanimoto, S.L.: Oct-trees and their use in representing three-dimensional objects. Computer Graphics and Image Processing 14 (1980) 249–270
- [26] Ricker, P.M.: A direct multigrid poisson solver for octtree adaptive meshes. The Astrophysical Journal: Supplement Series 176 (2008) 293–300
- [27] Danilewski, P., Popov, S., Slusallek, P.: Binned sah kdtree construction on a gpu. Technical report, Saarland University, Germany (2010) 15 Pages.
- [28] Yang, H.R., Kang, K.K., Kim, D.: Acceleration of massive particle data visualization based on gpu. In: Virtual and Mixed Reality, Part II HCII 2011. Number 6774 in LNCS, Springer (2011) 425–431
- [29] Zhang, Y., Liu, L., Lu, X., Li, D.: Efficient range query processing over dhts based on the balanced kautz tree. Concurrency and Computation: Practice and Experience 23 (2011) 796–814
- [30] Kravtsov, A.V., Klypin, A.A., Khokhlov, A.M.: Adaptive refinement tree: A new high-resolution n-body code for cosmological simulations. The Astrophysical Journal - Supplement Series 111 (1997) 73–94
- [31] Green, S.: Particle simulation using cuda. NVIDIA (2010)
- [32] Leist, A., Playne, D., Hawick, K.: Exploiting Graphical Processing Units for Data-Parallel Scientific Applications. Concurrency and Computation: Practice and Experience 21 (2009) 2400–2437 CSTN-065.
- [33] Akarsu, E., Dincer, K., Haupt, T., Fox, G.C.: Particlein-cell simulation codes in high performance fortran. In: Proc. ACM/IEEE Conf. on Supercomputing (SC'96), Washington DC, USA (1996)
- [34] Nakasato, N.: Implementation of a parallel tree method on a gpu. J. Comp. Science Online March (2011) 1–10
- [35] Schive, H.Y., Zhang, U.H., Chiueh, T.: Directionally unsplit hydrodynamic schemes with hybrid mpi/openmp/gpu parallelization in amr. Int. J. High Perf. Computing Applications Online November (2011) 1–16
- [36] et al., M.: Gpu octrees and optimized search. VIII Brazillian Symposium on Games and Digital Entertainment (2009) 73–76
- [37] Hoberok, Bell: Thrust: A parallel template library. http://www.meganewtons.com/ (2011)
- [38] Merrill, Grimshaw: Revisiting sorting for gpgpu stream architectures. In: PACT '10 Proceedings of the 19th international conference on Parallel architectures and compilation techniques. (2010)

21

3DNET - An ecosystem for the development, evaluation, and sharing of visualization workflows

S. Grimm¹, A. Paluszny², L. Parsonson³, R.Andrian⁴, W. Hernandez⁵,

L. Bourn¹, A. Bajwa¹, H. Hatzakis¹, L. Bai³

¹Biotronics3D, London, UK ²Imperial College London, UK ³University of Nottingham, UK ⁴The Institute of Cancer Research ⁵Universidad Central De Venezeula

Abstract—We present 3DNet: an integrated ecosystem that serves as a collaborative platform to combine efforts towards the development, evaluation, and sharing of visualization workflows. This system represents a paradigm shift by providing not only a web-based visualization toolkit but an instrument for standardized collaborative research that allows independent scientific entities to share the development effort and better combine their expertise. Two case studies illustrate key system design decisions as well as its usability and effectiveness: one example from medical imaging, and the other from reservoir engineering.

Keywords: Visualization, workflow, cloud, collaboration, ecosystem.

1. Introduction

We define visualization workflow as a configurable, structured set of semi- or fully-automated steps which implement a computational solution to a scientific problem, particularly those which involve the processing of image data by means of visualization, rendering, and analytics. At any given time, multiple scientists collaborate in the development and evaluation of these workflows: from the acquisition at a lab or a hospital, the automatic filtering, automatic or manual segmentation, to the computation of metrics and analytics. Participating entities are rarely in the same geographical location. When distributed research methodologies are utilized, management of resources raises issues, such as version control, data integration and security. Further challenges related to the development of complex visualization workflows include creating, editing and finetuning proposed approaches in an efficient manner. The execution of these workflows often generates vast amounts of data, and it becomes important to be able to track its provenance and how it was processed over the course of a research effort [4]. Thus, an integral solution for interdisciplinary collaboration, providing centralized storage, access to reusable visualization components, and offering coordination and communication capabilities for distributed development is essential.

Existing technologies often provide only a partial solution. For example, designing a visualization workflow may include the implementation of well-established standard algorithms. An attempt to curb this problematic is to employ third-party toolkits, such as VTK, ITK, ProtoVis, Pegasus, Triana, and Seg3D to reduce development time. WebGLot [14] goes a step further by allowing high-performance visualization in the browser. Another example is sense.us [8], a web-based tool which supports asynchronous collaboration for visualization, but does not address issues such as data storage, security and sharing. Other efforts addressing the problem of scientific workflows are the Kepler system [15] and Taverna: a tool for building and running workflows of services [9]. These technologies present partial remedies, but do not address the problem at hand from an integral perspective.

In this paper, we present 3DNet, a purely cloud-based answer to the problem of data sharing and collaboration in visualization. 3DNet provides an all-around solution to the problem of access and exploitation of existing visualization technology while offering a platform for swift visualization algorithm and workflow development. It allows researchers to make use of a feedback style of development whereby workflows can be modified following the analysis of results, including editing analysis parameters, and plugging in new analysis components with relative ease.

The rest of the paper is organised as follows. Section 3 discusses the used technologies and the motivation behind these choices, as well as giving an overview of the lessons we learned and the evaluation of the effectiveness of the ecosystem. Finally, Section 4 presents two specific case studies which illustrate the applicability of the ecosystem to solving problems in the areas of cancer research and reservoir engineering.

2. The 3DNet Ecosystem

The 3DNet ecosystem is an environment consisting of all components required to develop a visualization workflow, including collaborators and the technology with which they interact. Figure 1 shows a high-level overview of the 3DNet ecosystem and its components. The system has a visualization toolkit, a configurable workflow, plugins, a ubiquitous access system and a data management kernel. The ecosystem allows for an iterative, incremental methodology for the agile design and development of visualization workflows. Entities within the ecosystem interact with the workflow at any level, feeding it with data, running analysis modules, modifying experimental parameters and sharing results.



Fig. 1: The ecosystem allows for collaboration between different experts that can upload data onto the cloud and simultaneously create and edit the visualization workflow and analytics module. Output data, images and screen-shots are output and centrally stored, and are available to all participating experts.

2.1 Operations and Processes

Operations are actions that can be undertaken within the ecosystem by an entity. The ecosystem supports the following operations: (a) acquisition and import of data, (b) running of analysis modules, (c) storage and analysis of results, and (d) creation and modification of workflow, among others. Multiple operations can be grouped together, and this is referred to as a process. Processes allow entities to run multiple operations with minimal effort. This includes running operations concurrently, for instance to perform different analysis tasks on the same data set, and running operations in series, for instance to apply a set of filters sequentially. Processes form the core of the workflow, although a workflow need not be restricted to a single process. Multiple processes allow researchers to test and choose between different analytical techniques with ease, as well as experimenting with new combinations of operations on their data sets.

2.2 Life-Cycle

The ecosystem paradigm encourages an evolutionary lifecycle, driven by feedback and continuous development. The ecosystem allows researchers to accurately track the influx of data and the results of analyses performed, as well as maintaining a record of workflows created. A typical interaction is as follows: (1) Data is acquired and uploaded onto the ecosystem. (2) This data is accessed by analysts who create a workflow process with a particular goal in mind, using a basic combination of the appropriate analysis modules. (3) The process is run on the data, and the results analyzed for effectiveness. (4) At this point, the analyst can choose to use these results to modify the workflow process in some way, perhaps choosing a different combination of analysis modules, or modifying some input parameters, and then running the process again. (5) This process of evaluation and optimization continues indefinitely, or until a satisfactory result is achieved. (6) Administrators can monitor activities and progress, access reports and for more standardized protocols, can set permission rights to some of the entities to alter the workflows. Within this life-cycle it is possible, and often desirable, to include other users as reviewing parties, who can verify results of a particular analysis. Webbased development allows multiple parties to participate in the process of workflow development by manipulating single visualization components via a centralized, web-based solution. A web-portal is central to the solution, in this case, the 3DNet web-portal provides a centralized mechanism for login and access to the system.

3. 3DNet Technologies

The ecosystem exists within a cloud environment and can be accessed via web-browsers with a zero-weight client. Thus, the access is operating system-independent and installation free and collaboration requires no hardware or software setup. The entire system load is borne by clusters of servers, where each cluster is responsible for a different activity. Servers run out of a data center remote to the collaborators. Figure 2 shows the key components of the system:

- The Data Foundation Platform (DFP), which administrates the complete system in terms of information and database management, services, and security.
- The Visualization Workflow Platform (VWP), which handles all presentation, analysis and management of the information.
- The Development Platform Environment (DPE), which provides the tools for external innovators to implement and integrate modules into the system.

3.1 Data Foundation Platform

The data foundation platform handles all matters related to the administration and management of data at the system level. A data gateway provides infrastructure for the loading of data into the system which is then handled by the data import manager which incorporates acquisitions into the database. The latter interacts with the database manager which can classify the datasets based on their metadata, e.g. file headers in DICOM. DICOM data from a variety of application areas can be imported, e.g. magnetic resonance, computed tomography, digital microscopy, mammography, digital radiography, and ultrasound. Data can be exported in a transformed manner to DICOM and subsequently stored, downloaded, or shared. The database manager also implements multi-tier storage providing dataset caching as well as short- and long-term storage. The security and audit manager is in charge of controlling access to the database and tracking usage. 3DNet implements an on-demand public-key infrastructure based on strong encryption which is secure and transparent for the end-user via https and at the same time flexible to manage for IT engineers. The software core components highly optimize data flows to maintain a high security levels without compromising image quality.

3.2 Visualization Workflow Platform

The visualization workflow typically involves three phases: data exploration, analysis, and validation. Each phase has unique and specific requirements in terms of use, interaction and presentation of the data.

Data exploration involves understanding the data and identifying problems. It relies on the application of algorithms to manipulate the input data and display it in a manner which is convenient for interpretation. The visualization engine (VE) provides methods to manage the complexity inherent to high dimensionality. The VE includes provision of algorithms such as volume rendering, maximum intensity projection, thick slabs, and alike to the workflow. Furthermore, large amounts of data also require information to be presented relative to its importance and relevance. These are included in the VE in the form of modules for focus and context [6] and importance-driven rendering [25]. The VE offers modules for fusion and synchronization that can be incorporated to the workflow for datasets that require combination or comparison. Finally, layouts are defined to specify how the output of the VE modules is viewed in terms of size and location within the screen.

Data analysis relies on the extraction or derivation of information from the original data. The Analytics Engine (AE) comprises modules that manipulate data to extract information, including the application of analytical formulations and the extraction of models for simulation, e.g. [7]. AE includes modules such as segmentation, data filtering, contour extraction, integration with computational models, and alike. The analysis process often involves various steps which combine multiple of these methods in a specific order. Each module consists of a kernel, a set of rules, and a set of configurations. The kernel is applied to the data to execute Finally, the validation phase involves the aggregation and integration of data that emerges during the execution of the visualization workflow. The Management Engine handles such tasks by providing modules for semantics-based search, data mining, fetching related datasets as well as statistics, auditing, and alike. Furthermore, the ME is in charge of controlling the load balancing and the virtualization of the external modules, explained in Section 3.3.

3.3 Development Platform Environment

To provide the most flexible way of developing for 3DNet, two methods of integrating code were developed. The first, the development kit, provides a broad set of common tools for both analysis and visualization modules. The second, the module wrapper, allows developers who have already developed their module to integrate with the minimum of effort. The development kit offers developers the opportunity to create a module directly utilizing 3DNet's fundamental data types, removing the need for conversion between data formats. The development kit also provides a library of common components and functions to increase the speed of module development. The module wrapper was created with flexibility in mind. The idea was to create a method of seamlessly integrating externally developed code with the minimum of ease. This required building a scalable component which could include, if required, scalable instance management, and DICOM compliant networking capabilities. This would allow developers to wrap code that had not necessarily been developed with 3DNet in mind, and include developers who already had coding expertise in another technology that they wanted to use in the ecosystem.

The Visualization and Analytics Engines were developed using C++ in order to take advantage of the increased level of memory access offered and faster computation achievable. The entire core makes use of multi-threading and processing to support multiple users. In principle, parts can also be mapped to the GPU, but this has not been a primary goal due to the increased hardware cost. All functionality is wrapped in .NET providing access through the .NET CLR (Managed C++, C#, Visual Basic, among others). The toolkit handles all encrypted message processing with the various web services, and server side modules. Thus, the user codes directly against local classes and resources.

Server-side modules can be developed in any language, e.g. C++, Java, and C#, and can even be developed using higher-level toolkits such as Matlab, IDL, ITK, and VTK. Thus, existing modules can easily be ported to the ecosystem and be added to the designed workflow. Client-side



Fig. 2: The ecosystem allows for collaboration between different experts that can upload data onto the cloud and simultaneously create and edit the visualization workflow and analytics module. Output data, images and screen-shots are output and centrally stored, and are available to all participating experts.

modules can be developed in Silverlight following given design templates to maintain a common look-and-feel of the system. Silverlight is a browser- and platform-independent implementation of the .NET framework. Thus, it offers many of the advanced features of C# such as generics, anonymous delegates, and LINQ, as well as having the support of Microsoft Visual Studio for coding and Expression Blend for developing visual layouts. However, it is a new technology and therefore some advanced functionality is currently missing from the Silverlight API, making development difficult at times, exacerbated by poor documentation, and reliance on community samples. Nevertheless, we chose this as the underlying client technology as Silverlight shows the most promising technology for RIA application on a long term basis.

4. Case Studies

This section describes two study cases used to showcase the capabilities of the 3DNet ecosystem. The first is an example from cancer research, specifically the development of visualization workflows for clinical trials which involve automatic processing of large quantities of medical image data. The second describes the application of the ecosystem to reservoir engineering research: in particular the semiautomatic processing of high resolution CT data.

4.1 Cancer Research

Clinical trials are a fundamental step in research for cancer treatment. In drug evaluation trials, the studies are conducted and managed by clinical trial units (CTU) in the academic sector, or contract research organizations (CROs), and require a tight collaboration between scientists, clinicians and pharmaceutical companies. Nowadays, imaging is becoming a central component in the assessment of a drug treatment response [2]. Technological advances provide better, richer anatomical and functional images through a growing number of modalities and protocols, e.g. [13], [22]. By combining clinical inputs, electronic data and numerical analysis approaches, researchers can test hypothesis, derive quantitative measures, and correlate changes with clinically meaningful biological phenomena. Validation of both measures and analysis process is critical to rate the effectiveness of a drug, guarantee safety for the patients, and permit route to commercialization. However, the explosion of information raises challenges that require development of computational systems to support the exchange of vast amounts of imaging data between the various facilities engaged in the trial, and the management of complex workflows involving the multiple actors.

4.1.1 Collaboration Challenges

Clinical trials involve a large diversity of roles that need to collaborate and interact at all stages of a study, which in many instances include the following steps:

Study protocol: Pharmaceutical companies define the trial and imaging protocols.

Study population: Entrusted clinicians identify a study population for the type of drug to trial and cancer to investigate.

Data acquisition: An IT infrastructure needs to be set to collect and manage all data for the study.

Data validation and preparation: All acquired data go through a pre-assessment step to guarantee that they match the defined protocols, which is particularly crucial when the study involves multiple centers. Following the validation, pre-processing steps are sometimes required to prepare the data for analysis.

Data Analysis: Scientists use a combination of in-house and commercial analysis packages to extract quantitative measures from the data. The formulation of the hypothesis, implementation of the methods, testing and validation against clinical evidences relies on a back-and-forth communication between scientists and clinicians. The efficacy of the analysis heavily depends on parameters, intermediate data and results to be easily accessible to both ends at all time, to allow tuning, comparisons and optimizations of the methods.

Data Integration: CTUs and/or instigators of the trial aggregate all outputs from the complete study and perform

more statistical analysis to validate the whole methodology and procedure.

Unfortunately, the lack of systems that provide easy accessibility to information and knowledge sharing, which are key to create a truly interactive and efficient collaboration between all participants, slow down clinical trial processes, with consequences not only in hampering the path to health improvement, but also in financial terms for the instigators.



Fig. 3: Clinical trial workflow overview, key actors, and generated metadata - electronic records.

4.1.2 Management Challenges

The significant amount of data acquired and generated along the process and its management represent yet a crucial endeavour for CTUs.

A foremost aspect in managing clinical trial data is traceability through each step of the process. When multiple centers are involved in a clinical trial, or when part of a trial is outsourced, datasets are strongly recommended to be 'anonymized', with the instigator being able to associate them back to the original datasets. Valid and invalid datasets need to be identified and labeled. Pre-processing and analysis methods, as well as the parameters used, need to be logged so that the experiments can be reproduced. Although results are used in global statistical analysis to assess the overall effect of a drug, each single output needs to be traced back to the data from which it derived.

The records of all interactions and alterations made to the data, as well as the detailed reports of each aspect of data processing are electronic metadata that need to be managed alongside the images. These images do not only come from different modalities or different acquisition protocols to provide different types of information, but they can also encompass data acquired pre- and post-treatment that need to be combined to effectively assess treatment response. The various derived outputs are not only numerical values such as summary statistics, but also include regions of interest and multi-dimensional biomarker maps. Presenting and visualizing the profusion of information in an intuitive way that facilitates clinical interpretation is paramount to efficient and successful trial process.

The ability to represent the complete clinical trial workflow in a way that allows the various actors to check progress, and assess or re-evaluate any intermediate state, alternative method or optimization mechanism, is also fundamental to administrate, audit, and in turn validate the clinical trial process. But the need and benefit is also financial, as real-time process monitoring, recording and eventually reengineering is essential to improve efficiency and reduce redundancies and errors [12].

4.1.3 3DNet in Cancer Research

Adoption of the 3DNet ecosystem can provide valuable benefits in the process of validating pharmaco-kynetic (PK) models for dynamic contrast enhanced magnetic resonance images (DCE-MRI) [24]. 3DNet allows researchers and radiologists to access and share data from both clinic and research sites throughout the protocol experiments. Within the Development Platform Environment, in-house generic functional PK models written in IDL can be integrated into the system and tested against the trial data. Results (and parameters states) are stored as new DICOM series of quantitative maps that the clinicians can review as overlays on top of the original images, and comment, at any time. The validation procedure consists of iterative refining and optimization steps between researchers and clinicians, but also between the different researchers implementing various parts of the PK analysis, until the generic model, and its parameters, get tuned for a specific anatomical application (f.i. breast). During the complete process, data, results, reports and source code are always accessed from the same, centralized, single point of access: the cloud, which also represent a major advantage for the administrators in terms of monitoring, management and data aggregation.



Fig. 4: Implementing and adapting generic pharmaco-kynetic model to specific tumors types in drug treatment response quantification for (a) brain cancer, (b) liver metastasis and (c) breast cancer.

CO₂ can be injected into a rock at high pressure and great depth with the intent of trapping the nocive gas for thousands of years [16]. This concept, also known as CO_2 sequestration, is a very relevant research topic in reservoir engineering aimed at reducing CO₂ emissions to meet greenhouse emission targets and curb the imminent temperature increase due to concentration of these gases in the atmosphere. To this end, large efforts are invested in acquiring high resolution images of sandstones, carbonates and other rocks of current oil reservoirs, which make excellent candidates for systematic CO_2 sequestration as they are naturally contained geological systems at great depth and with the infrastructure required to inject fluids into the earth's subsurface. Samples are studied at the microscale and contain water, oil, and CO₂ at ambient and high temperature conditions. The idea is to understand how the geometry of the pore space of these rocks affects their macroscopic fluid properties and their ability to effectively trap CO_2 on a geological timescale [11].

4.2.1 Pore-Scale Visualization Workflow

In this study, the acquisition consists of micro-computed tomography (μ -CT) images with monochromatic x-ray (30 keV) datasets of sandstones with oil and brine generated at the synchrotron radiation beamline. The usable portion of the image is approximately 300x300x300 voxels and the nominal resolution is 9 microns (see Figure 5a). The visualization workflow has the following steps: crop to usable area (removing boundary effect areas); remove ring artifacts, including: identify ring center, apply transformation: to Polar coordinates, apply stripe removal: with combined wavelet - Fourier filtering [17], apply transformation: to Cartesian coordinates; adjust cropping; apply 3x3x3 Median/Otsu filter to remove noise [18]; (auto) adjust image color level; reduce color numbers / posterize; and segment CO2, brine and rock using thresholding. Figure 5 shows an example of an (a) initial, (b) filtered, and (c) segmented image. Post-processing of the segmented data involves: computing contact angles between oil bubbles, brine and rock; computing cluster size distribution; characterization and quantification of 3D oil clusters, extraction of pore network, extract pore size distribution, extract CAD geometry of pore space, running numerical simulations on extracted structures, and visualization of the simulation results. The visualization workflow often entails stitching multiple scripts together that trigger applications such as ImageJ and Gimp, with Javascript and C++ code which generate the images that can be overlaid with data using VTK. Analytics are usually visualized as graphs using tools such as MS Excel. All data transfers are made via email, or file exchange services, as well as all output and generated analysis.



Fig. 5: CT Image from a Sandstone Berea imbibed with water: (a) shows the original image - the bright spots are heavy metals present in the sample, (b) is the image after color adjustment, (c) illustrates thresholding on the filtered image.

4.2.2 Visualizing Shared Data

Visualization plays a key role in understanding the underlying mechanisms that control the feasibility of carbon capture and sequestration. However, it is not one of the core strengths of earth scientists, therefore, collaborations are usually required to process and analyze data acquired in the laboratory [21]. Sharing data entails confidentiality issues, as the data may belong to a different institution than the one that the scientist is affiliated to. For example, data may belong to another university, it may belong to an oil company or to the government. Thus, sending and transferring it might result conflictive and even illegal actions. Nonetheless, making data available for external analysis in a seamless and secure environment reduces the bureaucratic burden of the confidentiality and disclosure problematic. Furthermore, the transfer of large quantities of data becomes cumbersome, and introduces a further level of complexity, distracting the scientist of his true path. The problem of analyzing the acquired data further broadens to housekeeping of the files, tracking versions, controlling modifications, and storing large quantities of data. Moreover, as datasets become larger and numerous, the scientist is confronted with the problem of storing them, keeping backups, keeping them safe, and mobilizing them if it becomes necessary. The collaboration entails jointly building a visualization workflow that can be applied to the acquired data. The 3DNet ecosystem allows for the application of multiple analysis techniques in a single workflow. Users are able to assume multiple roles within the ecosystem, enabling them to create and edit visualization workflows, upload data, and analyze results as required.

5. Conclusion

We have presented an ecosystem that provides an infrastructure for the design, development, and evaluation of visualization workflows within the context of a collaborative environment. It is enabled by web technology and serves as a mechanism to connect remote parties working on common visualization workflows. It provides a secure and stable platform for data management and is specialized in the



Fig. 6: 3DNet in reservoir engineering

processing of images. Such a sytem is instrumental in the design and development of visualization workflows as it allows all involved specialists to collaborate in an efficient manner in generating prototypes that reflect novel research. As opposed to the development of web-tools, which brings incremental change into the way that visualization integrates into research across multiple disciplines, an ecosystem that embeds visualization into the cloud and provides an integral solution to researchers represents a step change. A true new methodology, available to the research community to interact, create and fine-tune visualization workflows, that answer challenges of inter-disciplinary collaboration and provides a competitive advantage to teams that choose to participate in it.

References

- [1] Seg3D: Volumetric Image Segmentation and Visualization. Scientific Computing and Imaging Institute (SCI).
- [2] H. Bliddal, M. Boesen, R. Christensen, O. Kubassova, and S. Torp-Pedersen. Imaging as a follow-up tool in clinical trials and clinical practice. *Best Practice and Research Clinical Rheumatology*, 22(6):1109 – 1126, 2008. Imaging and Musculoskeletal Conditions.
- [3] M. Bostock and J. Heer. Protovis: A graphical toolkit for visualization. *IEEE Transactions on Visualization and Computer Graphics*, 15:1121–1128, 2009.
- [4] S. B. Davidson and J. Freire. Provenance and scientific workflows: challenges and opportunities. In *Proceedings of the 2008 ACM SIGMOD international conference on Management of data*, SIGMOD '08, pages 1345–1350, New York, NY, USA, 2008. ACM.
- [5] E. Deelman, G. Singh, M. hui Su, J. Blythe, A. Gil, C. Kesselman, G. Mehta, K. Vahi, G. B. Berriman, J. Good, A. Laity, J. C. Jacob, and D. S. Katz. Pegasus: a framework for mapping complex scientific workflows onto distributed systems. *Scientific Programming Journal*, 13:219–237, 2005.

- [6] H. Doleisch, M. Gasser, and H. Hauser. Interactive feature specification for focus+context visualization of complex simulation data. In *Proceedings of the symposium on Data visualisation 2003*, VISSYM '03, pages 239–248, Aire-la-Ville, Switzerland, Switzerland, 2003. Eurographics Association.
- [7] K. Dykstra, R. Pugh, and A. Krause. Visualization concepts to enhance quantitative decision making in drug development. *The Journal of Clinical Pharmacology*, 50(9 suppl):130–139, 2010.
- [8] J. Heer, F. B. Viégas, and M. Wattenberg. Voyagers and voyeurs: Supporting asynchronous collaborative visualization. *Commun. ACM*, 52:87–97, January 2009.
- [9] D. Hull, K. Wolstencroft, R. Stevens, C. Goble, M. Pocock, P. Li, and T. Oinn. Taverna: a tool for building and running workflows of services. *Nucleic Acids Research*, 34(Web Server issue):729–732, July 2006.
- [10] L. Ibanez and W. Schroeder. *The ITK Software Guide 2.4*. Kitware, Inc., Nov. 2005.
- [11] S. Iglauer, S. Favretto, G. Spinelli, G. Schena, and M. J. Blunt. X-ray tomography measurements of power-law cluster size distributions for the nonwetting phase in sandstones. *Phys. Rev. E*, 82(5), 2010.
- [12] S. A. Khan, R. Kukafka, J. T. Bigger, and S. B. Johnson. Reengineering opportunities in clinical research using workflow analysis in community practice settings. AMIA - Annual Symposium proceedings / AMIA Symposium. AMIA Symposium, pages 363–367, 2008.
- [13] D. Klimov, Y. Shahar, and M. Taieb-Maimon. Intelligent visualization and exploration of time-oriented data of multiple patients. *Artificial Intelligence in Medicine*, 49(1):11 – 31, 2010.
- [14] D. Lecocq, M. Hadwiger, and A. Rockwood. Webglot: highperformance visualization in the browser. In SIGGRAPH '10: ACM SIGGRAPH 2010 Talks, pages 1–1, New York, NY, USA, 2010. ACM.
- [15] B. Ludäscher, I. Altintas, C. Berkley, D. Higgins, E. Jaeger, M. Jones, E. A. Lee, J. Tao, and Y. Zhao. Scientific workflow management and the kepler system. *Concurrency and Computation: Practice and Experience*, 18(10):1039–1065, 2006.
- [16] B. Metz, O. Davidson, H. de Coninck, M. Loos, and L. Meyer. Special Report on Carbon Dioxide Capture and Storage. Intergovernmental Panel on Climate Change, 2005.
- [17] B. Munch, P. Trtik, F. Marone, and M. Stampanoni. Stripe and ring artifact removal with combined waveletfourier filtering. *Optics Express*, 17:8567–8591, 2009.
- [18] N. Otsu. A thresholding selection method from gray-level histogram. IEEE Transactions on Systems, Man and Cybernetics, 9:62–66, 1979.
- [19] L. Parsonson, L. Bai, S. Grimm, A. Bajwa, and L. Bourn. Medical imaging in a cloud computing environment. In *International Confer*ence on Cloud Computing and Services Science. to appear, 2011.
- [20] W. Schroeder, K. Martin, and B. Lorensen. *The Visualization Toolkit, Third Edition*. Kitware Inc., 2007.
- [21] D. Semeraro, R. Wilhelmson, D. Bramer, J. Leigh, D. Porter, and J. Ahrens. Collaboration, analysis, and visualization of the future. In *Proceedings of the 20th Conference on IIPS*. Amer. Meteor. Soc., 2004.
- [22] Y. Shahar, D. Goren-Bar, D. Boaz, and G. Tahan. Distributed, intelligent, interactive visualization and exploration of time-oriented clinical data and their abstractions. *Artificial Intelligence in Medicine*, 38(2):115 – 135, 2006. Temporal Representation and Reasoning in Medicine.
- [23] I. Taylor, M. Shields, I. Wang, and A. Harrison. The triana workflow environment: Architecture and applications. In I. J. Taylor, E. Deelman, D. B. Gannon, and M. Shields, editors, *Workflows for e-Science*, pages 320–339. Springer London, 2007.
- [24] P. S. Tofts, G. Brix, D. L. Buckley, J. L. Evelhoch, E. Henderson, M. V. Knopp, H. B. Larsson, T. Y. Lee, N. A. Mayr, G. J. Parker, R. E. Port, J. Taylor, and R. M. Weisskoff. Estimating kinetic parameters from dynamic contrast-enhanced T1-weighted MRI of a diffusable tracer: standardized quantities and symbols. *J Magn Reson Imaging*, 10(3):223–232, 1999.
- [25] I. Viola, A. Kanitsar, and M. E. Groller. Importance-driven volume rendering. In *Proceedings of the conference on Visualization '04*, VIS '04, pages 139–146, Washington, DC, USA, 2004. IEEE Computer Society.

Design Space of Network Security Visualization

Xiaoyuan Suo

Dept of Math and Computer Science, Webster University, Saint Louis, MO, USA

Abstract--Recent years have seen a growing interest in the emerging area of computer security visualization which is about developing visualization methods to help solve computer security problems. In this paper, we analyze the design space of network security visualization. Our main contribution is a new taxonomy that consists of three dimensions – data, visualizations, and tasks. Each dimension is further divided into hierarchical layers, and for each layer we have identified key parameters for making major design choices. This new taxonomy provides a comprehensive framework that can guide network security visualization developers to systematically explore the design space and make informed design decisions. It can also help developers or users systematically evaluate existing network security visualization techniques and systems. Finally it helps developers identify gaps in the design space and create new techniques.

Keywords: Security Visualization, Design Space, Network Security.

1 Introduction

In this paper, we analyze the design space of network security visualization by developing a new taxonomy. Within this taxonomy framework, we identify key parameters and classes that define the structure of this design space. Using taxonomy to define a design space is a common method that has been used in the area of computer security [1, 2], information visualization [3-5], and computer-human interaction [3, 6].

Our in-detailed analysis intends to provide the fellow developers with a better understanding of the structure of the design space. This would help network security visualization developers understand the tasks at a more detailed level and understand how various visualization techniques address the tasks. First the analysis can help to explore the design space and make informed design decisions. Second, it helps developers or users systematically evaluate existing techniques and systems. Third, it helps developers identify gaps in the design space and create new techniques.

Our main contribution is a taxonomy that consists of three dimensions – data, visualizations, and tasks. Each dimension is further divided into hierarchical layers, and for each layer we have identified key parameters for making major design choices, (figure 1).



Figure 1. 3 major dimensions that determine the design space of network security; and their sub-hierarchical layers.

The *data* dimension consists of two layers: raw data and transformed data. The *visualization* dimension is divided into five layers: workspace, view, visual structure, and visual unit and visual variable. The *task* dimension consists of two layers: high level task and low level task, (figure 1).

In the next three sections, we will discuss each dimension of the design space in sequence. We will explain the relationship between different layers, identify key parameters, and list categories of possible design choices for each parameter. Since design choices made in one dimension often affect the design choices in other dimensions, we will also discuss the relationship of a parameter or category with parameters or categories in other dimensions.

2 Data

The data dimension can be further divided into raw data and transformed data layers. Raw data becomes transformed data through data transformation operations.

2.1 Raw Data

In network security visualization, raw data is the network traffic data. Although a wide variety of network data are presented in network security visualization systems, the majority of them focus on the raw network traffic data at the transportation layer and network layer. From our survey of the network security visualization literature, we have identified the most commonly visualized raw data as follows.

- Source and destination IP addresses
- Source and destination port numbers

- Time and date
- Protocols (e.g. TCP or UDP)

Some applications also make use of application layer data such as user information and application type.

2.2 Transformed Data

When raw data is processed to fit into the proper visualization media; this process is defined as data transformation. More specifically, most network security visualization programs obtain their data from various log files such as IDS logs, net flows, sys-log, firewall logs, etc. The designated data for the specific visualization purposes maybe extracted or normalized from the raw network traffic data.

Presumably, there are many data transformation methodology; the most common ones are listed as follows:

- Classifying
- Counting and aggregation
- Filtering
- Sorting
- Clustering

Classifying. Network traffic can be classified in different ways. For example, it can be classified by the application types [7], such as WWW, mail, or multimedia, and some systems visualize such data. Intrusion detection systems (IDS) try to classify network traffic into normal and malicious categories [8].

Counting and aggregation can help to reduce the size of the data and therefore allow visualization systems to display the data at multiple levels of detail [9]. In addition, certain user groups may not want their network traffic details to be exposed, and therefore the aggregated data is the only thing that is allowed to be visualized.

Filtering can help reduce the data size and dimensions, eliminate noises and duplications, and help users focus on the important data. Most network log files are filtered by their application programs. However often in times, data filtering can be done interactively by users as well [10-12].

Sorting is a very common type of data transformation. For example, users may want to sort network data by time or by source IP. Sorting is particularly useful in table based visualizations.

Clustering is often used as part of a data mining process. Again, the clustering can be done automatically by programs or semi-automatically with user intervention.

Relationship with the visualization dimension

Both raw network traffic data (e.g. IP address, port number) and transformed data (e.g. IDS classifications) may be presented in the same visualization. The selection of visualization techniques are largely influenced by the type and number of dimensions of the data. For example, some visualization techniques, such as parallel coordinates, are particularly useful for visualizing multidimensional data. The volume of data is also an important factor in making visualization design decisions. Another important design decision is how to map different data attributes to different visual structures, visual units, and visual variables.

Relationship with the task dimension

Data transformation is closely related to tasks. When data transformation is integrated with the visualization system, data transformation operations become tasks. That is, interactive data transformations become user tasks, while noninteractive data transformations become developer/program tasks. Dynamically linked or synchronized – if one view is changed, the other views will be automatically updated. Based on Roberts[13], the key concepts of view coordination include the scope of the correlation, initiator, and what is correlated.

Multiple views may share the *same data source* or use *different data sources*. In the former case, the same data is visualized in different ways. For example, one view may display the raw network traffic data, while the other view displays aggregated data generated from the same raw network traffic. Multiple views that share the same data source are usually coordinated.

Relationship with the task dimension

The arrangement of multiple views is a low level task. Such arrangement can be determined by developers (a developer task), by the program (a program task), or by users (a user task). A particularly interesting design is to allow users to snap-together multiple views and dynamically link them [14, 15].

3 View

A view can be a window or a frame within a window. It contains one or more visual structures. The key parameters for the view layers are *contents* and *viewpoint*.

Contents can be *dynamic* or *static*. Dynamic contents means that the data visualized in the view may be changed during runtime. The update of the view content is a low level task that can be performed automatically by program (e.g. network data or IDS data streaming) or manually by users (e.g. load a different IDS log file).

Viewpoint can also be dynamic or static. Dynamic viewpoint means that the viewpoint can be manipulated, mostly by users, while static viewpoint means that viewpoint is fixed. Dynamic viewpoint is particularly useful when the entire data set is too big to be visualized in a view, which is a

typical problem for network security visualization due to the enormous size of network data.

Relationship with the task dimension

Viewpoint manipulations are low level tasks. Typical examples include zoom, pan, and focus+context [3].

3.1 Visual Structure

A visual structure is made up of one or more visual units. This layer determines how visual units are organized. A large number of information visualization techniques belong to this layer. Examples include bar chart, scatter plot, color map, parallel coordinates, tree and graph, etc.

We identify three key parameters for the visual structure layer: *coordinate systems*, *relationship among visual units*, and *space filling*.

In network security visualization, the most commonly used coordinate systems are *Cartesian* coordinate systems (both 2D and 3D), *Polar* coordinate systems [16, 17], and *Parallel* coordinate system [18]. The majority of the network security visualization systems use 2D coordinates, while a few systems use 3D coordinates.

There are three types of relationship among visual units: *no connection, hierarchical connection, and non-hierarchical connection*. Hierarchical connections are used to represent tree data structures (such as attack tree), while non-hierarchical connections are used to represent general graph data structures (such as computer networks) [19].

The visual units may be *space filling* or *non-space filling*. Space filling means that the entire visual structure is occupied by visual units. Non-space filling means that there may be gaps between visual units. Space filling techniques generally make better use of the display space, but they may suffer from information overloading. Note that table based visualizations are considered as space filling techniques when we treat each table cell as a visual unit.

A visual structure may encircle other visual structures [20]. For example, a table may contain scatter plots or bar charts in its cells. A bar chart may be placed on top of a geographical map.

Relationship with the data dimension

The selection of visual structures is largely determined by the nature, the number of dimensions, and the size of the data set. For example, if there are relationships among data elements, then connections need to be established among visual units. Node-link diagram is often used to visualize network connections [21]. Multi-dimensional data set may be visualized using parallel coordinates or visual pivot table [22]. The space filling dense pixel map is often selected to visualize IP address space because it can visualize a large amount of information in a small space.

A typical design problem in the visual structure layer is how to map network data to the axes of the coordinate system. A common practice is to convert IP address into two numbers, each of them mapped to one axis of the Cartesian coordinate. Port numbers are often sequentially mapped to pixels in a dense pixel map, either by column or by row. Temporal data is typically mapped to a horizontal or vertical axis [6].

Relationship with the task dimension

The selection of visual structure is a low level task. Often visual structures are selected by developers and are not changed during run time. However, in automatic visualization generation systems, the visual structure can be selected by the program based on pre-defined rules. Or the visual structure can be defined by users at run time.

3.2 Visual Units and Visual Variables

Visual units are the building blocks of visualization. Some examples such as: *point, line, 2D shape (glyph), 3D object, text,* and *image*. Each visual unit is defined by, but may not limit to, seven visual variables [21]: *position, size, shape, value, color, orientation,* and *texture*. For each visual unit, a visual variable is either assignable or fixed. An empty cell means the visual variable is not applicable to the visual unit.

An assignable variable visual means that data attributes can be mapped to this visual variable, otherwise, if it is fixed, data attributes cannot be mapped to it. For example, for dense pixel maps, IDS classification can be mapped to pixel colors, but not to its size [6].

| | position | Size | shape | value | color | orientation | text |
|----------|------------|------------|------------|------------|------------|-------------|------------|
| position | assignable | fixed | fixed | assignable | assignable | fixed | |
| line | assignable | assignable | fixed | assignable | assignable | assignable | |
| 2D | assignable | assignable | assignable | assignable | assignable | assignable | assignable |
| 3D | assignable | assignable | assignable | assignable | assignable | assignable | assignable |
| text | assignable | assignable | fixed | assignable | assignable | assignable | |
| image | assignable | assignable | fixed | fixed | fixed | assignable | |

Table 2 Visual units and visual variables

Sometimes the selection of a particular visual structure would limit the choices of visual units. For example, if parallel coordinate is selected as visual structure, then lines should be used as visual units. Similarly, the selection of a particular visual unit may also limit the choices of visual structure.

Perhaps we can add a special visual unit – *Gestalt*, as in Gestalt psychological theory. Here, a *Gestalt* is defined as a group of visual units that can be easily perceived by humans as a "pattern" due to the Gestalt theory – that is, the laws of proximity, similarity, symmetry, continuity, etc. *Gestalt* is particularly important for network security visualization because one of the primary purposes of such visualization systems is to help users detect malicious or anomalous network traffic patterns. Such patterns are often visualized as a *Gestalt* of pixels, lines, or glyphs. For example, in current network security visualization systems as a cluster of lines or a group of closely packed pixels with the same color. If a malicious or anomalous pattern is not visualized as a *Gestalt*, then it is usually difficult to be detected by human.

Thus a fundamental challenge for network security visualization designers is "how to design the visualization so that the malicious or anomalous behavior can be perceived by users as *Gestalt*?" And the network security visualization systems should be evaluated by whether they can effectively convert malicious or anomalous behavior patterns to *Gestalts*.

Unfortunately, most of the current research in network security visualization still focuses on the low level details of how to map network data to visual units and variables, and the high level task of mapping malicious patterns to Gestalt has not received much attention.

Relationship with the data dimension

A basic design question a developer would face is how to map data items to visual units and how to map data attributes to visual variables. Again, the selection of visual units and visual variables is largely determined by the nature, the number of dimensions, and the size of the data set. For example, Chernoff faces [23], a 2D shape, is often selected to visualize multi-dimensional data. For large volume of data set, pixel is often selected as the visual unit because it allows more data to be visualized in a small display space.

The selection of visual variables is also influenced by the characteristics of visual variables. Bertin [24] has identified five characteristics: selective, associative, quantitative, order, and length. For example, color is selective but not quantitative, meaning it is appropriate to map categorical data (such as IDS classification) to color [6], but it is usually not appropriate to map quantitative data to color.

The mapping of data to visual units and visual variables is a low level task. This visual mapping can be pre-defined by developers, automatically performed by programs based on certain rules [24, 25], or manipulated by users at run time.

5 Task

In the field of information visualization, tasks are often defined implicitly as the interactions between users and the visualization system. Here we define tasks in a broader sense.

We divide the task dimension into two layers: high level task and low level task. For network security visualization, we define high level tasks as the tasks that deal directly with problem solving, and define low level tasks as the tasks that indirectly support problem solving.

5.1 High Level Tasks

Based on our definition, high level tasks can be further divided into several categories: *problem detection, problem identification, diagnosis, problem projection, and problem response.*

With only a few exceptions [26, 27], the current research on network security visualization has been focusing on low level tasks. The only high level task that has been effectively addressed so far is problem detection – that is, to detect malicious or anomalous behavior patterns through visualization. Current network security visualization systems support two types of problem detection: *signature based* and *anomaly based detections*.

In signature based detection, users know the visual patterns (signature) of the malicious behavior and try to look for the patterns in the visualization. Such visual patterns are specific to visualization design (particularly visual structure, visual units and visual variables). As a result the visual signature of a malicious behavior (e.g. denial of service attack) is usually different from system to system. In anomaly based detection, users establish a visual profile of the normal behavior and use it to find anomalous visual patterns.

In summary, we still do not have a good understanding of the high level tasks of network security professionals and how visualization techniques can assist in their work. Much research needs to be done in this area.

5.2 Low Level Tasks

Low level tasks are mainly about information gathering and presentation. A key parameter for low level tasks is who initiates the task. The *initiators* can be users, developers, or program.

It is important to differentiate the task initiators because it helps guide the design. The visualization of network security data is often the result of a combination of design choices made by developers, users, and programs. In majority of the network security visualization systems, developers make most of the design choices. But in some information visualization systems [13, 28, 29], users construct the visualization at run time. In systems that automatically generate visualizations, design choices are made by programs.

Therefore, we can further classify low level tasks based on the dimensions and layers they are associated with. Table 2 contains such a classification. Note that the tasks listed in the table are example tasks and more tasks can be added.

| Table 4 Low level task classification | | |
|---------------------------------------|--|--|
| | Low level tasks | |
| Raw data | Add, delete, or change data source | |
| Transformed | filter, aggregate, classify, sort, cluster | |
| data | | |
| Workspace | Add, delete, arrange, or coordinate | |
| | multiple views | |
| View | Zoom, pan, overview, focus+context | |
| Visual | Add, delete, or modify relations, define | |
| structure | visual structure | |
| Visual unit | Identify, locate, distinguish, categorize, | |
| | cluster, rank, compare, associate, | |
| | correlate, retrieve, find anomalies | |
| Visual | Change visual mapping | |
| variable | | |

Table 4 Low level task classfication

Table 5 provided a task classification of the existing systems.

6 Related work

Various taxonomies have been proposed in the field of information visualization [3, 5, 30-40]. Both Shneiderman [41] and Wehrend [39] proposed a two dimension taxonomy based on data types and low level tasks. However, they do not provide a classification for the visualization techniques. Tory and Moller [32] extended Shneiderman's work by dividing data into two categories, continuous and discrete, and introduced structure as a parameter. They also describe a data centered, two dimensional task taxonomy based on continuous/discrete data and data spatialization. Card and Mackinlay [3] analyzed the information visualization design space based on data type, data transformation, mark types, retinal properties (similar to visual variables), position in space time, view transformation, and widget - all of them are included in our taxonomy. Card, et al. [3] and Chi [30] described a pipeline based framework to classify visualization techniques. In general, our taxonomy is closer to the one proposed by Keim [5], who describe taxonomy in three dimensions: data, visualization techniques, interaction and distortion techniques. Keim provides a categorical classification for each dimension, but does not further divide them into layers or identify key parameters. Our taxonomy is more comprehensive and detailed than other taxonomies.

Our analysis of the data dimension is influenced by the pipeline model proposed by Card, et al. [3] and Chi [30], who

both classify data into raw data and transformed data. Our analysis of the data transformation operations is influenced by the work of Tang, et al. [42]

Our analysis of the visualization dimension is based a comprehensive survey of the literatures on computer security visualization, many of which are published in the proceedings of the first two International Workshop on Visualization for Computer Security (VizSEC). The five layer hierarchical framework is partially inspired by Zhou and Feiner [27]. However, Zhou and Feiner discuss their hierarchical layers in the context of automatic visualization generation. Unlike our work, they do not identify key parameters and categories for each layer. The analysis of the visual unit and visual variable layers are influenced by the work of Card and Mackinlay [3] and Bertin [24].

There has not been a comprehensive analysis of the tasks in network security visualization. Our analysis of the low level tasks is based on our experience and many previous works on general visualization tasks [3, 22, 24, 30, 37, 39, 43-45]. Our high level task analysis of the network security visualization also benefits from some previous works [26, 27, 42].

7 Conclusion

In this paper we analyze the design space of network security visualization. We define the design space in three dimensions – data, visualization, and tasks. We further divide each dimension into hierarchical layers and identify key parameters and categories and design choices. In the process, we try to identify gaps in the design space and point to areas that have not been sufficiently addressed.

This new taxonomy provides a comprehensive framework that can guide developers to systematically explore the design space and make informed design decisions. It will lead to a more structured design methodology. It can also help developers or users systematically evaluate existing network security visualization techniques and systems.

We will continue to refine our taxonomy as new network security visualization techniques and systems are developed. New parameters may be identified and new categories added. We are also working to identify visual design patterns for network security visualization, based on our literature survey and the analysis of the design space. Finally, we believe our new taxonomy can be extended to cover the emerging field of visual analytics [46].

References

- D. Brumley, L.-H. Liu, P. Poosankam, and D. Song, "Design space and analysis of worm defense strategies," in *Proceedings of the 2006 ACM Symposium on Information, computer and communications security*, 2006.
- J. Mirkovic and P. Reiher, "A Taxonomy of DDoS Attack and DDoS Defense Mechanisms," ACM SIGCOMM Computer Communications Review, vol. 34, pp. 39-53, 2004.
- [3] S. Card, J. Mackinlay, and B. Shneiderman, "Readings in Information Visualization: Using Vision to Think," Morgan Kaufmann, 1999.
- [4] A. Komlodi, P. Rheingans, U. Ayachit, J. R. Goodall, and A. Joshi, "A User-centered Look at Glyph-based Security Visualization," in *Proceedings of the Workshop on Visualization for Computer Security*: IEEE, 2005.
- [5] D. A. Keim, "Information Visualization and Visual Data Mining," *IEEE Transactions on Visualization* and Computer Graphics, vol. 8, pp. 1-8, 2002.
- [6] D. Bowman, "The Science of Interaction Design," in *ACM SIGGRAPH Course Notes*, 2000.
- [7] A. W. Moore and D. Zuev, "Internet Traffic Classification Using Bayesian Analysis Techniques," in Proceedings of the ACM International Conference on Measurement and Modeling of Computer Systems (SIGMETRICS), 2005.
- [8] K. Abdullah, C. Lee, G. Conti, J. A. Copeland, and J. Stasko, "IDS RainStorm: Visualizing IDS Alarms," in *IEEE Symposium on Information Visualization's* Workshop on Visualization for Computer Security (VizSEC), 2005.
- [9] S. Noel and S. Jajodia, "Managing attack graph complexity through visual hierarchical aggregation," in *Proceedings of the 2004 ACM workshop on Visualization and data mining for computer security* Washington DC, USA: ACM Press, 2004.
- [10] C. Muelder, K.-L. Ma, and T. Bartoletti, "A Visualization Methodology for Characterization of Network Scans," in Workshop on Visualization for Computer Security (VizSEC 2005), 2005.
- [11] H. Koike, K. Ohno, and K. Koizumi, "Visualizing cyber attacks using IP matrix," in *Workshop on Visualization for Computer Security*, October 26, Minneapolis, MN, USA, 2005.
- [12] A. Komlodi, P. Rheingans, U. Ayachit, J. R. Goodall, and A. Joshi, "A User-centered Look at Glyph-based Security Visualization," in Workshop on Visualization for Computer Security, Minneapolis, MN, USA, 2005, pp. 21 - 28.
- [13] J. C. Roberts, "On Encouraging Multiple Views for Visualization," in *Proceedings of the International*

Conference on Information Visualization (IV): IEEE, 1998.

- [14] C. North and B. Shneiderman, "Snap-Together Visualization: A User Interface for Coordinating Visualization via Relational Schemata," in Proceedings of the working conference on Advanced visual interfaces: ACM Press, 2000.
- [15] C. North and B. Shneiderman, "Snap-together visualization: can users construct and operate coordinated visualizations?," *International Journal of Human-Computer Studies*, vol. 53, pp. 715-739, 2000.
- [16] G. Fink, P. Muessig, and C. North, "Visual Correlation of Host Processes and Network Traffic," in *IEEE Visualization 2005, Workshop on Visualization for Computer Security (VizSEC 2005)*, 2005, p. 8.
- [17] G. Fink and C. North, "Root Polar Layout of Internet Address Data for Security Administration," in *IEEE Visualization 2005, Workshop on Visualization for Computer Security (VizSEC 2005)*, 2005, p. 8.
- X. Yin, W. Yurcik, and A. Slagell,
 "VisFlowConnect-IP: An Animated Link Analysis Tool For Visualizing Netflows," in *FLOCON* -*Network Flow Analysis Workshop (Network Flow Analysis for Security Situational Awareness)* Baltimore, MD, 2005.
- [19] D. Yao, M. Shin, R. Tamassia, and W. H. Winsborough, "Visualization of Automated Trust Negotiation," in *Workshop on Visualization for Computer Security*, Minneapolis, MN, USA, 2005, pp. 65 - 74
- [20] J. R. Goodall, W. G. Lutters, and A. Komlodi, "I Know My Network: Collaboration and Expertise in Intrusion Detection," in *Proceedings of the ACM Conference on Computer Supported Cooperative Work (CSCW)*, 2004.
- [21] Y. Livnat, J. Agutter, S. Moon, R. F. Erbacher, and S. Foresti, "A Visualization Paradigm for Network Intrusion Detection," in *Proceedings of 6th IEEE Information Assurance Workshop* US Military Academy, West Point, New York: IEEE, 2005.
- [22] S. T. Teoh, K.-L. Ma, F. Wu, and T. J. Jankun-Kelly, "Detecting Flaws and Intruders with Visual Data Analysis," *IEEE Computer Graphics and Applications*, vol. 24, p. 27, 2004.
- [23] S. T. Teoh, K.-L. Ma, and S. F. Wu, "A visual exploration process for the analysis of Internet routing data," in *IEEE Visualization Conference*: IEEE, 2003, p. 523.
- [24] J. Bertin, *Semiology of Graphics*: University of Wisconsin Press, 1983.

- [25] H. Chernoff, "Using faces to represent points in kdimensional space graphically," *Journal of the American Statistical Association*, vol. 68, pp. 361-368, 1973.
- [26] M. X. Zhou and S. K. Feiner, "Top-Down Hierarchical Planning of Coherent Visual Discourse," in *Proceeding of International Conference on Intelligent User Interface (IUI)*: ACM, 1997.
- [27] M. X. Zhou and S. K. Feiner, "Visual Task Characterization for Automated Visual Discourse Synthesis," in *Proceedings of the ACM Conference* on Human Factors in Computing Systems (CHI), 1998.
- [28] H. Koike, K. Ohno, and K. Koizumi, "Visualizing cyber attacks using IP matrix," in *Proceedings of the Workshop on Visualization for Computer Security* (*VisSEC*): IEEE, 2005.
- [29] A. D'Amico and M. Kocka, "Information assurance visualizations for specific stages of situational awareness and intended uses: lessons learned," in *Workshop on Visualization for Computer Security*, Minneapolis, MN, USA, 2005, pp. 107 - 112
- [30] E. H. Chi, "A Taxonomy of Visualization Techniques using the Data State Reference Model," in *Proceeding of IEEE Symposium on Information Visualization (InfoVis)*, 2000.
- [31] A. Komlodi, J. Chakraborty, E. Stanziola, and A. Sears, "Information Visualization Evaluation Review Bibliography," 2005.
- [32] C. Stolte, D. Tang, and P. Hanrahan, "Polaris: A System for Query, Analysis and Visualization of Multi-dimensional Relational Databases," *IEEE Transactions on Visualization and Computer Graphics*, vol. 8, 2002.
- [33] J. Lohse, H. Rueter, K. Biolsi, and N. Walker, "Classifying Visual Knowledge Representations: A Foundation for Visualization Research," in *Proceeding of IEEE Visualization Conference (VIS)*, 1990.
- [34] E. H. Chi and J. T. Riedl, "An Operator Interaction Framework for Visualization Systems," in *Proceeding of IEEE Symposium on Information Visualization (InfoVis)*: IEEE, 1998.
- [35] M. C. Chuah and S. F. Roth, "On the Semantics of Interactive Visualizations," in *Proceeding of IEEE* Symposium on Information Visualization (InfoVis): IEEE, 1996.
- [36] G. L. Lohse, K. Biolsi, N. Walker, and H. H. Rueter, "A Classification of Visual Representations," *Communications of the ACM*, vol. 37, pp. 36-49, 1995.
- [37] C. North and B. Schneiderman, "A Taxonomy of Multiple Window Coordinations," Technical Report CS-TR-3854, Department of Computer Science, University of Maryland1998.

- [38] S. F. Roth and J. Mattis, "Data Characterization for Intelligent Graphics Presentation," in *Proceedings of the ACM Conference on Human Factors in Computing Systems (CHI)*: ACM, 1990.
- [39] S. Wehrend and C. Lewis, "A Problem-oriented Classification of Visualization Techniques," in *Proceedings of the IEEE Symposium on Information Visualization (InfoVis)*: IEEE, 1990.
- [40] H. Shiravi, A. Shiravi, and A. A. Ghorbani, "A Survey of Visualization Systems for Network Security," *IEEE TRANSACTIONS ON VISUALIZATION AND COMPUTER GRAPHICS*, vol. 1, 2011.
- [41] B. Shneiderman, "The Eyes Have It: A Task by Data Type Taxonomy for Information Visualizations," in Proceedings of the IEEE Conference on Visual Languages: IEEE, 1996.
- [42] D. Tang, C. Stolte, and R. Bosch, "Design choices when architecting visualizations," *Information Visualization*, vol. 3, pp. 65-79, 2004.
- [43] J. C. Roberts, "Display Models for Visualization," in Proceedings of IEEE Conference on Information Visualization (IV), 1999.
- [44] D. A. Keim, "Visual Exploration of Large Data Sets," *Communications of the ACM*, vol. 44, pp. 39-44, 2001.
- [45] M. Tory and T. Möller, "Rethinking Visualization: A High-Level Taxonomy," in *Proceeding of the IEEE* Symposium on Information Visualization (InfoVis), 2004.
- [46] R. Amar, J. Eagan, and J. Stasko, "Low-Level Components of Analytic Activity in Information Visualization," in *Proceedings of the IEEE* Symposium on Information Visualization (InfoVis): IEEE, 2005.

Visualization and Clustering of Document Collections using a Flock-based Swarm Intelligence Technique

Richard H. Fowler, Raul A. Huerta, and Wendy A. L. Fowler

Department of Computer Science, University of Texas - Pan American, Edinburg, TX, USA

Abstract – Electronic availability of documents continues to increase, yet identifying documents relevant to the user remains a primary constraint in electronic document use. Visual representations of document collections can facilitate search by representing large collections of documents in a manner that is complementary to linear, text based representations. Visual representations can provide a means to make the overall structure of a collection comprehensible, as well as a mechanism to identify groups of useful documents and access relevant individual documents. The current work employs flock-based clustering to both organize documents and provide visual representations of documents. Reynolds' three rule flocking scheme is augmented with additional rules to provide document clustering. A unified visual representation supplies facilities for overview of the entire document collection, filtering documents, and retrieving individual documents. The system also utilizes visualization tools for individual cluster identification and exploration based on keyword search. Stereoscopic viewing is provided to enhance users' perception of 3D organization.

Keywords: visualization, information visualization, information search, clustering, swarm intelligence

1 Introduction

Electronic availability of documents continues to increase both with respect to types of access and breadth of coverage. Yet, identifying those documents that are relevant to the user for the task at hand remains the primary constraint in electronic document use. Most systems provide access using keyword search, whether using publicly available Internet search engines or digital libraries accessed by professionals. Visual representations of document collections during search can augment text-based techniques by representing large collections of documents in a manner that is complementary to linear, text based representations. Visual representations can provide a means to make the overall collection comprehensible, as well as mechanisms to identify groups of useful documents and access relevant individual documents.

In information visualization tasks the essential approach is to provide "overview first, zoom and filter, then details-ondemand" [24]. The approach to visualization of document collections we have taken follows that paradigm closely. The current work describes a system that provides a visual representation of a complete document collection formed using a biologically inspired flock-based [21] clustering technique [4, 5]. The technique provides a means to not only form clusters of documents, but also spatially order the clusters and documents within clusters. Together, these spatial orderings can provide the user both a global view of the document collection, as well as the ability to view relations at a more detailed intra-cluster level. Simultaneous cluster formation and spatial arrangement is efficient by eliminating the need for separate computational stages.

The system we have developed employs flock-based clustering to both organize documents with respect to content and to provide visual representations of documents. The basic clustering approach whereby additional rules are added to Reynolds' three rule flocking scheme [21] was augmented and tuned for the web based document sets with which the work was done. A single visual representation provides users facilities to gain an overview of the entire document collection, filter the document collection, and obtain information about individual documents. The system also visualization tools for individual provides cluster identification and exploration based on keyword search. Both desktop and large screen stereoscopic viewing facilities are provided to enhance the users' perception of three dimensional organization.

The following sections first present Related Work in swarm intelligence used to form clusters, focusing on flockbased techniques. This section also describes the metrics utilized to represent individual documents and document collections and extensions to Reynolds' model that have been employed with document collections. The next section describes the System that we have developed for flock-based visualization of document collections. Finally, Conclusions are presented and References listed.

2 Related Work

Clustering is the process of assigning a set of objects into groups, or clusters, so that the objects in the same cluster are more similar to each other than to those in other clusters [11]. For document retrieval the clustering of documents in a collection is widely used to facilitate user directed search through browsing [14]. Clustering techniques that make use of flock-based swarm intelligence techniques are relatively new and provide a useful adjunct to other techniques, such as hierarchical and nearest neighbor clustering.

2.1 Flock-based Modeling

Reynolds' seminal work in creating flock-like visual behaviors of birds [21] was designed to provide a technique for use in computer animation. As such, its goal was to provide a visual representation of flight and the grouping of individuals, i.e., flocking, that would appear realistic to viewers, rather than provide a biologically accurate model of behavior. Reynolds' agent-based technique was quite successful using a very limited number of rules. The basic flocking approach was quickly extended to a wide range of groups, including herds [7], schools of fish [9], crowds [29], unmanned air vehicles [30], and robots [31]. Further refinements have added additional rules to provide for flocklike behavior in the presence of predators, obstacles to avoid, and path following [22]. The approach has also been extended to multiple species flocking [5, 17] using techniques closely related to those for document clustering. The next sections present details of the original flocking model from which later extensions are derived, followed by a discussion of extensions to the basic model with a focus on clustering and document collections.

2.1.1 Simulating Flock Motion

Reynolds' original technique to create perceptually realistic flocking behavior for animation was based on providing rules for movement, or steering, for individual agents. Only three rules are used, illustrated in Figure 1, which are applied by each individual of the flock to steer its movement (direction and velocity):

- Alignment: Steer towards the average direction of movement of nearby agents
- Separation: Steer to avoid being too close to nearby agents
- Cohesion: Steer to move toward the center position of nearby agents.

Using these three rules, individual agents initially placed at any spatial location will adjust their directions of movement and velocities to form a single group that exhibits movement perceptually quite similar to a flock of birds or other social biological group, e.g., school of fish. These steering rules are applied considering only other agents within a small range around each individual, as shown by the circles in the figure, rather than to all other agents. A velocity vector is formed for each of the three rules, alignment, separation, and cohesion. These vectors are then used to determine the individual agent's velocity.



Figure 1. Only three basic rules of movement are necessary to create perceptually realistic flocking motion. The filled triangle represents the agent applying the rules, outline triangles show other agents, and the circle represents the maximum spatial range of rule evaluation.

The alignment rule orients the direction of movement for an agent with the direction of nearby agents. Nearby agents are considered those that are within the range of evaluation, $d(P_x, P_a) \le d_1 \cap d(P_x, P_a) \ge d_2$, where *d* is distance, *P* is position, *x* denotes another agent, and *a* the agent for which alignment is being determined. d_1 and d_2 determine the other agents that are considered in determining alignment and are pre-defined. d_1 represents the maximum range of rule evaluation, and d_2 is used to exclude the nearest other agent's velocity vector, \vec{v}_{align} , is aligned with the average of the velocity vectors \vec{v}_x for the *n* agents within the range of evaluation:

$$\vec{v}_{align} = \frac{1}{n} \sum_{x=1}^{n} \vec{v}_x \tag{1}$$

The separation rule keeps an agent from colliding with another agent. It does this by changing the agent's velocity vector, \vec{v}_{sep} , depending on distance from agents in the range of evaluation, $d(P_x, P_a) \le d_2$:

$$\vec{v}_{sep} = \sum_{x=1}^{n} \frac{\vec{v}_x + \vec{v}_a}{d(P_x, P_a)}$$
 (2)

The cohesion rule orients the movement of an agent toward the center of nearby agents. The agent's velocity vector, \vec{v}_{coh} , is oriented to the direction of the average spatial position of the agents within its evaluation range of evaluation, $d(P_x, P_a) \le d_1 \cap d(P_x, P_a) \ge d_2$:

$$\vec{v}_{coh} = \sum_{x=1}^{n} \overline{(P_x - P_a)}$$
(3)

Finally, the velocity vector for each agent, \vec{v}_a , is calculated by summing and weighting the velocity vectors calculated by the three rules:

 $\vec{v}_a = w_{align} \cdot \vec{v}_{align} + w_{sep} \cdot \vec{v}_{sep} + w_{coh} \cdot \vec{v}_{coh}$ (4) with w_{align} , w_{sep} , and w_{coh} pre-determined weights.
2.2 Extensions of Flock Based Modeling to Document Collection Clustering

By the agents' iteratively applying the three rules of alignment, separation and cohesion, a single group forms that moves in a perceptually good approximation of group movement of identical biological entities in nature. However, there are other cases of interest that the rules do not capture. For example, in order to create movement of herd animals across terrain, Gompert [7] augments Reynolds' three rules modeling bird flight with a fourth rule whereby each agents' position also determined by the elevation of the terrain at the agent's position. Such derivation of rule sets appropriate to the problem at hand is characteristic of agent based modeling [15].

In addition to single group modeling, another question that has been addressed using flock-based modeling is differentiation of the members of a single group of agents into distinct groups. This would occur, for example, when agents model birds, but the birds are of different species. In this case each species would form its own separate group, and each separate group would exhibit the movement patterns captured by Reynolds' three rules. Solutions to this problem add additional rules whereby agents, e.g., birds, are brought closer together or pushed farther apart depending on their similarity. One way to determine similarity values is by comparing feature vectors representing individual agents. Agents of the same species share many features and so have high similarity, which leads to their spatial positions moving closer. These groups of similar agents are also moved away from other, dissimilar, agents, which have also formed groups based on high inter-agent similarity. This general approach, in which additional rules are added to flock-based clustering, has been used to provide multiple groupings, or clusters, of individual interests [20], time varying data [17], arbitrary attributes [19, 26], and spatial data [6].

Clusters of documents can be identified using the same approach in which additional rules are added that consider similarity among feature vectors [5, 6]. For document clustering, additional rules can be added that consider the similarity of the topic content of documents, where the feature vector is a vector of terms used to describe a document's content. The next section describes techniques for determining document feature vectors and their similarity. In the following section techniques for deriving document clusters based on inter-document similarity are presented.

2.2.1 Document Feature Vectors and Similarities for Clustering

The most widely used measures of similarity among documents are based on the Vector Space Model [23]. This technique utilizes a document's words to transform each document to a feature vector representation that captures the document's content. Comparisons among documents' feature vectors are then employed to provide document similarities.

The Vector Space Model uses a list of indexing terms defined for a particular document collection. These terms can come from a fixed vocabulary or be derived for individual document sets. The text of each document is analyzed, and a vector, d, representing each document is created. d is of length equal to the number of indexing terms. Each element of d is given a weight, w, for each indexing term, i. Using this method, each document can be considered as a single point in a space of dimensionality equal to the number of index terms.

Typically, inter-document similarity considers the number of terms common to a document pair as represented by their term vectors, a measure of their content similarity, normalized by number of terms in documents. In order to increase the discriminative power of terms, terms in the term vector T are first weighted by the inverse frequency of occurrence in the term set, the $tf \times idf$ model [27]. The idea is that relatively rarely occurring terms are more useful in characterizing inter-document similarity, and, so, are differentially weighted. Equation 5 below provides the weight, w, for a term in the document term vector. tf is a term's frequency in the document collection, and tf_{ik} is the number of documents in the collection and n_k represents the number of documents containing term T_k . The weight, w, of term T_k is:

$$w = \frac{tf_{if} \times \log\left(\frac{N}{n_k}\right)}{\sqrt{\sum_{j=1}^{T_n} (tf_{if})^2 \times \log\left(\frac{N}{n_j}\right)^2}}$$
(5)

2.2.2 Flock-based Document Clustering

The earliest use of flock-based clustering was as one element of a hybrid clustering approach that used conventional techniques together with flock-based techniques for final determination of clusters [6]. At that time other swarm based techniques for clustering were also being explored, including ant colony optimization [13] and particle swarm optimization [16]. The first use of flocking-based techniques for document collections was by Cui et al. [5, 6]. Refinement of the approach continues, with recent work on efficient implementation demonstrating the clustering of 500,000 documents [32]. Contemporary with the work of Cui et al, Picarougne and colleagues developed a flock-based clustering system with visualization facilities [18, 19].

As noted, the basic approach of flock-based document clustering is to augment Reynolds' original three agent rules with additional rules. Cui et al. [5] utilized a document's term vector representation as the feature vector from which values for inter-document similarities were calculated, and these were then used in additional rules. To determine document clusters Cui et al. use two variations of rules to augment Reynolds' original approach to determining spatial placement of agents. The basic goal is to move similar documents together and dissimilar documents apart by adding two rules. In the first rule, strength of attraction, \vec{v}_{sim} , is proportional to the distance between the agents and the similarity between the agents' term values:

$$\vec{v}_{sim} = \sum_{x=1}^{n} \left(S(A, X) \times d(P_x, P_a) \right) \tag{6}$$

where S(A,X) is the similarity value determined using the term vectors, or feature sets, for agents A and X. Similarly, the strength of repulsion, \vec{v}_{dis} , is inversely proportional to the distance between the agents and the similarity between the agents' features:

$$\vec{v}_{dis} = \sum_{x=1}^{n} \frac{1}{S(A, X) \times d(P_x, P_a)}$$
 (7)

As before, the velocity vector for each agent, v_a , is calculated by summing the weighted velocity vectors calculated by the, now five, rules:

 $\vec{v}_a = w_{al} \cdot \vec{v}_{al} + w_s \cdot \vec{v}_s + w_c \cdot \vec{v}_c + w_{sim} \cdot \vec{v}_{sim} + w_{dis} \cdot \vec{v}_{dis} \quad (8)$ with w_{al} , w_s , w_c , w_{sim} , and w_{dis} pre-determined weights.

This basic approach was refined [6] by combining the two additional rules into a single rule that uses a variable threshold, T, by which attractive and repulsive forces can be manipulated based on agent feature similarity:

$$\vec{v}_{sim-threshold} = \sum_{x=1}^{n} \frac{(S(A,X) - T) \times \overline{(P_x - P_a)}}{d(P_x, P_a)}$$
(9)

There is now a single velocity vector based on document similarity, v_{ds} , that depends on the pre-defined threshold, *T*, to determine document similarity based agent movement. Such changes to rules are characteristic of extensions made to the basic flock based algorithm.

3 System

The system we have developed provides a visual representation of document collections based on flock-based clustering as one element of a document retrieval system. To perform the clustering, documents are first represented as term vectors. Then, using flock-based modeling that extends Reynolds' original approach, documents with high similarity move together and dissimilar documents move apart in a three dimensional space. These movements create spatial clusters of documents that share content, as represented by their term vectors, and clusters move apart spatially to fill the space that users view, as shown in Figure 2 for a set of web retrieved documents. In addition to the main viewing window, views from the top and sides of the document space are also available to provide orientation. Figure 3 shows the



Figure 2. Flock-based clustering for a document collection. User's are supplied tools to browse the collection, identify clusters of interest, and retrieve individual documents.



Figure 3. User has zoomed in on the two clusters at the lower right of the display. Further zooming allows selection of individual documents for display.

view after the user has zoomed in on two of the clusters. Further zooming allows the inspection and retrieval of individual documents.

Stereoscopic viewing is available to enhance perception of the three dimensional space in which documents are displayed. In various task-based studies user's performance has been shown to be enhanced through stereoscopy [1, 8, 28]. Additionally, the user can interactively change the view orientation, e.g., "spinning" and "jittering" the display, to discern further information about the nature of the spatial clustering [25].

As with other flock-based document clustering systems [5, 32], additional rules for agent movement based on document similarities are used by the system to augment Reynolds' three principle steering rules. A variant of Cui et al.'s threshold-based similarity rule [6] is used for clustering, providing a threshold for document separation that can be adjusted to tune cluster formation and separation for different document sets. The flock-based clustering employed by the system represents documents using the $tf \times idf$ method, using the most frequently occurring terms in document sets. Similarity between agents based on feature vectors is

calculated as $sim_{ij} = \sum_{k=1}^{T} w_{ik} \times w_{jk}$, where inter-document similarity, *sim*, for each document pair, *i*, *j*, for each of *T* elements in the term vector, is derived from term weights, *w*, for agent pairs. For efficiency it is calculated once prior to initiating flocking and stored, rather than recalculated on each iteration.

One tenet of information search is that users should be provided multiple paths of access to information [2]. Though the capabilities of the system center on flock-based clustering and its visual representation, conventional keyword search for individual documents is also available and provides a useful addition to cluster based search [3, 12]. The system's flockbased clustering visualizations provide one element of its document retrieval functions. Other elements are designed to support an iterative document retrieval process in which users' information needs are defined and met [10].

In addition to retrieving individual documents through keyword search, the system also provides facilities that combine keyword search and visual cluster representation. The goal is to help users know where to look and explore within the large visual representation in order to find clusters that contain documents likely to be relevant to the user's information needs. This cluster identification is accomplished by augmenting the cluster display by visually marking individual documents that match a keyword search. Users are then visually directed to those clusters with many keyword matches and which contain similar documents, as indicated by common cluster membership.

Another means by which flock-based cluster visualizations could provide capabilities integrated with a suite of document retrieval mechanisms is by using it to provide an alternative representation of search results provided by a conventional search engine. Internet search engines return results ranked with respect to relevance to a keyword based query. Were the query able to express the user's information need exactly and the retrieval mechanism able to then supply the documents that met the user's information need, then there would be little to desire in such a system. Unfortunately, this is not the case, due in part to users' inability to completely specify information needs in terms used by the retrieval system. Rather, it is more likely that only some degree of the user's need is met with initially retrieved documents. Further refinement of information needs and search vocabulary are parts of the iterative process information retrieval. By using the flock-based visualization system to provide clustered, versus sequential, ordering of documents within a retrieved set, the user could be provided a visual mechanism complementary to the sequence of documents to find relevant documents through exploration of the clusters of documents formed from a set of retrieved documents. Additionally, flock-based clustering is particularly efficient for incremental clustering [32], and the set of visually displayed document can be increased and maintained as the search continues.

4 Conclusions

The current work extends the use of flock-based clustering visualization in document retrieval through its integration with tools supporting the iterative information retrieval process. The system provides mechanisms for cluster identification by keyword query match to identify individual documents and show their location in the complete document collection, thereby enabling users to efficiently explore the document space. This exploration facilitates user information need specification, an important component in the retrieval process, as well as individual document retrieval. The system provides multiple paths to information items through its facilities for document collection browsing, cluster identification, and keyword based search.

5 Acknowledgements

This work was supported by National Science Foundation grant MSI-CIEC OCI-0636352 and the University of Texas – Pan American Computing and Information Technology Center.

6 References

[1] K. W. Arthur, K. S. Booth, and C. Ware, "Evaluating 3D task performance for fish tank virtual worlds," *ACM Transactions on Information Systems*, vol. 11, no. 3, pp. 239-265, 1993.

[2] M. J. Bates, "Subject access in online catalogs: A design model," *Journal of the American Society for Information Science*, vol. 37, no. 6, pp. 357-386, 1986.

[3] D. B. Crouch, C. J. Crouch, and G. Andreas, "The use of cluster hierarchies in hypertext information retrieval," in *Proceeding of Hypertext* '89, pp. 225-237, 1989.

[4] X. Cui, J. Gao, and T. E. Potok, "A flocking based algorithm for document clustering analysis," *Journal of Systems Architecture*, vol. 52, no. 8-9, pp. 505-515, 2006.

[5] X. Cui and T. E. Potok, "A distributed agent implementation of multiple species flocking model for document partition clustering," *CIA 2006, Lecture Notes in Computer Science*, vol. 4149, pp. 124-137, 2006.

[6] G. Folino and G. Spezzano, "An adaptive flocking algorithm for spatial clustering," *Parallel Problem Solving in Nature (PPSN) VII, Lecture Notes in Computer Science*, vol. 2439, pp. 924-933, 2002.

[7] J. Gompert, "Real-time simulation of herds moving over terrain," *Proceedings of Artificial Intelligence and Digital Entertainment*, pp. 149-150, 2005.

[8] N. Greffard, F. Picarougne, and P. Kuntz, "Visual community detection: An evaluation of 2D, 3D perspective and 3D stereoscopic displays," *Proceedings of 19th International Symposium on Graph Drawing, Lecture Notes in Computer Science*, vol. 7034, pp. 215-225, 2011.

[9] Y. Inada, "Steering mechanism of fish schools," *Complexity International*, vol. 8, pp. 1-8, 2001.

[10] P. Ingwerson and I. Wormwell, "Improved subject access, browsing and scanning mechanisms in modern online ir," *Proceedings of ACM SIGIR*, pp. 68–76, 1986.

[11] A.K. Jain, M.N. Murty, and P.J. Flynn, "Data clustering: A review," *ACM Computing Surveys*, vol. 31, no. 3, pp. 264-323, 1999.

[12] N. Jardine and C.J. van Rijsbergen, "The use of hierarchical clustering in information retrieval," *Information Storage and Retrieval*, vol. 7, pp. 217–240, 1971.

[13] N. Labroche, N. Monmarche', G. Venturini, "AntClust: Ant clustering and web usage mining," *Proceedings of Genetic and Evolutionary Computation Conference*, pp. 25– 36, 2003.

[14] A. Leuski, "Evaluating document clustering for interactive information retrieval," *Proceedings of the International Conference on Information Knowledge and Management (CIKM)*, pp. 33-40, 2001.

[15] C. M. Macal and M. J. North, "Tutorial on agent-based modeling and simulation," *Journal of Simulation*, vol. 4, pp. 151-162, 2010.

[16] V.D. Merwe and A.P. Engelbrecht, "Data clustering using particle swarm optimization," *Proceedings of IEEE Congress on Evolutionary Computation*, pp. 215–220, 2003.

[17] A.V. Moere, "Information flocking: time-varying data visualization using Boid behaviors," *Proceedings of the Eighth International Conference on Information Visualization*, pp. 409–414, 2004.

[17a] S. Momen, B. P. Amavasai, and N. H. Siddique, "Mixed Species Flocking for Heterogeneous Robotic Swarms," *EUROCON 2007 The International Conference on Computer as a Tool*, pp. 2329-2336, 2007.

[18] F. Picarougne, H. Azzag, G. Venturini, and C. Guinot, "On data clustering with a flock of artificial agents," *Proceedings of the 16th IEEE International Conference on Tools with Artificial Intelligence (ICTAI 2004)*, pp. 777-778, 2004.

[19] F. Picarougne, H. Azzag, G. Venturini, and C. Guinot, "A new approach of data clustering using a flock of agents," *Evolutionary Computation*, vol. 15, no. 3, pp. 345-367, 2006. [20] G. Proctor and C. Winter, "Information flocking: Data visualisation in virtual worlds using emergent behaviours," *Proceedings of Virtual Worlds*, pp. 168–176, 1998.

[21] C. Reynolds, "Flocks, herds, and schools: A distributed behavioral model," *Computer Graphics*, vol. 21, no. 4, pp. 25-34, 1987.

[22] C. Reynolds, "Steering behaviors for autonomous characters," *Proceedings of Game Developers Conference*, pp. 763–782, 1999.

[23] G. Salton, C. Yang, and A. Wong, "A vector space model for automatic indexing," *Communications of the ACM*, vol. 18, no. 11, pp. 613-620, 1975.

[24] B. Shneiderman, "The eyes have it: A task by data type taxonomy for information visualizations," *Proceedings of IEEE Symposium on Visual Languages*, pp. 336-343, 1996.

[25] B. W. van Shooten, E. M. A. G. van Dijk, E. Zudilova, A. Suinesiaputra, and J. H. C. Reiber, "The effect of stereoscopy and motion cues on 3D interpretation task performance," *Proceedings of the International Conference on Advanced Visual Interfaces*, pp. 167-170, 2010.

[26] E. Sklar, C Jansen, J. Chan, and M. Byrd, "Toward a methodology for agent-based data mining and visualization," *International Workshop on Agents and Data Mining Interaction (ADMI 2011)*, pp. 20-31, 2011.

[27] K. Sparck-Jones, "A statistical interpretation of term specificity and its application in retrieval," *Journal of Documentation*, vol. 28, no. 1, pp. 11–20, 1972.

[29] W. Tang, T. R. Wan, and S. Patel, "Real-time crowd movement on large scale terrains," *Proceedings of the Theory and Practice of Computer Graphics (TPCG'03)*, pp. 146-153, 2003.

[29] C. Ware and G. Franck, "Evaluating stereo and motion cues for visualizing information nets in three dimensions," *ACM Transactions on Graphics*, vol. 15, no. 2, pp. 121-140.

[30] N. R. Watson, N. W. John, and W. J. Crowther, "Simulation of unmanned air vehicle flocking," *Proceedings of Theory and Practice of Computer Graphics*, pp. 130-137, 2003.

[31] N. Xiong, J. He, J. H. Park, T. Kim, and Y. He, "Decentralized flocking algorithms for a swarm of mobile robots: Problem, current research and future directions," *6th IEEE Consumer Communications and Networking Conference (CCNC 2009)*, pp.1-6, 2009.

[32] Y. Zhang, F. Mueller, X. Cui, and T. Potok, "Dataintensive document clustering on graphics processing unit (GPU) clusters," *Journal of Parallel and Distributed Computing*, vol. 71, pp. 211-224, 2011.

Which Book Should I Pick?: Text Visualization Based on Readability and Genre

H. Kim¹, J. W. Park²

^{1, 2}GSAIM, Chung-Ang University, Seoul, Rep. of Korea

Abstract - This paper proposes readability visualization, genre visualization, and combined visualization to provide unconventional information for book selection. Data visualization was initiated for the practical purpose of delivering information, as it efficiently links visual perception and data so that readers are able to instantly recognize patterns in overcrowded data. In this interdisciplinary research we used the strength of data visualization, and this paper suggests three possible textual visualizations of a book, which may help users to find a desirable book, with the use of intuitive information out of a large volume of book data.

Keywords: Data Visualization, Information Visualization, Information Aesthetics

1 Introduction

Today, internet bookstores sell almost every single book in the world, and they provide a convenient book-searching function using keyword entry. Moreover, they deliver purchased books to the customers' home directly. Buying books over the internet has gradually become a part of our lives.

Internet bookstores provide information about books in order to help customers pick their desired books, in-formation such as a cover image, table of contents, and book reviews, in addition to basic information: author, editor, number of pages, a size, ISBN, and price. Customers also get information from other customers' opinions or statistical recommendations provided by internet bookstores, which are sometimes leading readers to follow other people's preferences blindly. In other word, these opinions are not information centered on the individual user but on the provider. Ideally, to choose a book, customers should read one or more chapters in advance of purchase to decide if a book's style and genre suit them and also to predict how much time and effort it will take from them to read it.

Many previous examples of text visualization including our recent projects[1][2] have generally focused on a book's contents and on plot analysis [3][4][5] but these approaches hardly provide practical clues needed for a reader to pick the "right book" for them, out of the thousands of recommendations on an Internet book store's web site.

Common readers' questions are not academic but practical: How much time will I need to read this book? Is it easy enough for a third-grade kid? Is there a genre called philosophical thriller? Is Twilight a romance or fantasy? Where I can find a romantic philosophical science fiction like Bicentennial Man? This paper suggests that intuitive text visualization reflecting readability and a user-oriented customizable genre would complement available book information and support reader book selection. (Fig. 1)



Fig. 1. The images of text visualization with an algorithm suggested in this paper: from the left, Harry Potter and the Prisoner of Azkaban (3rd book), Harry Potter and the Goblet of Fire (4th book), both by J. K. Rowling; The Critique of Pure Reason (English edition) by Immanuel Kant; and The Ethics (English edition) by Baruch de Spinoza.

2 Readability Visualization

2.1 Readability

A book's readability is one of the important considerations in book selection. "Readability" is the ease in which text can be read and understood, and it directly affects the difficulty of a book. The factors of readability measurement have been studied by various researchers [6][7][8][9][10][11][12][13][14] since Sherman L. A. (1888)[15] (Table 1).

The length of a sentence, number of sentences per paragraph, and amount of characters per word; these factors are closely related to the physical print space and the eye movements readability issue of "return sweep to next line[16]".

| Kitsfigon (1921) | syllables per word and length of a sentence |
|-----------------------------|---|
| Lively & Pressey (1923) | level of difficulty and frequency of words |
| Vogel & Washburne (1928) | number of different words in a paragraph and number of prepositions |
| Gray & Leary (1935) | grammatical complexity with various factors including average sentence length, percentage of easy words, and number of sentences per para- graph |
| Gunning Fog (1952) | number of complex words, length of a sentence, and percentage of three or more syllables words |
| Dale & Chall(1948) | length of a sentence with word count |
| Flesch(1948) | total syllables in 100 words and average number of words in a sentence |
| Bormuth(1975) | average characters in a word and average num- ber of words in a sentence |
| Harris & Jacobson (1975) | percentage of difficult words and average sen- tence length with word count |

Table 1. Readability Factors

2.2 Readability visualization with length factors

We visually simulated a human reading process using the length related readability factors. The longer word, sentence, and paragraph are less readable. Each character was converted to a visible point of which the brightness decreases gradually during reading process. If a character is punctuation, the brightness of the point slightly increases. As a result, a book with longer sentences/ paragraphs presents a dark and static image, but a book with concise sentences/ paragraphs shows bright and dynamic patterns (Table 2, Fig. 2).

```
set the base hue, saturation, brightness (0, 0, Max)
while(not end of the book) {
  get a next paragraph
  while(not end of the paragraph) {
    get a next character
    if (the character is punctuation or comma)
        increase brightness (sentence formula)
    else {
        decrease brightness (paragraph formula)
    }
    if (three or more syllables words) // optional
        decrease brightness (word formula)
    set a pixel on next location with assigned hue,
        saturation, brightness and draw
    }
    reset brightness to Max
}
```

'So you did, old fellow!' said the others.

We must burn the house down!' said the Rabbit's voice; and Alice called out as loud as she could, If you do. I'll set Dinah at you!'

Fig. 2. A sample readability visualization for Alice's Adventures in Wonderland by Lewis Carroll. To explain the algorithm, we accelerated the brightness transition ten times.

Fig. 3 shows a clear difference between two books. My Sister's Keeper has clear and concise sentence and paragraph style, as is often the case with most of the contemporary best sellers, but The Critique of Pure Reason is extremely hard to understand, even for well educated people not only because of the difficulty of contents but also because of the lengthy complex sentences and paragraphs. With this algorithm a reader can get an impression of overall readability instantly without reading the entire book. For example, the readability visualization of most children's books that have more dialogic contents will appear extremely bright and dynamic. Classic novels are darker than contemporary novels due to the lengthy writhing style. Most of philosophical, scientific, or academic books are much darker than books of others genres. The visualization matches with common knowledge of readability and comprehensibility. One of the rare exceptions to this is poetry. For example an image of a poem by T. S. Elliot is bright but the words in the poem are highly connotative and implicative. Sometimes readability does not accompany comprehensibility.



Fig. 3. Result of the algorithm 1; left is My Sister's Keeper by Jodi Picoult, right image is The Critique of Pure Reason (English edition) by Immanuel Kant

3 Genre Visualization

The readability visualization of previous chapter is helpful but an even more important function for book selection would be to visualize a book based on each individual's preference. This chapter suggests an intuitive visualization algorithm that visualizes a book based on genre; either conventional or custom genre.

The most of previous text clustering studies focused on finding dimensional closeness with various method such as digital signal processing with wavelet transform [17][18][19][20]. In this paper we used text analytic method with viewpoint of normal customers.

Usually genres are defined by a librarian with a historical categorization method or by a publisher with market preferences, but this algorithm provides unconventional categorization with a user's definition. The following is a brief explanation of each step:

a) Analyze as many books as possible to count the summed frequency of each word and to build an "Overall Word-Frequency Dictionary (OWFD)"-We processed four-thousand randomly selected digital books in our college academic library-,

b) Select books that represent a specific genre and make a "Genre Word-Frequency Dictionary (GWFD)",

c) Compare the "Genre Word-Frequency Dictionary (GWFD)" with the Overall Word-Frequency Dictionary (OWFD) and find extraordinary words of more frequency in the Genre Word-Frequency Dictionary (GWFD), and with the selected extraordinary words make a Genre-Identity Dictionary (GID).

The Genre-Identity Dictionary will be used for a similarity indicator on each genre, including a 'user-defined unconventional custom genre'. What follows is a detailed explanation of this process.

3.1 Overall Word-Frequency Dictionary (OWFD)

Table 3. Overall Word Frequency Dictionary(OWFD)

| Word | Frequency | Rank |
|------------|------------|--------|
| the | 11,583,621 | 1 |
| and | 5,586,383 | 2 |
| to | 5,320,279 | 3 |
| of | 4,923,275 | 4 |
| a | 4,669,669 | 5 |
| | | |
| zombielike | 18 | 98,770 |
| zulfi | 18 | 98,771 |
| zw | 18 | 98,772 |

We collected four-thousand digital books to make one OWFD (Overall Word-Frequency Dictionary). The dictionary contains frequency and rank of each word in all of the books in our sample. More than half a million words are in this dictionary, but the frequency distribution is not linear, so the 20% of high-ranked words covers more than 98% of the whole word-frequency. The lowest-ranked word group, which makes the database enormously inflated is filled with invented onomatopoeia, mimetic words, proper nouns (e.g. character names), and misspelled words. To reduce processing cost we selected the top 20% words (10,000 words) to build the final OWFD (Table 3).

3.2 Genre Word-Frequency Dictionary (GWFD)

We made a GWFD (Genre Word-Frequency Dictionary) the same way we made OWFD, but we only sampled books that are representative of each genre not whole collection. We chose some of conventional genres, which are fantasy, philosophy, and science fiction, and we also created a userdefined custom genre. To represent each genre, it is important to select the most suitable representatives of each genre; we chose several important books for each (Table 4).

Table 4. Selected genres and representative books

| Genre | Representing books | |
|-----------------|---|--|
| | - The Lord of the Rings, by J. R. R. Tolkien | |
| | - The Lion, the Witch and the Wardrobe, by | |
| Fantasy | C. S. Lewis | |
| | - Harry Potter and the Sorcerer's Stone, by | |
| | J. K. Rowling | |
| | - The Apology, Crito and Phaedo of | |
| | Socrates by Plato | |
| Philosophy | - Critique of Pure Reason by Immanuel | |
| | Kant | |
| | The Ethics by Benedict de Spinoza | |
| | - 2001: A Space Odyssey by Arthur C. | |
| | Clarke | |
| | - Childhood's End by Arthur C. Clarke | |
| Science Fiction | - Double Star by Robert A. Heinlein | |
| | - Starship Troopers by Robert A. Heinlein | |
| | - The Currents of Space by Isaac Asimov | |
| | - The Naked Sun by Isaac Asimov | |
| | - Jane Eyre by Charlotte Bronte | |
| Custom Genre | - Sense and Sensibility by Jane Austen | |
| | - My Sister's Keeper by Jodi Picoult | |

For the fantasy genre, we selected two books from two classic masters, Tolkien and Lewis, in addition to the first volume of the Harry Potter books, which is a best seller in the contemporary fantasy novel genre. For the philosophy genre, we selected three books: by Plato, Kant and Spinoza; and for the science fiction genre, we chose six books by the so called big three writers (Clarke, Heinlein, Asimov). The custom genre is a user-defined genre based on one user's preference, who is one of co-authors of this paper. We named the genre 'Kimyo's selection' after her nickname. She loved Jane Eyre, Sense and Sensibility, and My Sister's Keeper. These books show women writers' style, and they also speak out about social issues such as social class problems, child abuse, women's rights, and medical self-determination. These books are romantic and philosophical novels that cannot be described as a single conventional genre.

| Rank | Fa | ntasy | Phil | Philosophy | | Fiction | | Selection | |
|------|-----|--------|------|------------|-----|---------|-----|-----------|--|
| 1 | the | 17,529 | the | 23,407 | the | 21,024 | the | 14,587 | |
| 2 | and | 11,165 | of | 18,009 | to | 10,193 | to | 10,886 | |
| 3 | of | 7,202 | to | 10,544 | of | 10,032 | Ι | 10,791 | |
| 4 | to | 6,673 | and | 8,859 | and | 9,822 | and | 10,332 | |
| 5 | a | 6,243 | in | 8,842 | a | 9,472 | of | 8,486 | |
| 6 | he | 5,078 | is | 8,097 | Ι | 7,102 | a | 7,915 | |
| | | | | | | | | | |

 Table 5. Genre Word Frequency Dictionary(GWFD)

As we can see in Table 5, the words ranked by frequency are similar between the genre word databases because the basic and required words to build a sentence are same in each. Therefore it is not easy to find a difference with these databases.

3.3 Genre Identity Dictionary(GID)

The GWFD (Genre Word-Frequency Dictionary) has to be refined to represent a genre identity. We suggest an algorithm that compares between an OWFD (Overall Word-Frequency Dictionary) and a GWFD (Genre Word-Frequency Dictionary to extract extraordinary words for each genre.

We measured a rank difference of a word that occurs in both between two dictionaries. If a word had a higher rank in the GWFD than OWFD, we recorded the word into a database with a rank distance.



Fig. 4. The higher ranked words in the GWFD (Genre Word-Frequency Dictionary of fantasy genre); they form the Genre-Identity Dictionary

As Fig. 4 shows, there are words that are frequently found in a certain genres (e.g. wizard and wand in a fantasy genre). We

call these words "genre-identity words" and the database a "Genre-Identity Dictionary (GID)".

Table 6 shows genre-identity words in the GIDs (Genre-Identity Dictionaries) sorted by rank distance. Each genre has a special set of extraordinary words. Some of the words in the highest ranked group are proper nouns such as characters' names -except philosophy genre- but very soon we can find the unique genre-oriented words; wizard, journey (fantasy); bugs, discovery (science fiction); behavior, obliged, engagement (Kimyo's selection).

| Fable 6 | Genre | Identity | Dictionaries | GIDs |) |
|---------|-------|----------|--------------|------|---|
|---------|-------|----------|--------------|------|---|

| Rank | Fantasy | Philosophy | Science Fiction | Kimyo's Selection |
|------|---------|-------------|--------------------|----------------------|
| 1 | Hany | therefore | Baley | Kate |
| 2 | Frodo | conception | Daneel | Elinor |
| 3 | Gandalf | object | Rik | Anna |
| : | | | | |
| 41 | wizard | necessarily | Trantor | behavior |
| 42 | journey | regard | bugs | obliged |
| 43 | Legolas | proposition | discovery | engagement |
| | | | | |

3.4 Genre closeness calculation

We can calculate genre closeness of a book with GIDs (Genre-Identity Dictionaries). First, we made a frequencyranked word dictionary of the subject book that a customer wants to test, and we compared it against various GIDs. If a word in the dictionary of subject book is also found in a GID, we call it a "word hit". We also calculated the average rank distance of each paired word, and we named it "Word Average Distance". Naturally, more word hit and less word average distance implies that the subject book is closer to the genre.

Table 7 shows the results of the genre closeness between seven exemplary books and four predefined genres, including a custom genre 'Kimyo's selection'. Some of the subject books are not easy to describe as a single conventional genre. Alice's Adventures in Wonderland (by Lewis Carroll) is a children's book, but it also has profound symbolic messages. Solaris by (Stanislaw Lem) appears as clearly science fiction genre, but the book's message is philosophic; man's anthropomorphic limitations. In Jodi Picoult's first bestseller, Nineteen Minutes, she mixed family, morality, and many social controversies into a complex and twisted plot. We also tested Harry Potter and the Prisoner of Azkaban (by J.K. Rowling), The Analysis of Mind(by Bertrand Russell), An Inquiry into the Nature (by Adam Smith), and Little Women (by Louisa May Alcott).

| Book | Word Hit | | | Word Average Distance | | | | |
|-------------------------------------|----------|------------|------|-----------------------|----------|------------|----------|---------|
| Doon | Fantasy | Philosophy | SF | Kimyo's | Fantasy | Philosophy | SF | Kimyo's |
| Alice's Adventures in Wonderland | 274** | 188 | 108 | 229* | 116.03* | 116.34* | 118.13 | 117.45 |
| Solaris | 127 | 374** | 307* | 243 | 119.55 | 119.14* | 118.05** | 119.51 |
| Harry Potter3 rd Book | 350** | 128 | 120 | 208* | 115.06** | 121.31 | 119.28 | 120.74 |
| Nineteen Minutes | 86 | 145 | 115 | 270** | 119.41 | 121.50 | 120.17 | 118.46* |
| The Analysis of Mind | 82 | 880** | 355 | 467 | 116.81 | 90.79** | 114.30 | 113.36 |
| Inquiry into the Nature | 125 | 687** | 288 | 600* | 116.41 | 103.88** | 114.06 | 110.85 |
| Little Women | 190 | 267 | 153 | 525** | 119.07 | 118.66* | 119.74 | 118.54* |

Table 7. Genre Closeness for Subject Books

The genre closeness of selected books is not different from our expectation. A primary genre, the closest genre of each subject book, fits into the conventional genre classification of Internet bookstores, but what we have to pay attention to is the multi-genre closeness of mixed-genre books. Solaris shows closeness in two genres; science fiction and philosophy, as we suspected earlier. Alice's Adventures in Wonderland shows complex characteristic of fantasy (word hit), philosophy (word average distance), and the custom genre (word hit). Although The Analysis of Mind and An inquiry into the Nature both are conventionally classified as philosophical publications, only An Inquiry into the Nature shows closeness in the custom genre Kimyo's selection, which is related to social issues.

3.5 Genre closeness visualization

The result of genre closeness is accurate and clear, but showing many numbers, as seen Table 7, is not perceivable or practical. Moreover by increasing the number of subject books or predefined genres, the difficulty of similarity detection will be intensified. To solve this problem we suggest data visualization, which helps to find patterns intuitively based on genre closeness (word hit and word closeness), and we call the image "Figure of Genre Closeness". (Table 8, Fig. 5)



| place a 'word-frequency dictionary' of subject book left side | | | | | |
|---|--|--|--|--|--|
| (sorted with frequency) | | | | | |
| place a GID of a genre right side (sorted with rank distance) | | | | | |
| while(not end of 'word-frequency dictionary' of subject book) { | | | | | |
| get a next word from 'word-frequency dictionary' of | | | | | |
| subject book (left side) | | | | | |
| find the word from opposite GID (right side) | | | | | |
| calculate the distance between two words | | | | | |
| set the thickness of a line and transparency based on the | | | | | |
| distance (closer is thicker and brighter) | | | | | |
| draw a line between two point based on the thickness | | | | | |
| and transparency | | | | | |
| } | | | | | |



Fig. 5. Figure of genre closeness with a subject book, Harry Potter 3rd book (Harry Potter and the Prisoner of Azkaban) and a fantasy-genre-identity dictionary by connecting two same words in both dictionaries; the line thickness and transparency (alpha channel) are proportional to distance of each line.



Fig. 6. Figure of genre closeness

Fig. 6 shows results of this visualization. The most complex and bright image implies the closest genre of the subject book, so users are able to instantly describe the genre closeness; The Analysis of Mind as philosophy, Harry Potter and the Prisoner of Azkaban as fantasy, and Little Women as a custom genre. Moreover, this visualization presents mixed multi-genre closeness as well as the primary genre.

4 Combined Visualization of Readability and Genre

Previously we suggested readability visualization (Algorithm 1) and genre visualization (Algorithm 2). In this chapter, we will propose a practical visualization that combines the information of 'Readability Visualization' chapter and 'Genre Visualization' chapter, as each has own strengths. The suggested steps are as follows:

a) Define genre-hue based on color symbols,

b) Set the genre-hue as base hue in the readability visualization with hue rotation.

There is no right answer for the color-genre connection. We defined genre-hue based on the Table 9 below; based on symbol of color. This definition is temporary, to show the visual clearness of the algorithm.

| Table 9. Defined Genre Hue(bold words are closely related to |
|--|
| each genre) |

| Genre | Color | Representing books |
|----------------------|--------|---|
| Fantasy | Red | challenge, passionate, rage, active, exciting, dangerous, courage, happiness, love, loyalty, delight, power, sin, passion, vivid |
| Philosophy | Blue | dignity, sorrow, truth, harmony, happiness, sincerity, holy, intelligence, spirit, cold, stability, calm, hope, cold- hearted, melancholy, intelligent |
| Science Fiction | Purple | victory, modesty, tragedy, mystery, sacred, Exotic, splendid, sorrow, dignity, complicated, noble |
| Kimyo's Selection | Green | relaxation, happiness, beauty, hope, nature, safe, peace, calmness, Mild, longing.cozy, pleasant, abundant |

Table 10. Algorithm 3: Readability and Genre Visualization

```
Max)
while(not end of the book) {
   get a next paragraph
   set a saturation based on paragraph length
   while(not end of the paragraph) {
       get a next character
       if (the character is punctuation or comma)
           increase brightness (sentence formula)
       else {
           decrease brightness (paragraph formula)
       if(three or more syllables words) // optional
           decrease brightness (word formula)
       rotate hue // slow to fast
       set a pixel on next location with assigned hue,
       saturation, brightness and draw
   reset brightness and saturation to Max, hue to original
genre-hue
Adjust saturation for the clearness
```

Table 10 (Algorithm 3) is mostly the same as the previously suggested readability algorithm (Algorithm 1) except for the addition of a base genre-hue and a hue rotation. The base genre-hue sets a color tone of book visualization, and the hue rotation makes rhythmical color patterns, which emphasize book readability and genre with juxtaposition of complimentary color shadow. Additionally, saturation is decided by the length of paragraph; it also shows paragraph readability. The results of Algorithm 3 are shown in Fig. 7.



Fig. 7. Left images show readability visualization only (Algorithm 1); right images are visualization with the new suggested algorithm (Algorithm 3); a lengthy paragraph colored with de-saturated tone due to the saturation setup function in the algorithm. As a result, most of philosophy publications are not vivid as novels.

5 Conclusion and Future Work

This alternative method of book selection has the strength of conventional data visualization that finds patterns instantly over huge data with human perception. We tested our visualization and the results prove the effectiveness of the suggested algorithm. All testers (n=20) easily categorized the same genre based on hue. The subjects were also asked to evaluate readability visualization on 5-point scale (1=very easy; 5=very hard) of 10 books, and the result shows that most of them predicted readability correctly (r = 0.79, p<0.01), so we believe this research will help readers to choose a right book.

With customers' point of view this visualization could be very useful in the current online bookstore interface by showing essential information of the custom favorite genre and readability as the figure shows Fig. 8.



Fig. 8. Suggested example of the text visualization (Amazon.com format)

We only suggested one example of possible combined visualizations in this paper, but with readability and genre data there will be more visualization methods and interfaces, especially to support custom-genre applications. Those series of research will be our next task.

6 Acknowledgement

This work was supported by the National Research Foundation of Korea Grant funded by the Korean Government MEST, Basic Research Promotion Fund (NRF-2010-013-H00001).

7 References

[1] J. W. Park, G. Choe, "Visual Genealogy", Proceedings of ACM SIGGRAPH '09(New Orleans, Louisiana) Art Gallery, ACM, 2009.

[2] J. W. Park, "Information Aesthetics with Visual Genealogy Project", Leonardo, Vol. 44, No. 5, pp.464-465, 2011.

[3] W. B. Paley, "TextArc: Alice's Adventures in Wonderland", http://www.textarc.org/Alice.html, 2009.

[4] S. Posvec, G. McInerny, "The evolution of the origin of species",

http://www.visualcomplexity.com/vc/project.cfm?id=69, 2009.

[5] P. Steinweber, A. Koller, " Similar Diversity (Holly Bible)", http://similardiversity.net/, 2007.

[6] J. R. Bormuth, "Cloze Test Readability: Criterion Reference Scores," Journal of Education Measurement, Vol. 5, No. 3, pp. 189-196, 1968.

[7] E. Dale, J. S. Chall, " A Formula for Predicting Readability ", Educational Research Bulletin, Vol. 27, No. 1, pp. 11-20, 1948.

[8] R. Flesch, "A New Readabiltiy Yardstick", Journal of Applied Psychology, Vol. 32, No. 3, pp. 221-233, 1948.

[9] W. S. Gray, B. E. Leary, "What makes a book readable", The University of Chicago Press, 1935.

[10] R. Gunning, "The technique of clear writing", NY: McGraw-Hill International Book Co., 1952.

[11] A. J. Harris, M. D. Jacobson, "The Harris-Jacobson readability formulas", New York: David McKay, 1975.

[12] H. D. Kitson, " The mind of the buyer: A psychology of selling", Macmillan, 1921.

[13] B. A. Lively, S. L. Pressey, " A Method for Measuring the 'Vocabulary Burden' of Textbooks", Educational Administration and Supervision, No. 9, pp. 389-398, 1923.

[14] M. Vogel, C. Washburne, " An Objective Method of Determining Grade Placement of Children's Reading Material", The Elementary School Journal, Vol. 28, No. 5, pp. 373-381, 1928.

[15] L. A. Sherman, "Some observations upon the sentencelength in English prose", University studies (University of Nebraska (Lincoln campus)), Vol. 1, No. 2, pp. 119-130, 1888.

[16] G. W. McConkie, K. Rayner, " The span of the effective stimulus during fixations in reading", Attention, Perception, & Psychophysics, Vol. 17, No. 6, pp. 578-586, 1975.

[17] N. E. Miller, P. C. Wong, M. Brewster, H. Foote, "TOPIC ISLANDS—a wavelet-based text visualization system", Proceedings of the conference on Visualization '98, pp.189-196, 1998.

[18] S. Huang, M. O. Ward, E. A. Rundensteiner, " Exploration of Dimensionality Reduction for Text Visualization", Coordinated and Multiple Views in Exploratory Visualization 2005 (CMV 2005), pp. 63-74, 2005.

[19] W. Weber, " Text Visualization – What Colors Tell About a Text", Information Visualization 2007 (IV '07), pp. 354 - 362, 2007.

[20] N. Fang , X. Luo, W. Xu, "Single Text Combinatorial Semantic Field (STCSF): Construction and Visualization", The Third International Conference on Semantics, Knowledge and Grid 2007, pp. 188-193, 2007.

A GENERAL TAXONOMY FOR VISUALIZATION OF PREDICTIVE SOCIAL MEDIA ANALYTICS

Stacey Franklin Jones, D.Sc. ProTech Global Solutions Annapolis, MD

Abstract

The use of Social Media as a resource to characterize historical and current occurrences to reasonably predict future events is rapidly emerging. Effective visualization of Predictive Social Media Analytics (PSMA) data lends to meaningful and relevant translation, efficient interpretation and ease in use of results. A general taxonomy for the measurement of these attributes enable comparison of visualization techniques with the ultimate goal of identifying some of the "best" PSMA data representation approaches and export features for a broad range of applications and respective complexity.

Background

The quest to identify trends in large publicly available data realms such as that associated with Social Media and using them to foresee what may happen in the time to come requires robust analytic approaches. We are seeing the emergence of various applications that make use of algorithmic techniques to provide this predictive capability based on extraction and manipulation of huge volumes of data available via websites and other online means of communication used by people to share "social" information. These Predictive Social Media Analytics (PSMA) applications offer various features at different complexity levels to aid in making more informed decisions about marketing products and improving financial services. PSMA applications are also beginning to infiltrate the defense and intelligence space to assist in monitoring and forecasting world events. Their common goal is to generate output that facilitates understanding of complex data relationships over time. Since PSMA application operate on such massive amounts of data, visualization and other data representation of temporally significant instances with a reasonable degree of graphical clarity becomes no easy feat. Ultimately their value will be measured not only by the speed, sophistication and accuracy of their algorithms, but by how well they produce or provide interfaces that facilitate meaningful and relevant translations that lend to efficient interpretation and ease in use of results.

Amongst the numerous PSMA applications are three (3) that cover the general forecasting space of consumer product-focused, financial climate and defense intelligence forecasting. Educational use of PSMA is of particular interest to the author(s) but at this time is more a residual space related to the aforementioned. As such, for the purpose of a general survey of relevant PSMA visualization techniques the following applications are reviewed: IBM's SPSS predictive analytics software, SAS Social Media Analytics and Recorded Future. Admittedly so, the "predictive" power of the named PSMA applications do not claim to be equal. However, they generally incorporate output data representation patterns sufficient to develop a 'general' taxonomy for the measurement of key attributes.

IBM's SPSS predictive analytics software developers state that they can predict with confidence what will happen next so that you can make smarter decisions, solve problems and improve outcomes. They invite the user to create and share what is referred to as "compelling visualizations that better communicate [your] analytics results" [1]. IBM SPSS also states that with the use of their IBM SPSS Visualization Designer the user may "easily develop and build new visualizations that enable new ways to portray and communicate analytics" without extensive programming skills. Three (3) example visualizations available to the IBM SPSS user are: Network Graph, Scatterplot Matrix and Jungle Book Graph shown in Figures 1-3 [Reprint Courtesy of International Business Machines Corporation, © SPSS, Inc., an IBM Company].



Figure 3 Jungle Book Graph

SAS Social Media Analytics makes the claim that they provide the "power to know what lies ahead around the next corner" through "an integrated environment for predictive and descriptive modeling, data mining, text analytics, forecasting, optimization, simulation, experimental design and more. From dynamic visualization to predictive modeling, model deployment and process optimization, SAS provides a range of techniques and processes for the collection, classification, analysis and interpretation of data to reveal patterns, anomalies, key variables and relationships, leading ultimately to new insights and better answers faster." [2] Associated with the SAS analytics tool is an ability to illustrate trends and enable drilling down to a level of actual comments that contribute to the prediction(s). To illustrate these changes over time, they offer amongst other options, a dashboard of standard column, bar, line and pie charts and tables as a method for visualizing correlations of social media trends with the circumstances that triggered those events. An example of this approach which can incorporate marketing activities, product changes, world events and/or market conditions is shown in Figures 4 and 5.



Figure 4 SAS® SMA Superimposition A

Figure 5 SAS® SMA Superimposition B

Also offered via SAS® Visual Data Discovery is an interactive data visualization for analytics suite that includes, but is not limited to: scatter, bubble and 3-D contour plots with animation.

Recorded Future "strives to provide tools which assist in identifying and understanding historical developments, and which can also help formulate hypotheses about and give clues to likely future events." Recorded Future data may be accessed through a web services API using two export formats json or csv text. Developers are then referred to statistical or visualization software packages such as R, Spotfire or others that can make use of the export formats. A review of R presents a number of output formats such as box and whisper charts, pie charts, pairs plots, coplots, forest plots and common 3-D plots. Spotfire offers an "analytics' visual and exploratory experience" which includes "elegant and configurable visuals: Bar chart, Map chart, Line chart, Pie chart, Scatter plot, Combination Bar and Line chart, Cross Table, 3D Scatter Plot, Treemap, Heat Map, and Parallel Coordinates plot" capability in customizable dashboard configurations. In general, Recorded Future's analytics PSMA visualization capability is as robust as the application programming interface and the linked packages.

One could argue that the keys to the successful emergence of PSMA in terms of timely and informative use will be meaningful and relevant visualizations. However, what is common amongst the reviewed PSMA application visualization approaches is the utilization of 'old' and primarily two (2) dimensional data representations for a 'new' multi-dimensional phenomenon of analytics. Furthermore the temporal element of prediction introduces a complex dimension that maybe expanded and linked to places, event, entities, location, sentiment, behaviors and other such elements that lend to forecasting future events.

Why A Taxonomy?

As visualization techniques attempt to 'catch up' with the advancement of PSMA, we can expect a commensurate desire for comparison and subsequently a proliferation of best representations. A taxonomy that facilitates such a comparison is needed. The taxonomy need not be complex at this juncture of evolution as PSMA applications are bound to morph many times over with the improvement of the underlying search and analyses engines. However, there is a need to categorically address the different levels of analysis of the PSMA application output data. Three (3) general PSMA application visualization categories can be established based on the reviewing analysts' expertise and familiarity with data representations:

- 1. High expert data mining for intelligence, defense and other such datasets with high levels of complexity and many elements
- 2. Moderate strategic/competitive product introduction, political predictions, and other such datasets with medium levels of complexity and several elements
- 3. Low general consumer type for the common data observer which may be used in the educational environment and to communicate with the general public using a few elements

A General Taxonomy

A general taxonomy that characterizes the effectiveness of visualization of PSMA data would need to cover three (3) areas:

- 1. Meaningful and relevant data translation
- 2. Efficient interpretation
- 3. Ease in use of results

Effectiveness measurements would be taken across each of the three (3) PSMA application categories.

Meaningful and Relevant Data Translation

PSMA applications search and operate on large volumes of content from government sites, news sites, blogs posts, tweets and other such information available on the web. Predictive results of these massive processing tasks have the challenge of presentation in a manner that is informative and would ultimately lead to better decisions, problem solving and improved outcomes. As such, the measurement of how meaningful and relevant the results translation is paramount and is the first component of a general taxonomy in all three (3) of the PSMA application visualization categories.

Efficient Interpretation

The need to interpret predictive results timely is a component of each of the three (3) PSMA application visualization categories, but varies in sense of urgency. For the High category the ability to reasonably forecast a potentially threatening world event and mobilize any intervention would have a desirable range of minutes or hours. Whereas for the Medium or Low categories, a several days or perhaps even weeks is most sufficient for forecasting and preparing for a new product release or education focused event. Differing time requirements for interpretation are expected but are key elements for all three (3) categories. Consequently, the measurement of how efficiently the resulting predictive information can be interpreted is an additional component of the taxonomy.

Ease in Use of Results

All of the PSMA applications categories would benefit from data representations that lend to ease of use of the information intensive results. Two (2) dimensional charts and grams are limited in domain and range of visualization and are at least one generation behind the third generation search engine characterization of PSMA applications. Therefore, measurement of the extent to which the application offers results that are easily transferrable to next generation representations, perhaps that more animated and/or incorporate three (3) or better dimensions is desirable.

The following is a simplified General Taxonomy for Visualization of Predictive Social Media Analytics Grid Figure:

| | Meaningful/ | | |
|--------|-------------|----------------|----------------|
| | Relevant | Efficient | Ease of Use of |
| | Data | Interpretation | Results |
| | Very | Very | Very |
| High | Important | Important | Important |
| | Very | | Very |
| Medium | Important | Important | Important |
| | Very | Less | Very |
| Low | Important | Important | Important |

Next Steps

Validation of the general taxonomy or assessment of the PSMA applications visualization capabilities will need to be conducted. The metrics for the validation may have both quantitative and qualitative components and go across each of the three (3) visualization categories (i.e. High, Medium, Low). This validation can be initially conducted on the three (3) PSMA applications reviewed as they appear to represent the more evolved in the market. However, next steps should include a full survey of PSMA applications as they are developed and released as products. A call to and invitation to other investigators to address the first two (2) PSMA application visualization categories – specifically, the High and Medium – from a taxonomy perspective is made. The authors' immediate investigative interests will be primarily in the third category and more specifically as it has to do with predictive behaviors that affect the teaching/learning/scholarship environment of college students.

Conclusions

In conclusion, while the scope of this investigation and development of a General Taxonomy for Visualization of Predictive Social Media Analytics (PSMA) is not exhaustive in its review of application, it is representative of the evolving applications in the market today. It provides a framework for categorizing requirements based on type of use and information, and a respective set of general measurement – both quantitative and qualitative - classes for comparison and extracting best visualizations. The next steps would include validation of the general taxonomy and more broad and detailed investigation of a full survey of PSMA applications.

References

[1] SPSS Visualization Designer Charting your course just got easier http://www-01.ibm.com/software/analytics/spss/products/statistics/vizdesigner/ (last referenced May 15, 2012)

[2] SAS® Analytics Analytics delivering greater insight http://www.sas.com/technologies/analytics/ (last referenced May 15, 2012)

[3] Bollen, J., Mao, H., & Pepe, A. (2011). Modeling Public Mood and Emotion: Twitter Sentiment and Socio-Economic Phenome-na. In Proc. ICWSM 2011.

[4] Bollen, J., Mao, H., & Zeng, X. (2011). Twitter Mood Predicts the Stock Market. Journal of Computational Science www.elsevier.com/locate/jocs

[5] De Choudhury, M., and Sundaram, H. (2011). Why Do We Converse on Social Media? An Analysis of Intrinsic and Extrinsic Network Factors. In Proceedings of the Third SIGMM Workshop on Social Media (WSM 2011), in conjunction with ACM Multimedia 2011 (Scottsdale, Arizona, USA, November 28 - December 1, 2011)

[6] De Choudhury, M., Lin, Y.-R., Sundaram, H., Candan, K.S., Xie, L. & Kelliher, A. How Does the Data Sampling Strategy Impact the Discovery of Information Diffusion in Social Media? ICWSM '10: Proceedings of the 4th International AAAI Conference on Weblogs and Social Media. 2010

[7] Gruhl, D., Guha, R., Kumar, R., Novak, J., Tomkins, A. The Predictive Power of Online Chatter http://www.tomkinshome.com/site_media/papers/papers/GGK+05.pdf

[8] Murphy, T., Data Mining: Using Predictive Analysis And Social Network Analysis 02/09/2011 *New Tech Post* http://newtechpost.com/node/271

[9] Schmidt, J. (2007). Blogging practices: An analytical framework. Journal of Computer-Mediated Communication, 12(4), article 13. http://jcmc.indiana.edu/vol12/issue4/schmidt.html

[10] Truvé, S., Recorded Future : A White Paper on Temporal Analytics https://www.recordedfuture.com/assets/RF-White-Paper.pdf (last referenced May 15, 2012)

[11] Social Media As Predictive Analytics?

http://www.marketingvox.com/social-media-as-predictive-analytics-048020/ (Nov 3, 2010) (last referenced May 15, 2012)

[12] Spotfire About Spotfire <u>http://spotfire.tibco.com [On Line Demo(s)]</u> (last referenced <u>May 15, 2012</u>)

Figures Reprint Courtesy of International Business Machines Corporation, © SPSS, Inc., an IBM Company. SPSS was acquired by IBM in October, 2009; and permission of SAS Social Media Analytics

Visualization Tool for Apollonian Network and Packing Analysis

Lucas Vespa Department of Computer Science University of Illinois at Springfield Springfield, IL 62703

Abstract—Apollonian Networks, inspired by Apollonian Packings, is a new and fascinating area of research that has proven fruitful for modeling dynamic processes and realworld networks. Exploring these packings and networks further may provide solutions for a myriad of important problems. Given the brilliant physical properties of Apol-Ionian Packings and Apollonian Networks, it seems only natural to have a visualization tool specific to these structures. This work presents VAPNA, a Visual Apollonian Packing and Network Analysis software tool for exploring Apollonian Packings, Networks and their relationships. This tool currently generates two-dimensional networks within a packing and allows for analysis of paths and other aspects of the networks. VAPNA is designed to be an aid in furthering research in the area of Apollonianbased modeling and discovery.

I. INTRODUCTION

Apollonian Networks [1], [2], [3], [4], inspired by Apollonian Packings [5], [6], [7], are proving to be an interesting tool for modeling real networked systems. Apollonian Packings themselves have interesting recursive mathematical relationships, while Apollonian Networks have other interesting properties such as being euclidian and scale-free. Apollonian networks closely describe the behavior of granular packings and porous media, road systems, electrical supply networks and even brain behavior.

Apollonian packings are a fractal design of embedded tangent circles as shown in Figure 1(a). Apollonian networks can be represented by a graph as shown in Figure 1(b), with vertices formed by circles from a packing and edges formed by connecting tangent circles. Given the potential usefulness of Apollonian structures, it seems logical to have a visual tool that mimics the pen and paper process of discovery in Apollonian Packings and Networks.



Fig. 1. Apollonian Packing and Network generated by VAPNA.



Fig. 2. Deriving mutually tangent circles in an Apollonian Packing using the Descartes Circle Theorem. The curvatures (inverse of the radii) of any four mutually tangent circles (Descartes Configuration), satisfy the Descartes Circle Theorem shown in Equation 1.

In this work we present VAPNA, a Visual Apollonian Packing and Network Analysis software tool. Specifically, VAPNA allows simultaneous visual exploration of both packings and networks such that, in addition to analyzing packings and networks individually, their relationships may be examined as well. We believe this tool will be useful to those who examine and analyze Apollonian Network and Packing properties.

VAPNA allows a user to create packings, and networks within those packings, specifying parameters for how the networks are to be formed and what aspects of the structure should be visualized. Different aspects can be analyzed such as paths.

The remainder of the paper is organized as follows. Section II discusses how to build an Apollonian Packing and Network. Section III discusses the VAPNA software, including commands and examples. Section IV discusses related work and Section V concludes the paper.

II. APOLLONIAN PACKINGS AND NETWORKS

Constructing an Apollonian Packing begins with circles in a Descartes Configuration [8], which consists of four mutually tangent circles. Three mutually tangent circles with known radii are used to find a fourth mutually tangent circle. The four circles have radii (r1...r4). The curvatures (c1...c4), are the inverse of the radii and satisfy the following, according to the Descartes Circle Theorem:

$$\sum_{i=1}^{4} c_i^{\ 2} = \frac{1}{2} \cdot (\sum_{i=1}^{4} c_i)^2 \tag{1}$$



Fig. 3. VAPNA software screen shot.

An example of finding the curvatures of circles in a packing is shown in Figure 2. In this figure there is an outer circle with radius 1, two inner circles with radius $\frac{1}{2}$ and a fourth mutually tangent circle with radius $\frac{1}{a}$. The curvatures of these circles are (-1, 2, 2, a), and they satisfy the following:

$$-1^{2} + 2^{2} + 2^{2} + a^{2} = \frac{1}{2} \cdot (-1 + 2 + 2 + a)^{2}$$

Using this relationship, the value of a can be derived. Subsequently, the three circles with curvatures (-1, 2, a) can be used to find curvature b by the following:

$$-1^{2} + 2^{2} + a^{2} + b^{2} = \frac{1}{2} \cdot (-1 + 2 + a + b)^{2}$$

To create a packing, the process of finding mutually tangent circles continues in a space filling manner. Apollonian Networks can be created by connecting tangent circles to form a graph. However, in VAPNA, certain commands can be used to create networks with non-tangent circles as well.

III. VAPNA

Figure 3 shows the VAPNA software. Although most commands can be accessed from the menus, the commands should be executed by text input. VAPNA can visualize the full or partial packing, which represents the vertices of any network or path. It can also visualize the edges of a path or network. The following sections describe the VAPNA command set with examples.



Fig. 5. path command examples



Fig. 4. Core commands (net and path), along with associated modifiers.

A. Core Commands

VAPNA has two core commands, path and net. The commands use a curvature sequence (either static or iterative), to dictate construction of a graph, along with optional modifiers. The following is the basic usage of the core commands:

core_command(curvature_sequence,_ optional iterations) optional modifier_list

The core_command is either path or net. The curvature_sequence is used to designate the circles in the packing that represent the vertices in the Apollonian graph. The curvature_sequence can be static (e.g., 2;3;6;11;...) or generated iteratively by an equation (e.g., x + 2, generates 1;3;5;7;...). Several formats can be used for equation input at this time. On the case that an equation is used for the curvature_sequence, the number of iterations must be specified. Currently, infinite iterations are not supported.

The path and net commands have optional modifiers specified by the modifier_list. which adjust the core commands, giving the user more control. Figure 4 shows the modifiers and their possible usage in relation to the core commands.

B. Path and Modifiers

Unlike in a network, path does not connect all tangent circles in a curvature sequence. The path command connects tangent circles only in the sequence given. If there is a break in tangency between two curvatures in sequence, the path is broken and the tangency break is reported in the form of the two non-tangent curvatures.

path: The unmodified path command displays the lines (edges) of a path. Figure 5(a) shows an example for the following command: path(2;15;3;23;6;11;18;27).

path-lin and path-cir: The lin and cir modifiers for the path command display only the lines (edges) or circles (vertices) of a path and can be used together. Figure 5(b) shows an example of the following command: path(2;15;3;23;6;11;18;27) cir. Figure 5(c) shows and example of the following command: path(2;15;3;23;6;11;18;27) lin cir.

path-tri: The tri modifier completes a right triangle between all tangent circles in the path. Figure 5(d) shows and example of the following command: path(2;3;6) tri.

C. Net and Modifiers

The net command is used to create networks of tangent circles, and can also be extended to nontangent circles using certain modifiers.

net: The unmodified net command creates an Apollonian Network using tangent circles in the curvature sequence. Figure 6(a) shows an example of the following command: net(2;3;6;11;15;18;23).

net-seq: The seq modifier has similar functionality as the path command but, when used in



Fig. 6. net command examples

conjunction with net, extends to non-tangent circles. Figure 6(b) shows and example of the following command: net(2;3;11;27) seq.

net-mesh: The mesh modifier extends the net command to non-tangent circles. Figure 6(c) shows and example of the following command: net(all, 35) mesh. This command is a special case using a curvature sequence of all and a second argument of max_curvature. This special case generates a curvature sequence of all curvatures up to a maximum curvature value. The net-mesh command does not require this special case, it is just used to show and example of this special case curvature sequence.

net-seq-tri and net-mesh-tri: The tri modifier is used to complete right triangles between connected circles. The tri modifier can be used to extend both the net-seq and net-mesh commands.

net-union and net-intersect: The union and intersect modifiers can be used with the net command to create a network which combines a stored set of curvature sequences in either a set union or set intersection. Curvature sequences are stored in a history using the store command.

D. Other Commands

Other commands include methods to store and clear a history of curvature sequences, clear the display and change packing and network coloring.

store and history: The store command stores a curvature sequence in memory and is used as follows: *store(curvature_sequence)*. The history command clears all stored curvature sequences.

clear and clear-all: The clear command clears all paths and networks. The clear-all command clears networks, paths and packings.

color-lin and color-cir: The color command can be used to change the color of the networks and paths (lin), or the packings (cir). The command only changes the color for future display and not the current display. This is so that the user can layer network and packing colors.

IV. RELATED WORK

This section discusses work related to Apollonian Networks as well as graph and fractal visualization. To begin, several key works have helped to introduce Apollonian Networks as viable modeling tools.

Andrade et al. [1] and Doye et al. [2] both introduce Apollonian Networks. These works emphasize two dimensional Apollonian Networks and are the inspiration for VAPNA. Specifically, Doye et al. [2] use Apollonian Networks for describing energy landscapes whereas Andrade et al. [1] are concerned with representing dynamic relationships in physical elements. Zhang et al. [3] create an iterative method for generating high-dimensional Apollonian Networks. Emphasis in this work is given to Apollonian Networks of dimension greater than the two.

Apollonian Networks have interesting properties. Therefore the closest related work in visualization would be that which visualizes graphs with properties shared by Apollonian Networks. Along these lines, VanHam et al. [9] created a visualization tool for small-world graphs. Although Apollonian Networks have the property of being small world, this visualization tool is for general small-world graphs. Another similar visualization topic would be fractal visualization. The Fractal Science Kit [10] is a robust software suite for generating complex fractal visualizations. However, this tool is not network related and is not specific to Apollonian Packings.

V. SUMMARY AND FUTURE WORK

Apollonian Packings and Networks are interesting and useful structures that may be used to solve many problems. Having a tool to visually explore these structures which began from visual observations, is imperative for furthering research in this area. In this work we present VAPNA, a Visual Apollonian Packing and Network Analysis software which allows exploration of Apollonian Structures in a new an interesting way. Ideally this tool will help contribute to new discoveries.

In the future we plan to use more advanced network generation algorithms. For example, using a recursive algorithm to generate high-dimensional Apollonian Networks. We also plan to add better visualization support and different forms of analysis.

REFERENCES

- J. Andrade, H. Herrmann, R. F. Andrade, and L. Da Silva, "Apollonian networks: Simultaneously scale-free, small world, euclidean, space filling, and with matching graphs," *Physical Review Letters*, pp. 18702–18702, Jan. 2005.
- [2] J. P. K. Doye and C. P. Massen, "Energy landscapes, scale-free networks and Apollonian packings," Dec. 2006.
- [3] Z. Zhang, F. Comellas, G. Fertin, and L. Rong, "High dimensional apollonian networks," arXiv, Mar. 2005.
- [4] Z. Zhang, L. Rong, and S. Zhou, "Evolving apollonian networks with small-world scale-free topologies," *Phys. Rev. E*, Oct 2006.
- [5] I. Peterson, "Circle Game: Packing Circles Within a Circle Turns a Mathematical Surprise," *Science News*, pp. 254–255, Apr.
- [6] F. Varrato and G. Foffi, "Apollonian packings as physical fractals," *Molecular Physics*, pp. 2663–2669, 2011.
- [7] J.-C. Liu, C. Y. Wu, D.-C. Chang, and C. Y. Liu, "Relationship Between Sierpinski Gasket and Apollonian Packing Monopole Antennas," *Electronics Letters*, Jul.
- [8] J. C. Lagarias, C. L. Mallows, and A. R. Wilks, "Beyond the Descartes circle theorem," *Amer. Math. Monthly*, pp. 338–360, 2002.
- [9] Interactive Visualization of Small World Graphs, 2004.
- [10] L. Hilbert, "Fractal science kit fractal generator," http://www.fractalsciencekit.com/, 2011.

59

3-D Visualization of Simulink Physics Models Using Unreal Engine

Elise R. Haley, David J. Coe, and Jeffrey H. Kulick Department of Electrical and Computer Engineering University of Alabama in Huntsville, Huntsville, Alabama USA

Abstract - Model-based design (MBD) tools have evolved that enable simulation of complex cyberphysical systems. Improved visualization techniques are needed to help engineers better understand the results of these simulations. We present a toolkit that facilitates interactive visualization of Simulink models using the advanced graphics capabilities of the Unreal Development Kit. Our toolkit allows Simulink models to create and control actors within a UDK virtual world and allows transformations and events that occur within the UDK to be communicated back to the Simulink models for processing. This two-way communication allows the modeler to distribute the simulation between the two tools to take advantage of each tool's capabilities as needed. For demonstration, we present a Simulink model that controls the orbit and rotation of the solar system within a UDK virtual world.

Keywords: Simulink, Unreal, cyber-physical, 3-D visualization

1 Introduction

The decreasing cost of digital computation has facilitated the integration of programmable computing elements into a wide range of physical systems, and networking technologies have enabled the interconnection of these systems to create ever more complex cyber-physical systems. Engineers of these multi-domain, cyber-physical systems need better tools to help them design, integrate, and verify these increasingly complex systems. Model-based design (MBD) tools provide engineers the opportunity to evaluate the design of complex cyber-physical systems prior to implementation. MBD tools are particularly attractive since they facilitate exploration of alternative designs and removal of defects early in the lifecycle both at relatively low cost.

A challenge for any modeling and simulation tool is making the results accessible and understandable. In this paper we present a toolkit that simplifies the integration of Mathworks' Simulink MBD tool with Epic Game's Unreal Development Kit 3 (UDK) to aid in the visualization of simulation results. We also present a proof-of-concept application developed using our simulation-visualization framework.

2 Background

2.1 What is Simulink?

For many engineering disciplines, Simulink is the de facto standard MBD tool. Simulink is a graphical modeling and simulation tool that allows an engineer to assemble a model of a cyber-physical system by connecting a series of predefined subsystem blocks using virtual wires (signals) [1]. Blocks may represent fundamental model elements such as constants, gain blocks, and queues, or more complicated algorithms and subsystems such as Fast Fourier Transforms or Kalman filters. It is important to note that rather than being interconnected function stubs, the block models themselves are fully executable with actual outputs computed from the applied inputs. Abstraction via nesting of subsystem models into custom blocks may used to simplify modeling of large complex cyber-physical systems. Figure 1 below shows a Simulink model of a standard PID controller assembled from a combination of gain, summation, integration, and differentiation blocks.

The availability of domain-specific sets of modeling blocks (toolboxes) for control system design, digital signal and image processing, and other domains eliminates the need to implement fundamental domain elements and operations making Simulink a particularly attractive option for MBD.



Figure 1 - Simulink model of a Proportional-Integral-Derivative (PID) Controller

For example, rather than assembling a PID controller from basic elements as in Figure 1, a control system designer can just use the predefined PID controller block which includes an auto-tune button to assist in the selection of appropriate gain coefficients. These toolboxes allow a domain expert to quickly model and simulate complex systems. When no appropriate predefined block exists, an engineer may develop a custom block or write a custom function using either C or the MATLAB language and integrate the code with the block model. Simulink also includes the ability to synthesize C-language or Hardware Description Language (HDL) implementations directly from the model, and the ability to conduct hardware-in-the-loop simulations.

2.2 Why interface Simulink with UDK?

MathWorks has developed a 3D Animation toolbox for Simulink which is based on the Virtual Reality Modeling Language (VRML) standard [1], but its feature set and graphics quality are not yet competitive with that of an industry standard game development tool such as UDK which allows developers to create immersive, first person, 3D games with high quality graphics and animation [2]. The UDK's built-in networking support allows multiple-players to interact within the virtual world and also provides a convenient means of interaction with multiple Simulink models.

The versatility and advanced graphics capabilities of the UDK have already been extensively used in a variety of first-person games and nonentertainment applications. Militaries, for example, have explored the use of the UDK-based games for recruiting and training [3] and as platforms for exploring small unit infantry tactics [4]. UDK-based training aids have been developed for nuclear power plant workers [5] and to teach foreign languages and culture-specific gestures and body language [6]. The UDK has also been used to recreate crime scenes as part of a forensic investigation [7].

Notably, Mathis et al. presented a method of interfacing Unreal Tournament (UT) with MATLAB, Simulink, and USARSim to visualize operation of the SAMURAI nano air vehicle [8]. Vehicle dynamics modeled within Simulink were used to control the motion of the vehicle within UT. USARSim was used to gather image data for the user. MATLAB code was used to interface MATLAB and Simulink to USARSIM and UT via TCP/IP communication.

Our toolkit expands on previous work by enabling *two-way interactions* between Simulink models and UDK virtual worlds, which allows the modeler to distribute the simulation across Simulink and the UDK to take advantage of the strengths of each tool. Simulink, for example, would be used to model advanced algorithms, physics, and hardware that govern the behavior of actors (objects) within the virtual world, and the UDK would use its GPU processing capabilities to reposition UDK objects as needed, detect collisions between actors, apply appropriate lighting, etc. With our toolkit, the Simulink model can

• Dynamically create instances of actors from the UDK actor library



Figure 2 – Block diagram of Simulink-UDK toolkit.

- Modify the location, size, roll, pitch, yaw and other properties of the UDK actors
- Move objects via instantaneous translation of actors in the UDK model
- Move objects via interpolated motion of actors in the UDK model
- Respond to UDK events and actor interactions within the UDK that are reported back to Simulink as events, event streams, or actor property updates including
 - single event reports which aggregate multiple events associated with an object into a single, most recent event and
 - a stream of reports that report every interaction of an object with other actors in the scene via a vector of interaction events
- Modify behavior of a UDK actor via user inputs funneled from a Simulink/MATLAB GUI to the Simulink model

By allowing the Simulink model to completely control the building and animation of the UDK model, we hope to minimize the learning curve for the Simulink modeler. Below we provide a brief overview of the toolkit and present a proof-of-concept application.

3 Overview of toolkit

As shown in the block diagram appearing in Figure 2 above, our toolkit makes use of TCP for twoway communication and interaction between Simulink models and the UDK virtual world. To interact with an actor in the UDK, the Simulink model creates a TCP socket connection. Multiple socket connections can be created so that multiple Simulink models (or interactive users if desired) can interact with the UDK virtual world similar to the actions of a multi-person game. One advantage of this method is that computationally intensive Simulink models may be executed on separate hardware from the UDK to speed execution. Moreover, the UDK can perform basic tasks to maintain the virtual world, such as actor collision detection, allowing the Simulink model to focus on its computation until it is notified by the UDK of a collision or other critical event.

On the left side of the block diagram, you will see "Custom TCP Client 1", for example, listed within the Simulink block. The Simulink-side TCP send/receive code has two parts: a standard TCP stack and a message creation part which must specify the actor to affect as well as the operation that the UDK is to perform. Examples of operations include the "spawn" command which directs the UDK to create an instance of an actor predefined within the UDK actor library, the "set location" operation that directs the UDK move an actor instantaneously to a specified location, and the "move" operation that smoothly nudges an actor according to the provided velocity vector.

Similarly, the Unreal TCP communication code is also generic code that is customized for interfacing with specific classes of actors. The TCP server actor spawns one child connection handler per connection, to handle the actual send/receive operations. The connection handler can access any actor class you declare an instance for in the child's code. To spawn or manipulate a specific UDK actor object, you must instantiate that actor within its associated connection handler. At the start of the UDK/simulation run, all



Figure 3 – Simulink orbit model of solar system with details of Mercury orbit.

actors that may be instantiated by the Simulink model in the UDK application have to be predefined. The number, location, and properties of the actors can be changed on the fly by the Simulink model. UDK function delegates are used to forward UDK events such as collisions back to Simulink for processing by the relevant model.

Timing for the interactions between the Simulink models and the UDK model may be synchronous or asynchronous. In the asynchronous case, the Simulink models run as fast as they can and update the data on the UDK model as soon as possible. The UDK engine attempts to keep a "real time" view so that time progresses continuously. If synchronous interaction with the Simulink model is wanted, then both the UDK model and Simulink model can progress time step by time step in lock step. In cases where the Simulink model runs substantially slower than real time and a real time playback is desired, the control events from the Simulink model can be recorded in a file with time stamps. These files can be played back in real time after the simulation is over. The toolkit provides for interaction with one or more human users via the UDK or via the Simulink model. Various input/output devices are supported including keyboard, mouse, and force reflective joystick. In the case of the force reflective joystick, the force feedback values are computed by the Simulink model and presented back to the user via the standard Simulink support for force reflective devices.

4 **Proof-of-concept application**

At present we are developing a simple proof-ofconcept astrophysics application. Users can generate solar systems with varying numbers of planets, with varying orbits, masses, etc. Our goal is to develop a model that allows users to explore the concept of launch windows for space exploration. We are examining other models that are more amenable to the use of force reflective devices such as a pool table simulation as further demonstrations of our toolkit's capabilities.

In our proof-of-concept application, the Simulink model configures the rotation settings for Mercury and



Figure 4 – Screenshot of UDK virtual solar system driven by Simulink model.



Figure 5 – Simulink rotation model with MATLAB GUI for controlling direction of rotation.

Sun actors and implements the orbit of Mercury about the Sun. The UDK interface can process spawn, set rotation rate, and set position commands for all the actors. Figure 3 above contain a screenshot of the Simulink orbit model. Figure 4 below shows a screenshot of the virtual solar system as it appears within the UDK. Finally, Figure 5 shows a screenshot of the Simulink rotation model for a solar system actor, which includes a MATLAB GUI that is used to vary direction and rotational speed parameters of the Simulink model. One could use a similar technique to create "instructor" and "student" interfaces to the same simulation. For example, through the UDK interface, the workers undergoing training explore and interact with a virtual world (say a nuclear power station or chemical processing plant) to achieve a specific objective. Meanwhile, the instructor, monitoring the workers' progress via the UDK interface, can use the MATLAB GUI to trigger specific training scenarios that will require an appropriate response from the workers. The full physics model implemented within Simulink would provide the workers the ability to view via the UDK the consequences of their various decisions.

5 Conclusions and future work

As cyber-physical systems become more complex, it will be increasingly important to understand our models of these systems so that we can make better design and implementation choices. Our toolkit allows engineers to link their Simulink models to the UDK for first person, interactive visualization of simulation results. We extend previous work by facilitating two-way communication and interaction between Simulink and the UDK, and this allows for partitioning of the simulation between the two tools to take advantage of each tool's strengths and capabilities. As a result, each Simulink model may control spawning and manipulation of actors within the UDK, and relevant events occurring within the UDK virtual world can be forwarded back to Simulink for processing by the Simulink model. The use of TCP communications also allows multiple Simulink models to interact within the UDK while one or more human users interact with the virtual world through the UDK itself, as would first person gamers.

Future work includes the development of applications which can exploit the Simulink force reflective device interface, the investigation of mechanisms for switching on-the-fly between low quality and high quality visualization, and an exploration of synchronization mechanisms for timing critical simulations distributed between Simulink and the UDK.

6 References

[1] Simulink/Matlab by Mathworks, Inc.,

www.mathworks.com

[2] Unreal Development Kit 3.0 by Epic Games, http://udk.com

[3] Michael Zyda, John Hiles, Alex Mayberry, Casey Wardynski, Michael Capps, Brian Osborn, Russell Shilling, Martin Robaszewski, and Margaret Davis, "Entertainment R&D for Defense," *IEEE Computer Graphics and Applications*, January/February 2003, pp. 28-36.

[4] Jun Lai, Wei Tang, and Yihui He, "Team Tactics in Military Serious Game," *2011 Fourth Int. Symposium on Computational Intelligence and Design*, pp. 75-78.

[5] Z. Kriz, R. Prochaska, C.A. Morrow, C. Vasquez, H. Wu, and Rizwan-uddin, "Unreal III Based 3-D Virtual Models for Training at Nuclear Power Plants," *Proc. 1st Int. Nuclear and Renewable Energy Conference (INREC10)*, Amman, Jordan, March 21-24, 2010, pp. INREC10-1 – INREC10-5.

[6] Danna Voth, "Gaming Technology Helps Troops Learn Language," *IEEE Intelligent Systems*, September/October 2004, pp. 4-6.

[7] Xu Feng, Shan Daguo, and Yang Hongchen, "Simulation Research of Crime Scene Based on

UDK," *Information Science and Engineering* (*ICISE*), 2010 2nd International Conference on, 4-6 Dec. 2010, pp. 1-4.

 [8] Allison Mathis, Kingsley Fregene, and Brian
 Satterfield, "Creating High Quality Interactive
 Simulations Using MATLAB and USARSim", https://robotics.ucmerced.edu/Robotics/wspapers/IRO
 <u>S USARSimWS AMathis Final.pdf</u>

Nonlinear Model Structure Identification Based on Kernel Visualization

L. Keviczky and Cs. Bányász

Computer and Automation Research Institute of the Hungarian Academy of Sciences Control Engineering Research Group of the HAS at the Dept. of Automation and Applied Informatics Budapest University of Technology and Economics H-1111 Budapest, Kende u 13-17, HUNGARY

Abstract - The success of the nonlinear dynamic system identification strongly depends on the applied model structure. Nonlinear systems have almost infinite varieties of structures. This paper shows a simple structure identification technique based on image processing recognition method to distinguish between Hammerstein and Wiener models.

Keywords : nonlinear model, structure identification

1 Introduction

Nonlinear dynamic systems have infinite variety so it can not be expected that a unique optimal control solution exists for all of these complex processes. However, the history and development of nonlinear control systems show that one can expect relatively simple methods, which can be close or even reach the effectiveness of linear methodology. A very wide class of approaches are based on the usual Jacobean linearization [5] of a nonlinear system (*NS*). A large class of nonlinear systems can be made to have linear input-output behavior through a choice of nonlinear state feedback control law [10]. Other approaches assume special topology, when the structure of the *NS* makes the linearization possible. Such *NS* classes are e.g. the bilinear and the block-oriented factorable (cascade) systems [4].

The most well known factorable models are the simple Hammerstein and Wiener models.



Fig. 1 The simple Wiener and Hammerstein models

The simple Wiener model N_W (shown in Fig. 1a) is a cascade structure of a linear dynamic $Y^{W,dyn}$ and a nonlinear static $N^{W,stat}$ terms connected in series, i.e.

$$N_{\rm W} = N^{\rm W,stat} Y^{\rm W,dyn}$$
 or simply $N_{\rm W} = N^{\rm W} Y^{\rm W}$ (1)

The simple Hammerstein model $N_{\rm H}$ (shown in Fig. 1b) is a cascade structure of a nonlinear static $N^{\rm H,stat}$ and a linear dynamic $Y^{\rm H,dyn}$ terms connected in series, i.e.

$$N_{\rm H} = Y^{\rm H,dyn} N^{\rm H,stat}$$
 or simply $N_{\rm H} = Y^{\rm H} N^{\rm H}$ (2)

Note that the order of the nonlinear operators in the formulas is opposite to the order of the blocks shown in the figures. (This is a usual source of mistakes calculating transfer characteristics of open- and closed-loop nonlinear schemes.)

There exists many control methods published, including methods from the authors (e.g. [1], [7]), which can be applied for the Wiener and Hammerstein nonlinear dynamic model classes. These methods assume a proper process model identification procedure to obtain a good approximate model. A natural question always arises, whether the true process falls into these model structures.

Many nonlinear dynamic systems with input signal u(t) and output signal y(t) can be approximated in the vicinity of the working point at least by the so-called Volterra integral or the Volterra weighting function model [4]

$$y(t) = g_0 + \int_{\tau_1=0}^{t} g_1(\tau_1)u(t-\tau_1)d\tau_1 + \int_{\tau_1=0}^{t} \int_{\tau_2=0}^{t} g_2(\tau_1,\tau_2)u(t-\tau_1)u(t-\tau_2)d\tau_1d\tau_2 + \dots \quad (3)$$
$$= \sum_{i=0}^{n} \int_{\tau_1=0}^{t} \dots \int_{\tau_i=0}^{t} g_i(\tau_1,\dots,\tau_i)\prod_{j=1}^{i} u(t-\tau_j)d\tau_1\dots d\tau_i$$

Similarly to the discrete time description of the linear systems, a multi-dimensional convolution sum describes

the relation between the sampled input u[k] and output signals y[k]

$$y[k] = h_0 + \sum_{\kappa_1=0}^k h_1[\kappa_1]u[k-\kappa_1] +$$

+
$$\sum_{\kappa_1=0}^k \sum_{\kappa_2=0}^k h_2[\kappa_1,\kappa_2]u[k-\kappa_1]u[k-\kappa_2] + \dots \qquad (4)$$

=
$$\sum_{i=0}^n \dots \sum_{\kappa_i=0}^k h_i[\kappa_1,\dots,\kappa_i]\prod_{j=1}^i u[k-\kappa_j]$$

We call (3) the Volterra weighting function series. Another name for it is the Gabor-Kolmogorov series. The process is characterized by its Volterra kernels:

 $g_n(\tau_1,...,\tau_n)$ n = 0,1,2,... in the continuous time case and

 $h_n [\kappa_1, \dots, \kappa_n]$ $n = 0, 1, 2, \dots$ in the discrete time case.

(There are several methods to compute the relation between the continuous and discrete time kernels of a system.)

In the engineering practice only second order Volterra kernels are determined and used. The coefficients of a such a form are not difficult to estimate relatively simple model identification procedure, contrary to the above cascade models. Therefore it is very useful to perform a structure identification to determine which cascade model is the best for the measured process data.

While the Volterra kernels can easily be derived from the block oriented models, the structure and the parameters of the block oriented models cannot be computed in a trivial way from the estimated Volterra kernels.

The constant term g_0 is equal to h_0 , the first degree Volterra kernel $g_1(\tau_1)$ is equal to the weighting function

of the linear channel. The quadratic channel can be described as

$$y_{2}(t) = \int_{\tau_{1}=0}^{\infty} \int_{\tau_{2}=0}^{\infty} \int_{\tau=0}^{\infty} g_{4}(\tau)g_{2}(\tau_{1}-\sigma)g_{3}(\tau_{2}-\sigma)$$

$$u(t-\tau_{1})u(t-\tau_{2})d\tau d\tau_{1}d\tau_{2}$$
(5)

thus the quadratic kernel can be calculated by the following convolution integral

$$g_2(\tau_1,\tau_2) = \int_{\tau=0}^{\infty} g_4(\tau) g_2(\tau_1-\tau) g_3(\tau_2-\tau) d\tau$$
(6)

It seems there exists an unequivocal relation between the parameters of the block oriented model and the Volterra kernels. The transformation is unequivocal only in one direction because the weighting functions $g_i(\tau_i)$ of the block oriented models cannot be reconstructed from the identified Volterra kernels.

On the basis of the above considerations a so-called parametric Volterra model can be obtained for the second degree discrete-time Volterra kernel

$$y \left[k + d \right] = h_0 + \sum_{\kappa_1 = 0}^{\infty} h_1 \left[\kappa_1 \right] u \left[k - \kappa_1 \right] + \sum_{j=0}^{\infty} \sum_{\kappa_2 = 0}^{\infty} h_{2j} \left[\kappa_1 \right] u \left[k - \kappa_1 \right] u \left[k - \kappa_1 - j \right]$$
(7)

which is linear in the parameters. Here a usual time delay d was also introduced. If the parameters in this second order kernel are identified they can give information on the possible cascade structure for a second, more accurate identification method. So the first estimations of the kernel parameters \hat{h}_0 , \hat{h}_1 and \hat{h}_{2j} and their variances can be easily obtained from a linear regression.



Fig. 2 Plots of the Volterra kernels of a Hammerstein model with first-order lag term



Fig. 3 Plots of the Volterra kernels of a Wiener model with first-order lag term

2 Model Structures in the Parameter Space of a Volterra Kernel

The kernel parameter distributions of the two basic cascade models shown in Fig. 1 can be easily determined by computer simulation.

The Figs. 2 and 3 illustrates the Volterra kernels for the Hammerstein and for the Wiener models [4]. The Volterra kernels of Figs. 2 suggest that for Hammerstein model the quadratic kernels differ from zero only at the main skew (or secondary) diagonal. On the Fig. 3 the sections $h_2[\kappa_1,\kappa]$; $\kappa = \text{const}$ are proportional to each other, and the level lines are straight (subdiagonals), because $h_2[\kappa_1,\kappa_1] = h_1^2[\kappa_1]$.

It is possible to compute further kernel patterns for other cascade block-oriented nonlinear model classes.

3 Analytical Structure Indices

The quadratic kernels of the Hammerstein models differ from zero only at the main diagonal. This feature can be seen from the plot of $h_2[\kappa_1,\kappa_2]$ and from building the following index:

$$\alpha_{1} = \frac{\frac{2}{m(m+1)} \sum_{\kappa_{1}=0}^{m-1} \sum_{\kappa_{2}=\kappa_{1}+1}^{m} h_{2}^{2} [\kappa_{1}, \kappa_{2}]}{\frac{1}{m+1} \sum_{\kappa_{1}=0}^{m} h_{2}^{2} [\kappa_{1}, \kappa_{2}]}$$
(8)

In (8) the mean square values of the off-diagonal elements are divided by that of the main diagonal. Here and further on any characteristic measure of the size of the element can be used instead of the square value (e.g., absolute value). The main skew diagonal is the straight line for which $\kappa_1 = \kappa_2$ and *m* is the memory of the

kernels of the discretized model. It is easy to see that α_1 becomes zero only for the Hammerstein model.

The quadratic Volterra kernels along the main skew diagonal are proportional to the squares of the linear kernels in the simple Wiener models. Form the ratio of them as a discrete time function

$$\beta_1 \left[\kappa_1 \right] = \frac{h_2 \left[\kappa_1, \kappa_1 \right]}{h_1^2 \left[\kappa_1 \right]} \tag{9}$$

and its average value

$$\overline{\beta}_4 = \frac{1}{m+1} \sum_{\kappa_1=0}^m \beta_4 \left[\kappa_1\right] \tag{10}$$

Then the normalized deviation of (9)

$$\alpha_4 = \frac{\frac{1}{m+1} \sum_{\kappa_1=0}^{m} \left[\beta_4 \left[\kappa_1\right] - \overline{\beta}_4\right]^2}{\overline{\beta}_4^2}$$
(11)

becomes zero only for the Wiener model.

4 Pattern Recognation of the Kernel Shapes

The above analytical indices are exact measures for finding Wiener and Hammerstein model structures. However, it is very difficult to compute a statistical probe, how close they are to zero in case of noisy measurements. In the practice these measures are more or less heuristical norms for the model structure determination.

In our practice it was found that relatively simple pattern (shape) recognition algorithms work more robust than the analytical indices. Mathematicians typically define shape as an equivalence class under a group of transformations. This definition is incomplete in the context of visual analysis. This only tells us when two shapes are exactly the same. We need more than that for a theory of shape similarity or shape distance. The statistician's definition of shape, e.g. [3] or [6], addresses the problem of shape distance, but assumes that correspondences are known. Other statistical approaches to shape comparison do not require correspondences – e.g. one could compare feature vectors containing descriptors such as area or moments – but such techniques often discard detailed shape information in the process. Shape similarity has also been studied in the psychology literature.

There are several extensive surveys of shape matching in computer vision literature. There are basically two approaches: *feature-based*, which involve the use of special arrangements of extracted features such as edge elements or junctions; and *brightness-based*, which make more direct use of pixel brightness [9].

Assume that the variances of the estimated kernel parameters are denoted by σ_0 , $\sigma_1[\kappa_1]$ and $\sigma_{2j}[\kappa_1, j]$ and available from the linear regression performed using the model (7) linear in the parameters. Compute the relative significance factors:

$$\gamma_0 = \frac{\hat{h}_0}{\sigma_0} \quad , \ \gamma_1 \left[\kappa_1\right] = \frac{\hat{h}_1 \left[\kappa_1\right]}{\sigma_1 \left[\kappa_1\right]} , \ \gamma_2 \left[\kappa_1, j\right] = \frac{\hat{h}_{2j} \left[\kappa_1, j\right]}{\sigma_{2j} \left[\kappa_1, j\right]} (12)$$

For the second degree kernel transform the estimated parameters to a dot matrix \overline{K} putting 1 to those elements where the $\gamma_2[\kappa_1, j] \ge \delta_2$ and 0 elsewhere. Here δ_2 can be given heuristically or can be selected similarly than in the well-known Student's *t-probe* for linear model

identification. This simple method transforms the first model into the second *brightness-based case*.

For the Hammerstein model the special pattern is the skew diagonal of the kernel matrix. So it is possible to use a classical character recognition package, which has to find a "slash" (i.e. /) character for the transformed matrix \bar{K} .

For the Wiener model it is required to recognize special subdiagonals, which form a Toeplitz matrix structure. So it is possible to use a classical character recognition package , which has to find a "backslash" (i.e. \) character for all Toeplitz lower sub matrices in \overline{K} .

To the application of the *feature-based* method the feature extraction rule set is quite difficult to construct for these simplest cascade models. (This work, however, can not be avoided for higher order dynamics.) Instead, it is not difficult to construct another algorithm, which does not transform the estimated kernel to \overline{K} , but uses special masks, shown in Fig. 4 to ease the feature extraction.



Fig. 4 Filter masks for heuristic shape recognition

Using these masks it is not so difficult to form the feature extraction rule set and these rules can be applied in the combination with the analytical indices.



Fig. 5 Plots of the Volterra kernels of a Wiener model with second-order lag term

5 Generalization Possibilities

The above presented methods for Hammerstein models can be used without modification for the so-called generalized Hammerstein model, too, because the second degree kernels also concentrate on the main skew diagonal.

The generalized Wiener model with a first order lag can be also handled in the same way as shown above. Unfortunately more complex Wiener models (e.g. second order dynamics) need different feature set to be applied, because the kernel contour level shapes are considerably different (see Fig. 5 for the obtained kernel behaviors). The only practical solution is to collect a proper shape database for the different structure primitives. This is not difficult, because the database can be computed by computer simulation, however, very time consuming.

6 Conclusions

The purpose of the paper was to present the special feature sets of nonlinear dynamic cascade Wiener and Hammerstein models using the second degree Volterra kernel parameters.

The special shapes on the two-dimensional plots of the elements of the kernel matrix invokes the application possibilities of different pattern recognition algorithms for nonlinear structure identification purposes. The major contribution is that the presented methods are very good to distinguish the Wiener and Hammerstein model classes. Further structure determination, especially within the Wiener models, require more complex shape recognition algorithms and/or shape features extraction data bases.

7 References

- Bányász, Cs. and L. Keviczky. "An iterative nonlinear controller refinement scheme", *IFAC Symp. Nonlinear Control Systems NOLCOS'04*, Stuttgart, D, 173-178 (CD), 2004.
- [2] Belongie, S., J. Malik and J. Puzicha. "Shape Matching and Object Recognition Using Shape Contexts", *IEEE Trans. On Pattern Analysis and Machine Intelligence*, **24**, 24, 509-522, 2002.
- [3] Bookstein, F.L. Morphometric Tools for Landmark Data: Geometry and Biology, Cambridge Univ. Press, 1991.
- [4] Haber, R. and L. Keviczky. *Nonlinear System Identification Input-Output Modeling Approach*, Kluwer Academic Publishers, 1999.
- [5] Isidori, A. *Nonlinear Control Systems*, 3rd ed., Springer Verlag, London, 1995.
- [6] Kendall, D. "Shape manifolds, procrustean metrics and complex projective spaces", *Bull. London Math. Soc.*, **16**, 81-121, 1984.
- [7] Keviczky, L. and Cs. Bányász. "Generic two-degree of freedom control systems for linear and nonlinear

processes", J. Systems Science, 26, 4, 5-24, 2001.

- [8] Keviczky, L. and Cs. Bányász. "Optimal PID tuning for Wiener/Hammerstein nonlinear systems", 27th Int. Conf. of IEEE Industrial Electronics Society IECON'01, Denver, CO, USA, 746-751, 2001.
- [9] Nelles, O. *Nonlinear System Identification*, Springer, 2001.
- [10] Sastry, S. Nonlinear Systems Analysis, Stability and Control, Springer, 1999.

SESSION SIMULATION

Chair(s)

TBA
Computer Modeling and Simulation of Ground Penetrating Radar using Finite Difference Time Domain Code

J. Hebert^{1,2,}, T. Holzer D.Sc.¹, T. Eveleigh PhD¹, Dr. Shahryar Sarkani, PhD¹, and J. Ball PhD²

¹ George Washington University Washington, D. C.

² Naval Surface Warfare Center Dahlgren Division, Dahlgren, VA

Abstract

Modeling and simulation results of a system analysis of Ground Penetrating Radars (GPR) using Finite Difference Domain (FDTD) techniques are presented. Time Performance issues with GPRs need to be isolated in order to optimize of the radar's ability to detect and identify buried Using a system engineering approach, FDTD obiects. models were used to characterize the variables associated with the GPR to improve GPR detection process minimally affected by external sources of variability. These experiments make changes to the GPR's inputs while measuring the output response to identify issues and optimize performance. FDTD computer simulations produce idealistic environments that allow examination of the individual effects on the response. This paper provides a system engineering overview of the operations and processes of GPR systems and how MATLAB based FDTD computer simulations can be used to model and improve them. A plan for future work is presented.

Key Words: System, Modeling, Simulation, GPR, System Engineering, FDTD

1 Introduction

A system analysis is performed to isolate and understand the factors that affect a GPR's ability to detect and identify buried objects. The use of FDTD computer models and simulations in the systems analysis enable GPRs to be designed to do what they should do. Many factors affect a GPR's ability to detect and identify subsurface objects. A major factor is the electrical characteristics of the object and of the material in which the object is buried. Although several attempts to use carefully prepared test sites to make measurements to understand the magnitude of these performance limiting effects have been made, a comprehensive system engineering based investigation has not been reported to date. This is due to the time and expense of preparation of configurations for measurement based investigations, and the lack of any reported system engineering efforts applied to GPR.

The effects of changing simple variables such as surface and soil constituent properties on the GPR radar's output are not separable. Often the effects are unidentifiable in measurements made under field conditions. A synthetic data set produced by FDTD computer simulations allows the separation of input variables to better understand their effects on the output response of the radar. These simulations produce idealistic environments and test configurations to allow close examination of the individual effects of these variables on the response. This paper presents a version of FDTD code that has been implemented in MATLAB to model and simulate a GPR's performance. This version of the code is intended for use by researchers to observe, analyze, and understand how different system input variables affect the GPR and its performance.

Modeling the GPR as a system using FDTD techniques allows a set of specially designed experiments where deliberate changes are made to the input variables so that changes in the output response can be observed and performance limiting issues can be easily identified. Three MATLAB models and simulations are presented in this paper: (1) a model for calculating the Fresnel reflection and transmission coefficients for perpendicular and parallel polarity incident waves as a function of grazing angle, (2) a one dimensional (1D) FDTD model for comparison with the Fresnel model, and (3) a three dimensional (3D) FDTD model to allow simulation of the response to the GPR of changing various inputs. The amount of energy that is reflected at the boundary of two media (e.g., soil and buried target) with different permittivity is given by the Fresnel coefficient. The changes in output response were observed while controlling various input variables. The initial results of controlling the conductivity and permittivity of the soil and targets are presented. Conductivity is a measure of a material's ability to conduct an electric current. Permittivity relates to a material's ability to transmit (or "permit") an electric field.

David Montgomery states that one of the applications of experiment design is the identification of design parameters that work well over a wide range of conditions in order to determine the design parameters that most impact product performance [1]. Variables to be

considered in simulation of the GPR as a system are: (1) radiated waveform, (2) depth of penetration versus frequency, (3) transmitter antenna type, (4) height and grazing angle, (5) surface, soil, and target properties, (6) target characteristics, (7) clutter (8) moisture content, (9) interference (10) receiver antenna type, (11) signal collection resolution and rate, (12) signal processing techniques, and (13) optimizing response of all input variables to maximize detection and reduce false Our initial research focused on two input alarm rate. variables that were found to greatly affect the GPR's output response: conductivity and permittivity. A detailed examination of the response of the system to changing these two input variables allowed for optimization to obtain the most accurate possible output response. The system engineering goal of this modeling and simulation effort is to define what a GPR system should be rather than applying a classical approach of determining of what a GPR can be.

2 GPR System Analysis

This section discusses the motivation for using the SE tools of modeling and simulation in the development of the System Engineering GPR computer model and the selection of FDTD techniques to perform the system simulations and analyses. Surface penetrating active sensors (SPAS), such as GPR, and ultrasound have hundreds of real world applications for their ability to "see into" and characterize solid and semi-solid substrates. As such, they are highly desirable functional components for a growing number of advanced systems. The computer models and mathematics for surface penetrating active sensors can be quite involved with only a few sensor models developed for specific instruments, for specific applications, and/or for specific environments of use. To date, no general sensor system characterization models exist that can deterministically characterize sensor technology or examine the parametrics, and tune in a response to an intended environment of use and a desired target resolving capability. The ability to deterministically match system sensing needs to SPAS capabilities would be of great interest to the systems engineer. At present, it is very difficult for all but the most highly trained experts to know what SPAS capabilities might work under what given set of This paper presents an extensible approach conditions. towards the allocation of sensing performance requirements to determined SPAS solution technologies.

The goal of this system engineering analysis is to identify GPR system deficiencies and what can be done to improve the system's performance. Five important steps in the system engineering process include: (1) critical needs are identified, (2) current capabilities are assessed, (3) new or existing capabilities are explored, (4) prototyping or modeling and simulation are implemented and (5) final system deployed. The FDTD model for this research facilitates the system analysis required by steps 2, 3, and 4.

This approach could provide the systems engineer with a requirements-driven solution synthesis by better characterizing and populating the architectural trade space with valid SPAS alternatives that represent a range of possible SPAS solutions.

3 Radar Ground Penetrating

To analyze the GPR as a system we must first understand the components and functions of the GPR. This radar is used for the detection of objects buried below the surface. A GPR consists of a transmitting and receiving antenna, a source connected to the transmitting antenna, and signal processing equipment connected to the receiving antenna. The type of antennas, choice of the transmitted signal, and method of signal processing are all system variables that affect the output response and performance of the GPR As such each is a candidate for optimization as part of the GPR's system architecture and design.



Figure 1. Schematic drawing of typical GPR.

Figure 1 shows a GPR system and operating environment with the signals that are generated by the system. Filtering out the interference caused by the direct and the ground bounce signals in order to see the reflection of the return from the target may be necessary. The operating environmental variables that must be modeled in a GPR simulation include the two antennas, the electrical characteristics: permittivity, ε , conductivity, σ , and permeability, μ , of the air above the surface, the subsurface, and the target. Other variables include the height above the surface of the antenna, the separation distance between the antenna, and the depth of the target. Most of these variables are related or dependent on the other variables such that modeling them one at a time would cause unaccounted for errors in the output response. The best that can be done is to control the variables one at a time, while including all the variables in the GPR model and simulation. The research presented here includes a threedimensional, finite-difference time-domain (3D-FDTD) system analysis of the GPR that accounts for many of these variables simultaneously within the problem space.

| Material | ε _r : Davis | ε _{r:} Daniels | Velocity | Velocity | |
|------------------|------------------------|-------------------------|-----------|-----------|--|
| | and Annan | et al | (m/ns) | (ft/ns) | |
| | (1969) | (1995) | | | |
| Air | 1 | 1 | 0.3 | 0.96 | |
| Distilled water | 80 | | 0.03 | 0.11 | |
| Fresh water | 80 | 81 | 0.03 | 0.11 | |
| Sea water | 80 | | 0.03 | 0.49-0.57 | |
| Fresh water ice | 3-4 | 4 | 0.15-0.17 | 0.35-0.49 | |
| Sea water ice | | 4-8 | 0.11-0.15 | 0.28-0.35 | |
| Snow | | 8-12 | 0.09-0.11 | 0.35-0.50 | |
| Permafrost | | 4-8 | 0.11-0.16 | 0.40-0.57 | |
| Sand, dry | 3-5 | 4-6 | 0.12-0.17 | 0.18-0.31 | |
| Sand, wet | 20-30 | 10-30 | 0.05-0.09 | 0.57-0.70 | |
| Sandstone, dry | | 2-3 | 0.17-0.21 | 0.31-0.44 | |
| Sandstone, wet | | 5-10 | 0.09-0.13 | 0.35-0.49 | |
| Limestones | 4-8 | | 0.11-0.15 | 0.37 | |
| Limestone, dry | | 7 | 0.11 | 0.35 | |
| Limestone, wet | | 8 | 0.11 | 0.25-0.44 | |
| Shales | 5-15 | | 0.08-0.13 | 0.33-0.40 | |
| Shale, wet | | 6-9 | 0.10-0.12 | 0.18-0.44 | |
| Silts | 3-30 | | 0.05-0.13 | 0.18-0.44 | |
| Clays | 5-40 | | 0.05-0.13 | 0.16-0.44 | |
| Clay, dry | | 2-6 | 0.12-0.21 | 0.40-0.70 | |
| Clay, wet | | 15-40 | 0.05-0.08 | 0.16-0.25 | |
| Soil, sandy dry | | 4-6 | 0.12-0.15 | 0.40-0.49 | |
| Soil, sandy wet | | 15-30 | 0.05-0.08 | 0.16-0.25 | |
| Soil, loamy dry | | 4-6 | 0.05-0.08 | 0.40-0.49 | |
| Soil, loamy wet | | 15-30 | 0.07-0.09 | 0.22-0.31 | |
| Soil, clayey dry | | 4-6 | 0.12-0.15 | 0.40-0.49 | |
| Soil, clayey wet | | 10-15 | 0.08-0.09 | 0.25-0.31 | |
| Coal, dry | | 3.5 | 0.16 | 0.53 | |
| Coal, wet | | 8 | 0.11 | 0.35 | |
| Granites | 4-6 | | 0.12-0.15 | 0.40-0.49 | |
| Granites, dry | | 5 | 0.13 | 0.44 | |
| Granites, wet | | 7 | 0.11 | 0.37 | |
| Salt, dry | 5-6 | 4-7 | 0.11-0.15 | 0.37-0.49 | |

Table 1. Relative permittivity, ε_r , and EM velocity for selected geological materials

One variable that has a large impact on a GPR's performance is the permittivity. Table 1 [2] shows the relative permittivity and electromagnetic wave velocity for common subsurface materials. The amount of energy that is reflected at the boundary of two media with different permittivity is given by the Fresnel coefficient. For air to soil with permittivity, ε_r , and permeability, μ_r , the index of refraction (Fresnel reflection coefficient) is described by:

$$n = \sqrt{\frac{\epsilon \mu}{\epsilon o \mu o}} = \sqrt{\epsilon_r \, \mu_r} \tag{1}$$

This relationship is used to illustrate the changes in the electromagnetic waves at the interface of two materials with different permittivity and permeability in the results section below. One observes that electromagnetic waves pass through the earth and the receiving antenna records the timing and magnitude of the arriving energy. A GPR image is actually an image directly related to the dielectric properties of the subsurface. The dielectric constant controls the velocity and the path of electromagnetic waves, including those reflected off objects below the surface.

3 FDTD Technique

FDTD techniques relate the surface currents and charges in a problem space that are modeled by Maxwell's curl equations which are:

$$\nabla \times \mathbf{E} = -\frac{\partial \mathbf{B}}{\partial t}$$
(2)
$$\nabla \times \mathbf{H} = \mathbf{J} + \frac{\partial \mathbf{D}}{\partial t}.$$
(3)

These equations are used to develop a solution approach known as the finite difference formulation. A detailed development of the equations for the three dimensional version of the FDTD code is presented in a thesis by Williford [6]. Although the FDTD approach can be carried out both in the time and frequency domain, the model used for this research implements the time domain formulation. FDTD models the propagation and interaction of an electromagnetic wave in a region of space that may contain any object. This method is different from the integral equation method in that it analyzes the interaction of the incident wave with a portion of the structure at a given instant in time rather than solving the entire problem at once. Yee [6] first suggested the FDTD formulation for solving Maxwell's two curl equations (1) and (2), stating that the derivatives in these equations could be expressed as differences of the field values between neighboring positions, both temporally and spatially. These difference equations vield the values of the field at a given location in time and space if the values at all positions in the problem space are known at an earlier time.

The solution of an electromagnetic interaction problem by the FDTD technique is straight forward. For our system model, the problem space is divided into a lattice of uniform sized cells. As shown in Figure 3, the gridding procedure involves placing the components of the electric (E) and magnetic (H) fields around a unit cell and evaluating the field components at alternate half-time steps.



Figure 3. 3D Yee cell [13]

By alternating between the E and H fields, a central difference expression can be developed for both the space and time derivatives that maintains a higher degree of accuracy than either a forward or backward difference formulation. The problem solution proceeds by time-stepping throughout the problem space, repeatedly solving the finite difference form of Maxwell's two curl equations. In this fashion, the incident wave is tracked through the problem

space as it intercepts and interacts with the targets, at layer interfaces, and with other objects in the problem space.

Yee [6] developed the FDTD algorithm in 1966 as a method to compute the waveforms of pulses scattered from infinitely long, rectangular cross section, conducting cylinders Rymes [7] used FDTD to analyze data from direct lightning strikes to a NOAA C- 130 aircraft. This code was later modified and used by Hebert and Sanchez-Castro [8] to analyze the data from inflight lightning strike measurements by a CV-580 aircraft and by Williford [5] to explore the validity of different boundary conditions using FDTD to model an F-16 aircraft. Williford, Jost, and Hebert [8] found using FDTD absorbing boundary conditions with FDTD produced better results than the perfectly electrically conducting (PEC) reflective boundary conditions originally used by Yee [6] but at the cost of longer run times.

Based upon these efforts, FDTD has been shown to be useful for the modeling and analysis of electromagnetic interaction with systems. These codes are easily adapted to a variety of materials in the problem space leading directly to their choice to analyze GPR data. In addition, nonlinearities and time-varying quantities can be represented in the problem space grid, if the needed equations can be written at the appropriate location. In addition, FDTD codes written in MATLAB are easily adapted to parallel processing and multiprocessor systems.

4 3D-FDTD Models and Simulation

The FDTD code calculates the solutions to Maxwell's Equation in their differential form. FDTD solutions are simple and depending on the choice of time steps and grid lengths provide extremely accurate representations of the interaction of electromagnetic waves and materials with different constituent properties. Modeling using FDTD techniques allows the observation of changes in response due to changing input variables without the expensive cost of physical experiments.

There are many versions of the 3D-FDTD code. Some are readily available for download on the internet. Commercial versions of the code and versions that are reported in scholarly journals come in packages are not open source and are not available for researchers. For this reason, a GPR model and simulation program implementing FDTD techniques was developed in MATLAB.

Many different algorithms exist for target detection and identification, noise and interference suppression, removal of direct and air wave effects, and correction of attenuation losses. The input data for the research and comparison of these algorithms is provided by the FDTD techniques implemented in the MATLAB code.

Previous researchers have successfully used 3D-FDFD techniques to investigate some aspects of a GPR's performance [5, 11]. While helpful, these studies produced

only limited results. Under some physical soil conditions, the recognized landmine signature possesses high quality contrast while under other conditions no signature is detected. Fritzsche [3] demonstrated via modeling that GPR signals at 900 MHz would be strongly attenuated in moist soil. Trang [4] found through simulations and experiments with a GPR signals operating at 600-800 MHz that nonmetallic mines were easier to detect in moist soil.

The FDTD computer model implemented as part of this research facilitates the analysis of complex dielectric constant of soil and attenuation of GPR signals. In addition, the system model is capable of plotting the complex dielectric constant of soil coupled with the attenuation of GPR signals versus soil physical properties.

To predict the performance of electromagnetic sensors sub-systems, it is common practice to use models that estimate the soil's characteristics including dielectric properties. Trang found that no current model exists to completely describe all the electrical properties of a soil type [4]. Measurements to baseline GPR operational performance made at many sites worldwide are helpful but still leave a great deal unknown due to uncertainties caused by factors such as soil composition, layering, clutter, rock and other undesired artifacts recorded in the measurement. Alternatively, the FDTD computer models and simulations allow the variables associated with GPR system to be researched and characterized.

FDTD techniques model many variables that are controllable while some variables are not. Using FDTD synthetic data allows one to control what might otherwise be undefined or uncontrollable variables. The system engineering goal for the simulations is to find bounds for the input values of the uncontrollable variables which make the systems performance predictable and manageable. Thus a GPR system design can be optimized to effectively handle a wider variety of operational conditions

Figure 2 shows a B-mode image of buried pipes. A B-mode image is produced by sweeping a narrow beam while transmitting pulses and detecting echoes along a series of closely spaced scan lines. The algorithm for B-mode image simulation and processing includes calculation of the amplitude and two-way time delay of a signal reflected from each layer of a multi-layered media; simulation of echo signals, clutters, speckle and impulse noise; construction of synthetic range profile; and image formation.



of

al

Meters

Figure 2. Example of a B mode plot

Belli, Rappaport, Udall, Hines and Wadia-Fracetti [10] provided an excellent example of a subsurface tunnel modeled in FDTD, as illustrated in Figures 4 and 5.



Figure 4. (a) 3D tunnel geometry, and (b) Detail of y - z plane indicating sensor location when $\theta = 0^{\circ}$ [11]



Figure 5. 3D-FDTD simulated B-scan contours (air-filled tunnel buried in sand with backgrounds removed) [10]

These simulations show how measured GPR data can be faithfully modeled in FDTD and how FDTD simulations can be used to model a GPR system's performance. It shows that the FDTD model produces Typical B-scan contours and the extracted hyperbolas for the tunnel example that can be seen in Figure 5 with the background reflection at the air/sand interface removed. Four angles are selected for 3D B-scan simulation: 0°, ≈23.96°, 45° and $\approx 53.13^{\circ}$. The hyperbolas extracted from the B-scan simulations were compared to a library of hyperbolas generated by 2D FDTD to determine the angle of the GPR waves travel path. By comparing the angles from the simulations with measured data, these angles were found to produce the B-scans that most closely match the measured ones. The results are summarized in Table 2. The determined angles are well matched to the actual angles. Again, and as expected, the case of $\theta = 45^{\circ}$ results in the largest error in determined θ .¹⁰

Table 2. Tunnel Example Correlation Results

| 3D simulation angle, Ø | Best 2D correlation | Maximum error (Distance from tunnel in <i>s</i> -direction | Mean error |
|---|------------------------|--|---------------|
| 0° | 0° | 180.0 ps at 2.25 m | 73.9 ps |
| arctan (4/9) ≈ 23.96° | 24° | 93.8 ps at 2.63 m | 38.0 ps |
| 45° | 24° | 152.1 ps at 3.39 m | 47.8 ps |
| Arctan (4/3) $\approx 53.13^{\circ}$ | 54° | 535.9 ps at 4.0 mm | 206.2 ps |

5 Simulation Results

Dependence on Frequency: System analysis begins by selecting one input and determining its effect of the system's performance. If one extends the analysis of system inputs to the effects of frequency on the depth and resolution like that presented by GST^{11} , the results shown in Table 3 show the relationship between resolution, "blind" zone and reflection depth with reference to the antenna used. The simulated measurements are made in a media whose relative dielectric permittivity, $\varepsilon_r = 4.0$ and the specific attenuation is 1 to 2 dB/m. Reflection depth is the detection depth of a flat boundary with reflectance equal to 1.

Table 3. Frequency Dependence [10]

| Parameter | Antenna | | | | | | | | |
|---------------------|--------------|-------------|--------------|---------|------|-------|-------|--|--|
| Frequency (MHz) | 2 | 900 | 500 | 300 | 150 | 75 | 37 | | |
| Resolution (m) | 0.06- 0.1 | 0.2 | 0.5 | 1.0 | 1.0 | 2.0 | 4.0 | | |
| "Blind" zone (m) | 0.08 | 0.1- 0.2 | 0.25- 0.5 | 0.5-1.0 | 1.0 | 2.0 | 4.0 | | |
| Depth (m) | 1.5-2 | 3-5 | 7-10 | 10-15 | 7-10 | 10-15 | 15-30 | | |

Controlling Conductivity: The 1D-FDTD model allows one to investigate the effect of controlling one variable at a time. Figure 6 shows the results of a FDTD simulation where the specific conductance, σ , of the media is controlled and set to 5.0 Siemens/meter, the relative electrical permittivity, ε_r , set to 1.0, the frequency set to 2 GHz, and with a grid dimension of dx = 0.75 cm or 20 divisions per wavelength. The figure shows the attenuation of the fields in time.

Controlling Permittivity: Another example of system analysis by controlling one variable at a time is the constituent property of permittivity. Permittivity is a property that describes the ability of the media to store electric charge. It can also affect the frequency, wavelength, or amount of energy that is transmitted or reflected.



Figure 6. Example of Controlling Sigma.

The table presented within Figure 7 presents the relative permittivity of a number of common earth media. A graphic showing the boundaries and reflections from layers of different permittivity is also included. The reflection and transmission of the electromagnetic waves at each earth

media layer interface depends upon the difference of the permittivity of each layer. The signal received by the GPR receive antenna sub-system is a mixture of the reflection and delays propagating through the multi-layer paths. A representative profile for the different layers is presented.





Figure 7. Controlling Permittivity



Figure 8. Fresnel reflection and transmission coefficients

The reflection and transmission coefficients for two layers with relative permittivity's of $\varepsilon_{r1} = 2.0$ and $\varepsilon_{r2} = 4.0$ is shown in Figure 8. The reflection and transmission amplitude coefficients are shown for both perpendicular and parallel polarizations of EM waves incident from normal to 90 degrees. For $\varepsilon_{r1} = \varepsilon_{r2}$, there is total transmission and no reflection.

Figure 9 shows the ability of the 1D-FDTD simulation to model the effects of different values of permittivity on the propagation of electromagnetic waves. The specific conductance of the media is set to $\sigma = 0$ Siemens/m and the value of permittivity is controlled at $\varepsilon_r = 1.0$ and $\varepsilon_r = 10$. The media is nonmagnetic with permeability equal to free space, μ_0 . The simulation shows how ε_r affects both the frequency and the speed of Both graphs show 12 nano-seconds of propagation. propagation. The higher the ε_{r} , the slower the wave propagates. This delay gives insight into how deep a reflecting target might be if the ε_r is known or a method to determine the ε_r if the depth of the reflecting object is known.



Figure 9. Effect of permittivity on propagation.

This exercise allows one to understand GPR physical processes better by controlling variables that are modeled in the FDTD model. It demonstrates the ability of the model to perform a bistatic polarimetric simulation of the GPR. Using a simple FDTD model and simulation with perfectly matching boundary conditions, a FDTD simulation of rods at half a meter depth was performed. The homogeneous media show the expected result that polarized electromagnetic waves induce larger currents in the direction in which the wave and rod are oriented. Exposed to a polarized EM wave in the x direction, the x-directed rod has larger induced currents in the x-direction, while the ydirected rod has a strong tendency to induce currents in the ydirection if the EM wave is polarized in the y direction. This explains why GPR migration algorithms, developed on a matched-filter response basis, are used to both detect and determine the shape of a buried pipe like object.

Gurel et al presents an excellent example of prism modeling [12]. In Figure 10, the FDTD model simulates two conducting prisms of 21 x 21 x 16 cells that are buried five cells under the ground, and separated by twenty cells. The Ascan waveforms are calculated and presented next to B-scan results. In Figure 10, the scattering results for a cavity and a dielectric object, with a permittivity of $\varepsilon_r = 1.0$ and $\varepsilon_r = 8$, respectively, are presented. The two targets are buried twenty cells apart and five cells under the ground that is modeled with a relative permittivity of $\varepsilon_r = 4.0$. Figure 10 illustrates the typical A-scan and B-scans expected and the ability of the FDTD model to simulate the GPR performance. In Figure 11(a) the targets are dielectric object and a cavity in the ground. Note that the amount of reflection from the two objects closely follows the Fresnel reflection and transmission coefficients illustrated in Figure 7 for layers with the values of $\varepsilon_r = 4.0$ for the soil and $\varepsilon_r = 8$ for the dielectric object and $\varepsilon_r = 4.0$ for the soil and $\varepsilon_r = 1$ for the void. This results in the return from the cavity being larger than the return from the dielectric object. The results of this FDTD simulation are consistent with those using Fresnel reflection and transmission coefficients to calculate the reflection from the objects.

In the second simulation, the dielectric object is replaced by a conducting prism. The reflection from the perfectly conducting prism is nearly 100% and much larger than that of the cavity.



Figure 10. Two perfectly conducting prisms buried 5 cells under the ground and separated by 20 cells [12]



Figure 11. Two objects buried 5 cells under the ground and separated by 20 cells. (a) a cavity and a dielectric object and (b) a cavity and a perfectly conducting prism [12].

These TDFD models and simulations clearly show how the researcher can vary the media and targets buried in the media and systematically evaluate the GPR's performance. These experimental results yield the conclusion that the FDTD technique can be used to accurately simulate the GPR measurements and to faithfully analyze GPR data.

6 Conclusions and Future Plans

The initial results presented in this paper demonstrate the ability of the 3D FDTD method to model and simulate the effects of several media on the propagation of GPR signals. It is an important step in the system engineering analysis to identify GPR system deficiencies and what can be done to improve the system's performance. GPR computer models and FDTD simulations provide insight into how GPR systems including their signal processing algorithms perform to detect and identify objects buried under the ground. The TDFD model and simulations described in the paper allow the researcher to vary the media and targets buried in the media and systematically evaluate the GPR's performance. The effectiveness of the algorithms for data acquisition, signal processing and image processing for target detection and identification can be evaluated. Results from this modeling demonstrate the possibility of future use of this methodology for algorithm development and refinement that will better characterize and expand the trade space with valid GPR alternatives . The approach of simulating various input variables for an existing GPR using relatively simple 3D FDTD calculations has been demonstrated. The experimental results obtained lead to the conclusion that the FDTD techniques can be successfully used for analysis and parameter optimization of the basic signal processing algorithms in GPR.

Future planned research includes: accounting for the humidity and the inhomogeneity of soils on a GPR's performance to allow the development of robust highperformance detection algorithms. This includes the modeling of objects other than simple pipes and prisms such as multiple targets, dielectric targets in both homogeneous and anisotropic media. In the research to define appropriate solutions, FDTD has the computational ability to faithfully model a large variety of problem spaces. The propagation and detection of buried objects will be further investigated to obtain a better understanding of how the physical GPR components and processes affect the ability to detect and identify buried objects. Finally, the simulations will be expanded to the antenna-to-air and air-to-ground interfaces in order to better understand the interference paths of direct and ground bounced signals on the signals received from the reflections below the ground.

References.

[1]Montgomery, D. C. *Design and Analysis of Experiments*, 7th edition. Hoboken, NJ: John Wiley and Sons, 2009.

[2] Baker, G. S., Jordon, T. E., Pardy, J.. "An Introduction to ground penetrating radar (GPR)." *The Geological Society of America*, (Special Paper 432, 2007): 2

[3] Fritzsche, M. "Detection of buried landmines using g 3round penetrating radar." *Proceedings of the SPIE.* 2496 (1995): 100-109.

[4] Trang, A.H.. "Simulation of mine detection over dry soil, snow, ice, and water." *Proceedings of the SPIE*. 2765 (1996): 430-440.

[5] Williford, C.F. Comparison of Absorption and Radiation Boundary Conditions Using a Time-Domain Three-Dimensional Finite Difference Electromagnetic Computer Code. (1985): Air Force Institute of Technology, Wright-Patterson AFB, OH.

[6] Yee, K.S.. "Numerical Solution of Initial Boundary Value Problems Involving Maxwell's Equation in Isotropic Media," *IEEE Transactions on Antennas and Propagation*. AP-14 (May 1966): 302-307.

[7] Rymes, M.D. T3DFD User's Manual: Final Technical Report, August 1979--June 1981.

[8] Hebert, J.L. Sanchez-Castro, C. Implementation of a Three-Dimensional Finite Difference Electromagnetic Code for Analysis of Lightning Interaction with a FAA CV-580 Aircraft, Wright-Patterson AFB, OH. May 1987.

[9] Williford C., Jost, R.J., Hebert, J.L. "Comparison of Absorption and Radiation Boundary Conditions in 3DFD Code." *Proceedings of 1986 International Aerospace and Ground Conference on Lightning and Static Electricity.* (1986): 44-1--44-10.

[10] Belli, K., Rappaport, C., Udall, C., Hines, M., Wadia-Fascetti, S. "Use of 2D FDTD Simulation and the Determination of the GPR Travel Path Angle for Oblique B Scans of 2D Geometries." *Geoscience and Remote Sensing Symposium 2010 IEEE International.* (2010).

[11] "An Introduction to GPR" Global Scan Technologies. April 2012. http://www.gstdubai.com/Download/GPR.pdf.

[12] Gürel, L. "Simulations of Ground-Penetrating Radars Over Lossy and Heterogeneous Grounds" *IEEE Transactions* on *Geoscience and Remote Sensing*. 39.6. (2011): 1190-1197.
[13] Hakim, A. "The dual Yee-cell FDTD scheme." April 2012. http://ammar-hakim.org/sj/je/je7/je7-dual-yee.html#.

Simulation of the human lumbar spine during stooping and squatting using a forward kinematics model

K. Patel¹, A. Ghasempoor¹, and M. Abdoli-Eramaki²

¹Department of Mechanical Engineering, Ryerson University, Toronto, Ontario M5B 2K3, Canada ²School of Occupational and Public Health, Ryerson University, Toronto, Ontario M5B 2K3, Canada

Abstract - Working with human subjects to measure internal loadings of lumbar spine is difficult. Numerical modeling makes it possible to study response of spine to a variety of loading conditions. This paper reports the development of a rigid body model of spine using a forward kinematics The model consists of accurate approach. CAD representations of vertebrae interconnected by intervertebral discs providing six degrees of freedom at each joint. The vertebrae are articulated using Virtual Muscles controlled by PID controllers individually. The forward kinematics approach resolves the redundancy issues associated with other lumbar spine models which utilize inverse kinematics. The proposed model predicted similar trend for trunk rotations, compression forces and moments as described in the literature. The comparison between the proposed model and predictions from existing dynamic models showed good agreement.

Keywords: Simulation, Lumbar, Spine, Muscles

1 Introduction

Lumbar spine is a complex structure of the body. Since musculoskeletal disorders in lower back are prevalent in active and inactive occupations prediction of the dynamic response of spine under various loading conditions is an important consideration [1]. Both active and inactive work can increase the instability of spine to some degree leading to injuries and pain disorders of lower-back. Instability in spine could be due to the sudden perturbation [2], prolonged sitting (fatigue) [3], whole body vibration [4], awkward posture of spine [5], heavy lifting [6], unsupported sitting or standing [7] as well as larger applied forces [8].

Computational modeling of lumbar spine has been an active area of research for some time. Modeling of the spine using the computerized model is necessary for fully assessing reaction forces and motions in lumbar spine. Sensors cannot easily be introduced into spine of a human subject due to the risk of degeneration of the other discs [9]. Modeling also helps researchers in understanding the internal mechanisms of lumbar spine at extreme load and motion levels where physical testing is limited due to ethical reasons.

Majority of lumbar spine modeling can be categorized as finite element (FE) or rigid body modeling based. Finite element methods have been used by a number of researchers and show better results for small movements of spine [10]. It is more difficult to use FE based models for large movements and high levels of loading [11]. Asymmetric geometry of the spine in FE models produces asymmetric forces for right and left facet joint [12] which contributes to model instability and makes it difficult to model the behaviour when larger compressive load are present [13].

Rigid body models consider vertebrae as rigid parts, massless springs and dampers to model the behaviour of discs, joints and ligaments. Keller and Colloca [14] developed a model to analyze Posteroanterior (PA) motion response of lumbar spine. They reported good results for static and dynamic response of the lumbar spine using greatly simplified rigid body model when compared with in vivo measurements. The model was, however, limited to PA motion response in one single direction due to immobility of pelvis and thorax as well as the assumption of linear force-velocity relationship of the model. Kassem et al. [15] expanded this model to simulate the three dimensional motion responses. Their model was still limited to linear, static and dynamic mechanical responses and did not take in to account the geometry of spine. De Zee et al. [16] developed a detailed rigid body model using the AnyBody Modeling system using inverse dynamics to study muscle, ligament and reaction forces.

The models reviewed so far, whether FE based or rigid body based, use inverse kinematics approach to investigate static and dynamic motion responses as well as internal forces. Inverse kinematics is a technique in which the desired position and orientation is assumed to be known or is planned. Using this kinematics relationship, internal forces required for achieving this set of positions and orientations is calculated. Because of redundancy problem while acquiring the desired position and orientation using inverse kinematics,



Figure (1): Model of Lumbar spine including springs and dampers.

one could have multiple results or no results. The redundancy is caused by the fact that there are not enough equilibrium equations available to determine all the muscle forces [17]. In the case of multiple plausible results, optimization is typically used to pick the best set of answers. Using the inverse kinematics technique to find the motion responses and internal loads in spine, many researchers have encountered problems with results since optimality criteria are somewhat arbitrary. Another drawback of this technique is that it is not suitable for very high or very low loadings [18]. Studies have shown that inverse kinematics technique is also not reliable for intervertebral translation [19].

In this paper a forward kinematic model is proposed which eliminates the redundancy problem associated with inverse kinematics modeling of lumbar spine. The model mimics the natural process for activating muscles in that a high level command, for example pick a load, is processed through a series of controllers which activate the right muscles to the right level for performing the task. This eliminates the need for any assumptions and there will be only one solution for the final result. The paper focuses on muscles forces, internal loads on the intervertebral discs as well as motion response of the lumbar spine. The feasibility of this approach is demonstrated by comparison of results with data from



Figure (2): Lumbar spine model with front, back and side muscles attached

literature.

2 Lumbar spine model

Majority of the lumbar muscle tendons are attached to the vertebral bodies and muscle strength and forces mainly depend on muscle's length and physiologic cross section area (PCSA). In order to establish accurate location for attaching tendons, three-dimensional models of vertebrae were obtained (Zygote, Media Group, Inc, USA). In the proposed model



Figure (3): A subsystem of feedback controller block



Figure (4): Lifting activities: (a) Squat lift, (b) stoop lift with load of 0N and 180N

vertebrae are considered rigid bodies (Keller et al., 2002b). Figure 1 shows the scanned CAD model of lumbar vertebrae and pelvis. Each intervertebral disc was modeled with six degree of freedom (30 DOF in total). Non-linear springs and dampers were used to model the dynamic characteristics for each degree of freedom (Figure 1).

Ligaments wrap around vertebral bodies and act as shock absorbers and also restrict the excessive flexion, extension and shear of the vertebrae under excessive loads [20]. Ligaments were modelled using torsional springs and dampers similar to rigid body model develop by Keller et al,[1]. The spring rate of the rotational spring was taken to be 570 N.m/rad where as damping coefficient of the rotational damper was taken to be 150 N.m/ (rad/s) [1,14].

The model presented in this paper is a direct kinematics representation. As such, muscle forces are applied to vertebrae to produce the desired motion. "Virtual muscle" [21] was used to simulate the input-output relationship between muscle activation calculated by PID controllers, length and the corresponding muscle forces. Fifteen local muscles and two global muscles (one front and one back) are connected to lumbar spine (Figure 2). This number was selected as a compromise between the accuracy of model predictions and complexity of the model. An equivalent PCSA was calculated and used for these muscles so that the overall effect would be similar to original muscle forces. The lengths of the muscles were continually determined by tracing the attachment points from the CAD model. The local muscles were connected from transverse process of each vertebra to iliac crest and spinous process of each vertebra to the sacrum (10 and 5, respectively) [16]. The global front muscle was connected from 12th rib of the spine to the sacrum while global back muscle was connected from 8th rib of the spine to the sacrum. All local muscles and global front muscle were connected straight whereas global back muscle path was connected as via-point to reflect the realistic muscle path [16]. The physiological cross section area (PCSA) of the



Figure (5): Trunk Rotation for squat and stoop lift with load of 0 N and 180 N

global muscles, rectus abdominis (RA) and Longissimus thoracis pars lumborum (LGPL) were taken as 5.67 cm^2 and 12.10 cm^2 [18].

PID (proportional integrator and derivative) controllers were used to regulate the activation and therefore, the force being produced by each muscle. Figure 3 shows the structure of five sub systems for the feedback controller block, each with 'n' muscles. The forces produced by the muscles cause the lumbar spine to respond dynamically. This rotation is compared with the desired lumbar spine position and the difference between the current and desired position is calculated as error. These errors are fed in to the PID controller, which adjust the muscle activation accordingly until the desired lumbar spine position is achieved.

3 Case study: Stooping and Squatting

The main objective of this investigation was to develop and validate the forward kinematics lumbar spine model. To show the feasibility of the direct kinematic modeling, stoop and squat lifts were studied. Model predictions of internal loads were compared with dynamic model from literature [22]. To make comparison of internal loads in lumbar spine possible, the total body weight of 74 kg (mean value of 11 subjects) and external load of 180N was considered. The body weight of each segment was calculated using anthropometric data [23].

Compression forces, net external moments and passive ligamentous moments were directly obtained from the model, whereas the moments generated by muscle forces were calculated by taking the difference between the net external moment and passive ligamentous moments.

To examine dynamic loading situations, Stoop and squat lifts were carried out without load and with a load of 180N. Simulation was performed for five seconds and four types of data were obtained: trunk rotation, compression forces and moments. In both cases of stooping with and without load, hip, knee and ankle joints were stabilized during the lift while



Figure (6): Compression force on discs between L5-S1 for squat and stoop lift with load of 0N and 180N

pelvis was rotated to 45 degrees in order to allow the lumbar spine to rotate to its maximum amount. The allowable rotation for each joint was around five degrees. During simulation, relative rotation between L1 and L2 was slightly more than five degrees whereas the other four joints rotated less than five degrees. The upper body was considered rigid. In squatting, hip was rotated 30 degrees, knees 90 degrees and ankles 60 degrees. The pelvis rotation was limited to 26.3 degrees in squatting in forward direction (Figure 4). Trunk rotation for the squatting and stooping with load of 0N and 180N is shown in Figure 5. In stooping, the figure shows the maximum trunk rotation of 70 degrees at 2.5 seconds with or without load. In squatting, the trunk rotation for lifting the load is higher than the one with no load.

For stooping movement, all five lumbar discs were allowed to rotate a maximum of five degrees, allowing the whole trunk to be fully extended. After lifting a load of 180N, the trunk rotated only a few additional degrees. For both stooping and squatting after picking up the load, an upward shift was observed compared with no load situation. These are due to the gravitational force applied on the body which requires more muscle force to pull the body from the fully extended position to standing straight.

The effect of these two movement patterns is clearly indicated in intradiscal loadings shown in Figure 6. The maximum compression forces were recorded on disc L5/S1 for all activities. The maximum compression force on the same disc for squat and stoop lift with load of 180 N was estimated to be 4119 N and 4909 N, respectively. Compression forces for all discs are shown in the table 1 for both stoop and squat lift and for both loading conditions (0N and180N).

When lifting a load, the muscle forces are adjusted by nervous system through varying the simulation frequency or changing the number of motor units that are activated at a given time [23]. This process may take up to one third to half a second depending on the load. In present case study, time duration of 0.2 second was used to lift the load of 180N. A



Figure (7): Net external moment at the L5-S1 level for stoop and squat lift without load and with load of 180N

similar rate has been used in the work of Bazrgari et al. [22]. Also the compression forces for both stoop and squat lifts without loads were higher for later half of the cycle (from forward flexion to upright position) compared to first half of the cycle (from upright position to forward flexion). This unsymmetrical loading result from gravitational force applied to trunk from forward flexion

For zero external loads, the maximum passive ligamentous moments at L1-L2 was recorded as 66Nm and 23Nm for stoop and squat lift, respectively. For the same disc, moments for stoop and squat lifts with load of 180N were recorded as 86Nm and 46Nm, respectively. These occurred at the time of maximum trunk rotation. The maximum external moments were measured on the disc between L5-S1, where stoop lift showed greater moment compared to squat lift. Same trend was observed even when carrying the load of 180 N for both movements (stoop and squat lifts). The net external moments were higher by ~40% and ~26% in stoop lift as compared to squat lift while carrying the load of 0N and 180 N, respectively. Similar results were obtained for passive ligamentous moments where the moments were higher by ~52% and ~34% in stoop lift compared to squat lift for no load and 180N external load. Furthermore, the muscle moments were higher for stoop after lifting the load of 180 N (higher by 3%) compared to squat lift. This indicates that the moments were primarily carried by passive ligamentous spine but after applying the external loads, some portions of the moments were resisted by the muscles. Figures 7-9 show the predictions of the external moments, muscle moments, and passive ligamentous moments. The passive ligamentous and muscle moments were higher in stoop lift due to the larger intersegmental rotation of the lumbar spine. In both stoop and squat lift cases with load in hands, the passive ligamentous, muscle and external moments were increased abruptly as the load reached its maximum value of 180N. After reaching its peak point, external moments and muscle moments curve decrease gradually while the passive ligamentous moments decrease abruptly.

| Lumbar discs | Squat (0 N) | | Stoop (0 N) | | Squat (180 N) | | Stoop (180 N) | |
|--------------|----------------|----|----------------|----|------------------|----|------------------|----|
| | C | М | C | М | С | М | С | М |
| L1/L2 | 1088 | 23 | 1627 | 66 | 2003 | 46 | 2150 | 86 |
| L2/L3 | 1326 | 13 | 2074 | 36 | 2404 | 23 | 2693 | 48 |
| L3/L4 | 1517 | 10 | 2587 | 22 | 2870 | 18 | 3402 | 31 |
| L4/L5 | 1820 | 12 | 2953 | 23 | 3325 | 20 | 3889 | 31 |
| L5/S1 | 2292 | 20 | 3602 | 42 | 4119 | 38 | 4909 | 56 |

Table (1): Maximum internal loads measured in spine for different cases at various levels; Compression force, C (N) and Passive ligamentous moment, M (N.m)

4 Discussion

Compression forces and moments on the disc between L5-S1 and trunk rotation for stoop and squat lifts with load of 180N and 0N were compared with Bazrgari et al.[22]. The trunk rotations were found well within the range except for squat lift with ON loads as well as some minor differences in stoop lifts. Rotation obtained for the trunk was lower compared to the trunk rotation for squat lift with load of 0N from Bazrgari et al. The difference in trunk rotation for squat lift with ON was mainly due to the lumbar rotation. The proposed model estimated ~50% lower lumbar rotation while the pelvis rotations were predicted approximately the same as predicted by Bazrgari et al. The trunk rotation in the squat lift was increased by approximately five degrees after lifting the load which was not seen in the literature. In stoop lifts and squat lifts the thorax and pelvis rotations with load and without load for first half cycle were overlapping in proposed model whereas in Bazrgari et al. the measurements of the thorax and pelvis rotations were shifted toward left after lifting the load. The main reasons for the shift in curves (Fig. 6) were due to the lack of instructions in the manner of lifting the load for subjects and the fact that the measurements were based on the mean values of 11 subjects.



Figure (8): Portion of the moment resisted by the muscles on disc L5-S1 during lifting

Compared to Bazrgari et al. the compression force on the disc L5/S1 was found higher in all activities for 0N and 180N, whereas these forces on other discs were lower except for stoop lift with 0N (Table 1). The maximum passive ligamentous moments in the proposed model were measured on disc L1-L2 and were found higher compared to literature for all activities whereas the net external moments and moments resisted by the muscles were measured higher on disc L5-S1 but were found to be lower compared to Bazrgari et al.

Both the proposed model and Bazrgari et al. showed sudden increase in compression forces and moments after lifting the load in both activities (stoop lifts and squat lifts). In the Bazrgari et al. the compression forces and moments increased sharply to their maximum values and decreased suddenly while in the proposed model, the compression forces increased sharply to a point and then gradually reached their maximum value. In the proposed model, results for moments also increased sharply and reached its maximum point but decreased gradually. For net external moments, the maximum difference was calculated to be 15% or lower for all activities in current model. The compression forces and passive ligamentous moments are shown in the Table 1. The difference in compression forces and moments could be attributed to differences in lumbar rotation at each segment as well as shoulder and forearm rotation.

5 Conclusion

The model presented in this paper can be used to determine the muscle forces, compression forces, moments and motion responses of the lumbar spine using the forward kinematics technique which eliminates the need of optimization based on loading assumptions. The model can easily be modified to evaluate different postures such as standing, lifting, sitting or bending in either direction. Muscles play a vital role in stabilizing the spine so it will be appropriate to include all muscles in a complete model of spine. The proposed model included fewer numbers of lumbar muscles with larger physiological cross section areas to reduce the computational load without reducing the required muscle force to perform a particular activity. The results for stoop and squat lifts were



Figure (9): Passive ligamentous moment during squat and stoop lifts with load of 0N and 180N

found in good agreements in comparison with other dynamic model predictions presented in the literature. The model showed that adopting squat lifting over stoop lifting can reduce muscle forces and internal loads; and increase spine stability. The forces and moments were larger in stoop lifts compared to squat lift mainly due to the larger trunk rotation not because of larger lever arm.

6 **References**

[1] Keller, T.S., Colloca, C.J., Béliveau, J., 2002b. Forcedeformation response of the lumbar spine: A sagittal plane model of posteroanterior manipulation and mobilization. Clinical Biomechanics 17, 185-196.

[2] Bazrgari, B., Shirazi-Adl, A., Parnianpour, M., 2009. Transient analysis of trunk response in sudden release loading using kinematics-driven finite element model. Clinical Biomechanics 24, 341-347.

[3] Søndergaard, K.H.E., Olesen, C.G., Søndergaard, E.K., de Zee, M., Madeleine, P., 2010. The variability and complexity of sitting postural control are associated with discomfort. Journal of Biomechanics 43, 1997-2001.

[4] Pankoke, S., Buck, B., Woelfel, H.P., 1998. Dynamic FE model of sitting man adjustable to body height, body mass and posture used for calculating internal forces in the lumbar vertebral disks. Journal of Sound and Vibration 215(4), 827-839.

[5] Hess, J.A., Kincl, L.D., Davis, K., 2010. The impact of drywall handling tools on the low back. Applied Ergonomics 41, 305-312.

[6] Wai, E.K., Roffey, D.M., Bishop, P., Kwon, B.K., Dagenais, S., 2010. Causal assessment of occupational lifting and low back pain. Results of a systematic review. Spine Journal 10, 554-566.

[7] Damecour, C., Abdoli-Eramaki, M., Ghasempoor, A., Neumann, W.P., 2010. Comparison of two heights for

forward-placed trunk support with standing work. Applied Ergonomics 41, 536-541.

[8] Harvey, J., Tanner, S., 1991. Low back pain in young athletes. A practical approach, Sports Medicine 12, 394-406.

[9] Wilke, H., Neef, P., Hinz, B., Seidel, H., Claes, L., 2001. Intradiscal pressure together with anthropometric data - a data set for the validation of models. Clinical Biomechanics 16, S111-S126.

[10] Arjmand, N., Shirazi-Adl, A., Parnianpour, M., 2008b. Trunk biomechanics during maximum isometric axial torque exertions in upright standing. Clinical Biomechanics 23, 969-978.

[11] Arjmand, N., Shirazi-Adl, A., Parnianpour, M., 2008a. Relative efficiency of abdominal muscles in spine stability. Computer Methods in Biomechanics and Biomedical Engineering 3, 291-299.

[12] Kuo, C., Hu, H., Lin, R., Huang, K., Lin, P., Zhong, Z., Hseih, M., 2010. Biomechanical analysis of the lumbar spine on facet joint force and intradiscal pressure – a finite element study. BMC Musculoskeletal Disorders 11, 1-13.

[13] Shirazi-Adl, A., Parnianpour, M., 2000. Load-bearing and stress analysis of the human spine under a novel wrapping compression loading. Clinical Biomechanics 15, 718-725.

[14] Keller, T.S., Colloca, C.J., 2002a. A rigid body model of the dynamic posteroanterior motion response of the human lumbar spine. Journal of Manipulative and Physiological Therapeutics 25, 485-496.

[15] Kassem, A.H., Sameh, A., Keller, T.S., 2004. Modeling and simulation of lumbar spine dynamics. Fifteenth IASTED International Conference on Modeling and Simulation, Marina Del Ray, California, USA, March, 1-3.

[16] De Zee, M., Hansen, L., Wong, C., Rasmussen, J., Simonsen, E.B., 2007. A generic detailed rigid-body lumbar spine model. Journal of Biomechanics 40, 1219-1227.

[17] Rasmussen, J., Tørholm, S., de Zee, M., 2009. Computational analysis of the influence of seat pan inclination and friction on muscle activity and spinal joint forces. International Journal of Industrial Ergonomics, 39, 52-57.

[18] Arjmand, N., Gagnon, D., Plamondon, A., Shirazi-Adl, A., Larivière, C., 2009. Comparison of trunk muscle forces and spinal loads estimated by two biomechanical models. Clinical Biomechanics 24, 533-541.

[19] Sun, L.W., Lee, R.Y.W., Lu, W., Luk, K.D.K., 2004. Modelling and simulation of the intervertebral movements of the lumbar spine using an inverse kinematic algorithm. Medical and Biological Engineering and Computing 42, 740-746.

[20] Adams, M.A., Bogduk, N., Burton, K., Dolan, P., 2002. The Biomechanics of Back Pain. London, U.K., Churchill Livingstone an imprint of Elsevier Science Ltd.

[21] Song, D., Cheng, E., Brown, I., Davoodi, R., Loeb, G.E., 2008. Virtual Muscle 4.0.1: Muscle for MATLAB.

[22] Bazrgari, B., Shirazi-Adl, A., Arjmand, N., 2007. Analysis of squat and stoop dynamic liftings: Muscle forces and internal spinal loads. European Spine Journal 16, 687-699.

[23] Winter, D.A., 2005. Biomechanics and Motor Control of Human Movement. John Wiley & Sons, Hoboken, New Jersey.

High Performance Monte Carlo and Time-Stepping Dynamics for the Classical Spin Heisenberg Model on GPUs

K.A. Hawick and D.P. Playne

Computer Science, Institute for Information and Mathematical Sciences, Massey University, North Shore 102-904, Auckland, New Zealand email: { k.a.hawick, d.p.playne }@massey.ac.nz Tel: +64 9 414 0800 Fax: +64 9 441 8181

April 2012

ABSTRACT

The Heisenberg model of classical spins makes use of both Monte Carlo stochastic dynamics as well as time-integration of its equation of motion. These two schemes have different parallelisation strategies and tradeoffs. We implement both algorithms using a data-parallel approach for Graphical Processing Units (GPUs) and we discuss the resulting performance on various combinations of single and multiple GPU. In addition to studying Monte Carlo dynamical update schemes, we use our fast simulation code to explore the scaling and time correlations of a largescale Heisenberg model system using a high-order numerical integration algorithm, which enables study of accurate spin wave phenomena and time-correlation functions. We also discuss various graphical rendering models to appropriately visualise the spin vectors inside an interactive Heisenberg spin simulation.

KEY WORDS

Heisenberg model; classical spin; Monte Carlo dynamics; time-integration dynamics.

1 Introduction

Simulation of complex systems is a powerful means of investigating phase transitions [6] and critical phenomena [16]. Visualising the approach to criticality of such system is also important to help develop an intuitive understanding of simulation model parameters. A high-performance visual simulation also aids in navigating through the model parameter space to identify interactively those areas that are worth more exhaustive simulation and collection of statistical measurements from appropriate numerical experiments.



Figure 1: A visualization of a sample 3D Heisenberg simulation.

A great deal of work has been done on models such as the Ising model [10, 15] which is based upon applying a stochastic Monte Carlo based dynamics on a system of spins modelled by individual bits. The Heisenberg model [1] of classical spin system [13, 24] is interesting because its more realistic continuous individual spins can be simulated dynamical using the Monte Carlo method [3] but also using a more realistic timeintegration method based on proper equations of motion. This combined approach means that the Heisenberg model is more appropriate for studying dynamic growth, decay and relaxation properties [19], since the simulation time can be more readily identified as a "real" time variable rather purely as an artifact of the Monte Carlo algorithm.

In this present paper we make use of Graphical Processing Units (GPUs) and NVidia's Compute Unified Device Architecture (CUDA) software language and framework to develop very fast numerical simulations of large Heisenberg spin model systems in two and three dimensions [22]. We also develop various graphical rendering models using OpenGL software to visualise and explain the dynamically evolving spin vector field as both Monte Carlo dynamics and time integration dynamics are applied. Figure 1 shows a sample rendering of a three dimensional Heisenberg system, simulated on a 64^3 lattice.

Although other workers have reported simulation work of Monte Carlo simulations of the Heisenberg model [5, 26], it is also important to be able to simulate its time-integration dynamics [2]. A typical numerical experiment consists of rapidly quenching an hot random initial spin pattern using the Monte Carlo dynamics, followed by a carefully controlled fine grained time-stepping on the system using a suitable numerical integration scheme to obtain temporal measurements.

We have developed automatic code generation techniques that has allowed us to generate high order time integration software to solve the equations of motion for the spins to tenth order accuracy. This high degree of accuracy and associated stability allows us to investigate key measured properties such as time correlation functions over longer periods of simulated time than would otherwise be feasible. These measurements then have the potential to be compared with quantities obtained from experiments on real physical magnetic systems [18].

Our article is structured as follows: In Section 2 we summarise the (classical) Heisenberg model of spins and the dynamical schemes that we can apply to it. In Section 4 we present some visual renderings of the spin system as well as a discussion of some computational performance measurements of our software running on various individual and multiple GPU systems. In Section 5 we give us discussion of some of the phenomena we observed and we offer some conclusions and areas for further work in Section 6.

2 Heisenberg Model Simulations

The classical Heisenberg model is essentially a continuous spin version of the Ising model.

The Heisenberg system is realised using a ddimensional array (such as a cubic lattice in 3-D or square lattice in 2-D) where each lattice site has a spin on it. A spin comprises a vector (s_x, s_y, s_z) where each component is a scalar, normalised so that $s_i \in [0, 1]$. In practice, the normalisation of the spin vectors to be unit vectors implies that |s| = $\sqrt{(s_x^2 + s_y^2 + s_z^2)} \equiv 1$ and this means there are effectively only two degrees of freedom. Using spherical trigonometry, the unit spin vector can thus be represented by two angles $\theta \in [0, \pi], \phi \in [0, 2\pi]$.

The energy function (the Hamiltonian) of the classical Heisenberg model system is:

$$\mathcal{H} = -\mathcal{J} \sum_{i,j} \mathbf{s_i} \cdot \mathbf{s_j} \tag{1}$$

Where the summation is over the nearest neighbours of the lattice site and the negative sign (with a positive value of \mathcal{J} means we get ferro-magnetism - ie alignment of the spins for strong couplings (low temperatures). The dot product is between two neighbouring spins and essentially contributes towards the total energy when two spins couple or closely align in direction with one another.

We can write the equation of motion for the spins in the form of a differential equation as:

$$\frac{d\mathbf{s_i}}{dt} = \mathbf{s_i} \times \sum_{\mathbf{j}} \mathbf{J}\mathbf{s_j}$$
(2)

Where $\mathbf{A} \times \mathbf{B}$ is the cross product of the two vectors \mathbf{A} and \mathbf{B} , and so the differential equation is actually a vector equation – with a separate component for each of the x, y, z parts of $d\mathbf{s_i}$.

This can be transformed into a finite difference equation using the standard techniques such as Euler (poor stability) or Runge-Kutta (a lot better stability) and we in fact use a tenth order algorithm due to Hairer [8], that has a very high degree of accuracy and stability. We thus obtain an explicit formula for the change Δs in each spin in terms of its prior value and the values of its nearest neighbours.

$$\mathbf{s}_{\text{new}} = \mathbf{s}_{\text{old}} + \Delta \mathbf{s}$$
 (3)

This algorithm allows us to update the Heisenberg spins very carefully and with a realistic and meaningful time that can be compared with temporal measurements of real magnetic systems. However in practice the numerical experiments that are typically performed on the Heisenberg system involve quenching a hot initially random arrangement of spins to a finite temperature. Time integration is not well suited to this quenching process as it is too slow. In practice a Monte Carlo stochastic algorithm is used to step the system forward in pseudo-time to thermal equilibrium, and the time-stepping algorithm can then be applied. Monte Carlo thermalisation using one of the standard algorithms such as Metropolis [17] of Glauber [7] are described in detail elsewhere [12]. In summary, these algorithms work as follows. At each Monte Carlo step each spin is considered in turn (usually in random order). A new direction for the spin is generated randomly and the energy consequences ΔE of this change computed. If the spin change would decrease the energy then the spin change is immediately accepted. Otherwise the change is accepted according to the Boltzmann probability factor $\exp(-\Delta E/k_B T)$ where k_B is Boltzmann's constant - which we take to be unity for our purposes, and T is the temperature, which is effectively just the reciprocal of the coupling parameter J.

3 GPU Implementations

Graphics Processing Units (GPUs) have been shown to be a very effective processing architecture for regular lattice simulations such as the Heisenberg model. Originally designed for rendering real-time graphics for computer games, GPUs have evolved into highly parallel architectures and are being increasingly used for scientific applications. In previous work, GPUs have been used for processing the Ising spin model [11] as well as scalar and vector [21] models.

All simulations discussed in this paper have been executed on Fermi architecture NVIDIA GPUs. Fermi GPUs contain up to 16 multiprocessors which each contain 32 scalar processors or SPs. Each multiprocessor contains on-chip memory which allows information to be shared between SPs, all the multiprocessors also have access to global memory which is the main storage area of the device.

These simulations perform all memory access through global memory which is automatically cached on Fermi devices. This memory access has been shown to provide the best performance for this access pattern [21]. For more details on GPU architectures and implementing lattice-based simulations on GPUs see [20, 21].

The implementation of the Heisenberg model is different to previous work in that each spin is represented by a three-dimensional vector and requires two phases of computation - the equilibration phase and the spin update phase. The equilibration phase is computed using a Monte-Carlo method while the spin update is computed by integrating the Heisenberg equation of motion over time. Each of these phases must be parallelised for them to be computed on a GPU.

3.1 Equilibration Phase

The equilibration phase of the Heisenberg simulation requires the use of a parallel Monte-Carlo. The Metropolis algorithm does not parallelise well as it can lead to race-conditions so instead the checkerboard or red-black update is used. The checkerboard update pattern ensures that no two neighboring lattice cells are changed during the same update, meaning no race-conditions can occur. This red-black checkerboard pattern is shown in Figure 2.

| 0 | 1 | 2 | 3 | 4 | 5 | 6 | 7 |
|----|----|----|----|----|----|----|----|
| 8 | 9 | 10 | 11 | 12 | 13 | 14 | 15 |
| 16 | 17 | 18 | 19 | 20 | 21 | 22 | 23 |
| 24 | 25 | 26 | 27 | 28 | 29 | 30 | 31 |
| 32 | 33 | 34 | 35 | 36 | 37 | 38 | 39 |
| 40 | 41 | 42 | 43 | 44 | 45 | 46 | 47 |
| 48 | 49 | 50 | 51 | 52 | 53 | 54 | 55 |
| 56 | 57 | 58 | 59 | 60 | 61 | 62 | 63 |

Figure 2: The checkerboard update pattern.

Each update will read the value of a cell and it's nearest neighbors and compute the energy of the configuration (*E*1). It will then randomly generate another spin and compute the energy of this alternative configuration (*E*2). The change in energy can then be calculated ($\Delta E = E2 - E1$) and used to either accept or reject the new configuration. The new configuration is accepted if $\Delta E < 0$ or with probability $e^{-(\Delta E)/k_bT}$.

This update process requires the use of a random number generator (RNG) to create a new random spin configuration and also to determine if the spin should be accepted when $\Delta E > 0$. For this simulation the RAN random number generator discussed in Numerical Recipes [23] is used. This is a fast random number generator with relatively low storage requirements and has been shown to pass appropriate statistical tests [14]. The random number generation process is parallelised by creating a different RAN RNG for each lattice cell and appropriately initialized. This way each lattice cell has an independent stream of random numbers.

The checkerboard update pattern ensures that no two neighboring lattice cells are changed during the

same update, unfortunately GPUs give the best performance when sequential threads access sequential memory addresses because these memory transactions can be coalesced into a single transaction. This difference between algorithmic requirements and optimal GPU access patterns can be overcome by reordering or crinkling the lattice. The lattice can be rearranged so that cells that are updated at the same time are stored sequentially which allows the memory access to be as efficient as possible. This process is thoroughly described in [9] and the lattice from Figure 2 is shown in it's crinkled form in Figure 3.

| 0 | 2 | 4 | 6 | 1 | 3 | 5 | 7 |
|----|----|----|----|----|----|----|----|
| 9 | 11 | 13 | 15 | 8 | 10 | 12 | 14 |
| 16 | 18 | 20 | 22 | 17 | 19 | 21 | 23 |
| 25 | 27 | 29 | 31 | 24 | 26 | 28 | 30 |
| 32 | 34 | 36 | 38 | 33 | 35 | 37 | 39 |
| 41 | 43 | 45 | 47 | 40 | 42 | 44 | 46 |
| 48 | 50 | 52 | 54 | 49 | 51 | 53 | 55 |
| 57 | 59 | 61 | 63 | 56 | 58 | 60 | 62 |

Figure 3: The checkerboard update pattern shown on a crinkled lattice.

3.2 Spin Update Phase

In the spin update phase of the simulation, the equation of motion for each spin is computed and integrated over time to compute the spin at the end of the time step. Every spin is updated each time step which means that the best memory access pattern is provided by the uncrinkled lattice. Also the spin update phase does not require a random number generator.

There are a number of different methods that can be used to integrate the equation of motion (equation 2) over time. In this simulation (as with previous work) the explicit methods from the Runge-Kutta family of integration methods are used. These methods are used because they parallelise well and the higher-order methods can provide very good stability and accuracy. The higher order methods become increasingly complex to implement and for this reason we make use of code-generation techniques to produce template code. The code generator can create integration code for lattice-based simulations from a Butcher tableau [4,21].

These methods all integrate the motion of the spins

by calculating a number of intermediate stages. The derivatives of these intermediate stages are combined to calculate the final spin configuration. Higher order methods are more expensive in terms of memory storage and computational intensity. However, these more stable higher-order methods can often simulate systems with larger time-step and lead to an overall reduction in computation time.

4 **Results**

One important feature of the Heisenberg model is the presence of a phase transition at the critical temperature T_c which can be seen in Figures 6 and 7. These figures show the temperature dependent behavior of three Heisenberg systems in two-dimensions (Figure 6) and three-dimensions (Figure 7).

The first system has a temperature higher than the critical temperature $(T > T_c)$ and exhibits random 'hot' behavior. As the probability of accepting a new random spin approaches 1 the system will become completely random. The temperature of second system is near the critical temperature $(T \approx T_c)$ and shows quite different behavior to the first system. There are clear structures forming in this system while maintaining an element of randomness. The final system has a temperature well below the critical temperature $(T < T_c)$ and will simply minimize the energy of the system.

The effect of this temperature dependent behavior can be easily seen by plotting the energy of the system. Figure 4 shows the energy of 1024x1024Heisenberg systems averaged over 20 runs with T={0.1, 0.2...10}. The checkerboard update is used for 1024 steps until the system reaches equilibrium and then the Runge-Kutta 4th Order integration method is used to integrate the motion of the spins over time. It can be easily seen from the plot that the energy for the different systems quickly reaches equilibrium and remains stable at this value.

The temperature of the system also affects how quickly the correlation of subsequent systems decays. The correlation of a system at time t to a previous state at time 0 is calculated as the sum of the dot products of each current spin $s_{t,i}$ and the initial spin $s_{0,i}$. This can be written as:

$$C_t = \sum_i s_{0,i} \cdot s_{t,i} \tag{4}$$

The correlation of the two-dimensional Heisenberg model in the spin update phase for T =



Figure 4: The energy of the two-dimensional Heisenberg model for $T = \{0.1, 0.2...1.0\}$ vs time (ln scale).

 $\{0.1, 0, 2...1.0\}$ is shown in Figure 5. The Runge-Kutta 4^{th} Order integration method with a time step of h = 0.01 is used to evolve the system. The correlation of the system evolution shows different behavior depending on the regime of the temperature.



Figure 5: The correlation of the two-dimensional Heisenberg model for $T = \{0.1, 0.2...1.0\}$.

5 Discussion

We found that the time -integration method is computationally too inefficient to achieve a thermally equilibrated system in a reasonable number of time steps. The Monte Carlo dynamics can be implemented to rapidly move the quenched system to thermal equilib-



Figure 6: A series of three two-dimensional Heisenberg systems - below, near and above the critical temperature.



Figure 7: A series of three three-dimensional Heisenberg systems - below, near and above the critical temperature.

rium, but otherwise it does not yield good time measurements that can be easily related to real time behaviour in real magnetic system. Development of a code that can apply both dynamic schemes was therefore necessary.

We experimented extensively with different rendering schemes. The scheme used for the illustrations in this paper are based on a simple mapping of colour hue and value to the two degrees of freedom - namely the two angles ϕ , θ that precisely define each unit spin vector.

Another model - the clock model [25] – is similar to the classical Heisenberg system, except it has only two components in each spin vector. Consequently the clock model spins can be specified using only one angle and it could be visualised using an arrow or some simpler colour mapping.

There is scope for further experimentation with different rendering schemes. Some thresholding of the spin values by direction, could be used to look at a partial subset of the spins in a 3D hyperbrick at once. The problem of rendering a four dimensional system effectively is an open one. Although real magnetic systems are only 3 dimensional simulating a fourdimensional model is useful since it allows the dimensional dependence for various structural and phase transitional properties to be studied. This remains an open problem for the present however.

6 Conclusions

We have described the classical Heisenberg model of spin systems and our implementation of it on twodimensional and three dimensional lattice systems using Graphical Processing Units. Using appropriate memory mappings and data structures we obtained a sufficiently fast implementation of the Heisenberg simulations that we are able to explore its properties in near interactive time.

The Heisenberg system is somewhat more difficult to render than the Ising model. We experimented with various graphical renderings of the spin system using colour hue and value to map to the two independent degrees of freedom of the unit spin vectors of the evolving system.

We implemented both Monte Carlo dynamics as well as a high order time integration dynamical scheme. We were able to use these to quench the system from an initial random state and subsequently integrate it carefully to obtain time correlation functions respectively. There is scope for a more detailed study of spin wave phenomena using this simulation and metrical analysis. We also expect to be able to adapt our simulation to study damaged and frustrated spin models which have a direct bearing on comparisons with properties of new magnetic materials.

References

- Anderson, P.W.: New approach to the theory of superexchange interactions. Phys. Rev. 115, 2–13 (1959)
- [2] Bernaschi, M., Parisi, G., Parisi, L.: Benchmarking gpu and cpu codes for heisenberg spin glass overrelaxation. Computer Physics Communications 182, 1265–1271 (2011)
- [3] Binder, K. (ed.): Monte Carlo Methods in Statistical Physics. Topics in Current Physics, Springer-Verlag, 2 edn. (1986), number 7
- [4] Butcher, J.C.: Numerical Methods for Ordinary Differential Equations. No. ISBn 978-0-470-72335-7, Wiley, second edition edn. (2008)
- [5] Campos, A.M., Pecanha, J.P., Pampanelli, P., de Almeida, R.B., Lobosco, M., Vieira, M.B., de O. Dantas, S.: Parallel implementation of the heisenberg model using monte carlo on gpgpu. In: Proc. Computational Science and its Applications (ICCSA 2011). vol. LNCS 6784, pp. 654–667 (2011)
- [6] E.N.Miranda, N.Parga: Dynamical phase transitions in the classical heisenberg model. J.Phys.A: Math. Gen 24, 1059–1064 (1991)
- [7] Glauber, R.: Time dependent statistics of the ising model. J. Math. Phys II 228(4), 294–307 (1963)
- [8] Hairer, E.: A Runge-Kutta Method of Order 10. J. Inst. Maths. Applics. 21, 47–59 (1978)
- [9] Hawick, K.A., Playne, D.P.: Hypercubic Storage Layout and Transforms in Arbitrary Dimensions using GPUs and CUDA. Concurrency and Computation: Practice and Experience 23(10), 1027–1050 (July 2011)
- [10] Hawick, K., Leist, A., Playne, D.: Cluster and fast update lattice simulations using graphical processing units. Tech. Rep. CSTN-104, Computer Science, Massey University (November 2009)
- [11] Hawick, K., Leist, A., Playne, D.: Regular Lattice and Small-World Spin Model Simulations using CUDA and GPUs. Int. J. Parallel Prog. 39(CSTN-093), 183– 201 (2011)
- [12] Hawick, K.A.: Domain Growth in Alloys. Ph.D. thesis, Edinburgh University (1991)
- [13] Joyce, G.S.: Classical heisenberg model. Physical Review 155, 478–491 (1967)
- [14] L'Ecuyer, P.: Software for uniform random number generation: distinguishing the good and the bad. In: Proc. 2001 Winter Simulation Conference. vol. 2, pp. 95–105 (2001)
- [15] Leist, A., Playne, D., Hawick, K.: Interactive visualisation of spins and clusters in regular and small-

world Ising models with CUDA on GPUs. Journal of Computational Science 1, 33-40 (2010), www.elsevier.com/locate/jocs

- [16] M.E.Fisher: The theory of equilibrium critical phenomena. Rep.Prog.Phys. 30, 615–730 (1967)
- [17] Metropolis, N., Rosenbluth, A.W., Rosenbluth, M.N., Teller, A.H., Teller, E.: Equation of state calculations by fast computing machines. J. Chem. Phys. 6(21), 1087–1092 (1953)
- [18] M.Steiner, J.Villain, C.G.Windsor: Theoretical and experimental studies on one-dimensional magnetic systems. Advances in Physics 25(2), 87–209 (1976)
- [19] Muller-Krumbhaar, H., Binder, K.: Dynamic properties of the Monte Carlo Method in Statistical Physics. J. Stat. Phys. 8(1), 1–24 (1973)
- [20] NVIDIA® Corporation: CUDATM 3.1 Programming Guide (2010), http://www.nvidia.com/, last accessed August 2010
- [21] Playne, D.P.: Generative Programming Methods for Parallel Partial Differential Field Equation Solvers. Ph.D. thesis, Computer Science, Massey University (2011)
- [22] P.Peczak, Ferrenberg, A.M., D.P.Landau: Highaccuracy Monte Carlo study of the three-dimensional classical Heisenberg ferromagnet. Phys.Rev.B 43(7), 6087–6093 (Mar 1991)
- [23] Press, W.H., Teukolsky, S.A., Vetterling, W.T., Flannery, B.P.: Numerical Recipes The Art of Scientific Computing. Cambridge, third edn. (2007), iSBN 978-0-521-88407-5
- [24] R.E.Watson, M.Blume, G.H.Vineyard: Classical Hesienberg magnet in two dimensions. Phys.Rev.B 2(3), 684–690 (1970)
- [25] Stanley, H.E.: Introduction to phase transitions and critical phenomena. Oxford Science Publications (1987)
- [26] Weigul, M., Yavorskii, T.: Gpu accelerated monte carlo simulations of lattice spin models. In: Proceedings of the 24th Workshop on Computer Simulation Studies in Condensed Matter Physics (CSP2011). Physics Procedia, vol. 15, pp. 92–96 (2011)

Motion Simulation of the Modular Walking Robot MERO using Force and Attitude Sensors

¹Ion ION, ²Grigore STAMATESCU

¹Technology of Manufacturing Department, ²Automatic Control and System Engineering Department, "POLITEHNICA" University of Bucharest Romania

Abstract. A large number of vehicle models have been developed for their mobility characteristics over irregular surfaces because, more and more applications requiring movement on a natural, unarranged terrain, made the feetmovement solution become more and more attractive The modular mechatronic mobile system MERO* (MEchanism Robot-*Pelecudi Ch et.al.) by reconfiguring their architecture can carry the heavy loads on the irregular surfaces. Modern methods of drawing up machines and robots necessarily include a simulation stage of their functioning. These activate the functioning simulation that encompasses several rules and specifications whose enactment generates behavior data and the instructions operating on the pattern's description variables. Architect of the reference structure of walking robot has three two-legged modules. Every leg has three freedom degrees, a slip sensor and tactile sensor to measure the contact which consists of lower and upper levels. In this paper presents some aspects of stability displacement of walking robot MERO

Keywords: Simulation, Walking robot, Modular walking robot.

1 Introduction

Teams worldwide have been focused on goals such as creating an autonomous walking robot equipped with functions like handling objects, locomotion, perceiving, navigation, learning, judgment, information storage and intelligent control, and that can carry out tasks like altering the multitude of the parts belonging to a dynamic universe. Scientists of all times have been permanently mesmerized and have studied the simplest but the most important movement, namely the mechanical movement of humans and animals. Humankind is so much anthropomorphism addicted that it is almost impossible for it to conceive or imagine automatic systems, even having artificial intelligence, and that are not anthropomorphic. The access of man to dangerous areas where his safety is jeopardized made the scientific research approach topics of various purposes and conceive devices that through their performances aim at covering different fields. The architecture of these systems is quite different and depends on their purpose and destination .For example, walking robots protect the environment better because their contact to the ground is discrete, which substantially diminishes the surface to be crushed, the robot's weight can be optimally distributed on the contact surface through controlling the forces. The variation of the distance from the

ground allowed the robot to step over young trees or other vegetation growing in the area it moves on. The walk is defined by the manner the waking robot moves between two points, under specific circumstances. To achieve and guide a walking robot requires thorough knowledge about all walking possibilities because choosing the number of legs and their structure depends very much on the selected walk. The selection of the walk type depends on several elements such as: shape and consistency of the ground the robot walks on walk stability driving and controlling the movements of the elements of the walking systems, speed and mobility movement requires.[1],[13],[14] To use the moving robot as a means of transport, several parameters that characterize its dynamic features can be changed within a wide enough range. Thus, for instance, the occurrence of the supplementary load aboard would change the weight, the center of gravity position, the body inertia moments. The wind and other different forces may act upon the robot and their influence can hardly be anticipated. The action of any kind of similar disturbances might be a cause that produces considerable deviations from the established moving track of the robot. It is very difficult to select the type of walk, mainly during real walking. Therefore, it is necessary that the ground surface to be defined before selecting the walk.[3],[10] The walking robot's steps are a sequel of movements of the legs, coordinated with a succession of movements of the body for the purpose of moving the robot from one place to another.

2 Movement simulation of the modular walking robot MERO by Denavit – Hartenberg formalism

Let us a modular walking robot consists of three modules. [9],[5],[7] Each module has two 3-DOF legs, symmetrically arranged on the platform axis (fig.1). The legs on the right – onto the movement direction are

superscript marked with 2i, i = 1, 6 whereas the legs on the left with 2i-1. Each platform of the rear modules is connected to the platform of first module by a 3-DOF kinematic chain with two links and three rotational pairs. The axes if these pairs are concurrent and perpendicular two by two.

In order to carried out the movement simulation of a leg, a coordinate axes system is attached to each link, with the Denavit – Hartenberg rule [2],[12] This formalism may not only simplify the problem formulation, but can also yield considerable advantage in the solution of simulation problem.

The pairs of each leg are numbered consecutively from *A* which is pair number 1 to *C* which is pair number 3.

The Denavit - Hartenberg systems attached to each link are subscript numbered as the pairs respectively. The platform is designed as link number (0) and the remaining links are numbered consecutively. All pairs of the leg mechanism are rotational and actuated ones.

The characteristic axis Z_i of each pair should be defined. The positive sense of each of these axes is defined arbitrarily. If the axes Z_i and Z_{i-1} are skew with respect to each other, then there is one common perpendicular between them. The perpendicular is designed as the X_i axis. If the Z_i and Z_{i-1} axes are parallel, the X_i axis may chosen as any common perpendicular. The positive direction of the X_i axis is designed as proceeding from Z_{i-1} to Z_i . If the Z_{i-1} and Z_{i-1} intersect, the positive sense of X_i axis is arbitrarily.

When the X_i axes are all defined, then are define both the Y_i axes and the origin of each right hand coordinate system. So, a coordinate system defined is attached to each link. The parameter a_i is defined as the distance from O_iZ_i to $O_{i+1}Z_{i+1}$ axes, measured along $O_{i+1}X_{i+1}$. Because of the orientation of the $O_{i+1}X_{i+1}$.axis, a_i is always positive.

The parameter α_i is defined as the angle between the positive $O_i Z_i$ and the positive $O_{i+1} Z_{i+1}$ axes, as seen from positive $O_{i+1} X_{i+1}$.

The parameter θ_i is the angle between positive $O_i X_i$ and the positive $O_{i+1} X_{i+1}$ axes, as seen from positive $O_i Z_i$.

The parameter s_i is defined as the distance from O_iX_i to $O_{i+1}X_{i+1}$ axes, measured along the O_iZ_i axis.

Under this definition, the Denavit – Hartenberg transformation matrix \mathbf{A}_{i}^{j} has the well-known form:

$$\mathbf{A}_{i}^{j} = \begin{vmatrix} 1 & 0 & 0 & 0 \\ a_{i}^{j} \cos \theta_{i}^{j} & \cos \theta_{i}^{j} & -\cos \alpha_{i}^{j} \sin \theta_{i}^{j} & \sin \alpha_{i}^{j} \sin \theta_{i}^{j} \\ a_{i}^{j} \sin \theta_{i}^{j} & \sin \theta_{i}^{j} & \cos \alpha_{i}^{j} \cos \theta_{i}^{j} & -\sin \alpha_{i}^{j} \cos \theta_{i}^{j} \\ s_{i}^{j} & 0 & \sin \alpha_{i}^{j} & \cos \alpha_{i}^{j} \end{vmatrix}$$
(1)

To would the walking robot's moves is assumed that:

The kinematical length of the binary link (1) is null and it is connected to the platform (0), by pair *A* and to the link (2) by pair *B*; the axis of pairs *A* and *B* are perpendicular.

the binary link (2) is connected to the link (1) by the pair *B* and to link (3) by the pair *C*; the axis of pairs *B* and *C* are parallel.

The A_3^j matrix performed the coordinate transformation of a point belonging to link (3) from $O_4^j X_4^j Y_4^j Z_4^j$ system to



Fig. 1 The Denavit –Hartenberg axis system attached of modular walking robot it is suggested support of technological equipments, for a leg mechanism RRR

 $O_3^j X_3^j Y_3^j Z_3^j$ the system attached to link (2). In a similar manner, the coordinates of lower end point *P* belonging to link (3) from $O_4^j X_4^j Y_4^j Z_4^j$ system to $O_1^j X_1^j Y_1^j Z_1^j$ system attached to the platform (0) is performed by the equation

$$\begin{vmatrix} 1 \\ X_{1P}^{j} \\ Y_{1P}^{j} \\ Z_{1P}^{j} \end{vmatrix} = \prod_{k=1}^{3} \mathbf{A}_{k}^{j} \begin{vmatrix} 1 \\ X_{4P}^{j} \\ Y_{4P}^{j} \\ Z_{4P}^{j} \end{vmatrix}, j = 1, 2.$$
(2)

This matrix equation described the geometrical model of the leg 1 and 2 of the walking robot. The goal of the direct kinematic analysis is to calculate the position, velocity and acceleration of the end point P, in terms of the variables pair

 θ_i^j , $i = \overline{1, 3}$. In inverse kinematic analysis, matrix equation

(2) is solved with respect to the variables pair θ_i^J , $i = \overline{1, 3}$.

The positions of the point P and the positions of the platform with respect to the reference coordinate axes system OXYZfastened to the ground are considered as known. Therefore, position of the point P with respect to the platform coordinate axes system are known.

The movement of the legs the rear modules are controlled by the following equations:

$$\begin{vmatrix} 1 \\ X_{1P}^5 \\ Y_{1P}^5 \\ Z_{1P}^5 \end{vmatrix} = \mathbf{A}_1^3 \mathbf{A}_2^3 \mathbf{A}_3^3 \mathbf{A}_1^5 \mathbf{A}_2^5 \mathbf{A}_3^5 \begin{vmatrix} 1 \\ X_{4P}^5 \\ Y_{4P}^5 \\ Z_{4P}^5 \end{vmatrix} ,$$

$$\begin{vmatrix} 1 \\ X_{1P}^6 \\ Y_{1P}^6 \\ Z_{1P}^6 \end{vmatrix} = \mathbf{A}_1^3 \mathbf{A}_2^3 \mathbf{A}_3^3 \mathbf{A}_1^6 \mathbf{A}_2^6 \mathbf{A}_3^6 \begin{vmatrix} 1 \\ X_{4P}^6 \\ Y_{4P}^6 \\ Z_{4P}^6 \end{vmatrix} ,$$

$$\begin{vmatrix} 1 \\ X_{1P}^7 \\ Y_{1P}^7 \\ Z_{1P}^7 \end{vmatrix} = \mathbf{A}_1^4 \mathbf{A}_2^4 \mathbf{A}_3^4 \mathbf{A}_1^7 \mathbf{A}_2^7 \mathbf{A}_3^7 \begin{vmatrix} 1 \\ X_{4P}^7 \\ Y_{4P}^7 \\ Z_{4P}^7 \end{vmatrix} ,$$

$$\begin{vmatrix} 1 \\ X_{1P}^7 \\ Y_{1P}^7 \\ Z_{1P}^7 \end{vmatrix} = \mathbf{A}_1^4 \mathbf{A}_2^4 \mathbf{A}_3^4 \mathbf{A}_1^7 \mathbf{A}_2^7 \mathbf{A}_3^7 \begin{vmatrix} 1 \\ X_{4P}^7 \\ Y_{4P}^7 \\ Z_{4P}^7 \end{vmatrix} ,$$

$$\begin{vmatrix} 1 \\ X_{1P}^8 \\ Y_{1P}^8 \\ Z_{1P}^8 \end{vmatrix} = \mathbf{A}_1^4 \mathbf{A}_2^4 \mathbf{A}_3^4 \mathbf{A}_1^8 \mathbf{A}_2^8 \mathbf{A}_3^8 \begin{vmatrix} 1 \\ X_{4P}^8 \\ Y_{4P}^8 \\ Z_{4P}^8 \end{vmatrix} .$$

$$(3)$$

Each of these matrix equations is equivalent with three nonlinear equations and has six unknowns, namely variables of the pairs. In the inverse kinematic analysis three out of six unknowns must be imposed from independent conditions.

Because of these variables' particular values, the Denavit -Hartenberg [2],[12] transformation matrices have the following simpler form. Maintaining stability is a special problem that occurs while the robot walks, when one or more legs are in the transfer phase. When all the legs are in the support phase, it is obvious that the projection of the center of gravity is within the support polygon. If one or more legs are in the transfer phase, the geometry of the support polygon changes and it occurs the risk that the protection of the center of gravity moves outside the support polygon. Solutions to such situations depend on how the modular walking robot is configured. In (Fig.2) where are shown a sequence of the computer simulation of the successive gait of the modular walking robot, in C++ and in (Fig.3) animation gait in Studio Max Animation are shown. The gait described in the following section was developed with the purpose of easing the simulation and providing a better understanding of different walk strategies of leg RRR. in different configurations Two main configurations were taken into consideration: RRR - 3 module assembly one module in the front and two modules in the back, in a triangle shape, in (Fig.3).



Fig.2. Computer simulation gait for modular walking robot MERO 1



Fig.3. Computer animation for applications with modular walking robot MERO 1

The software application was built using Visual Studio .NET integrated development environment, was written in C# language and was mainly built around DirectX display SDK (software development kit), package responsible for the onscreen representation of the modular robot configuration. There are two distinct windows that start when the application is launched: The control window and the visual window in (Fig 2). The control window supplies the user with every tool needed to simulate different aspects of the robot unlike the visual one that only shows the graphic representation of the before mentioned robot.

3 Reaction forces modelling in walking robots control

Distribution of reaction forces from the supports of the legs is one of the important problems that must be solved in order to organize movements complicated on landscape for the legs of walking robots with relief complicated. Friction cones of the supports are circular and can be oriented arbitrarily and the points of support - whether they are more than three - are not covered by a plan.

When the number of support points is greater than three, the problem of determining the distribution of forces it is statically indeterminate. To calculate these forces is necessary to know the layout and terrain feature and the positions, dimensions and materials of components of the robot (kinematic and organological size, modules of elasticity).

Choosing an element of this set is based on the travel restrictions imposed on the system. In some cases, the set of admissible solutions can be null if data movement is impossible to achieve, or may be formed from a single element and then there is only a possibility of realization of movements, without being able to take into account additional restrictions.

For a given configuration of the system displacement, the forces of reaction of the supports are determined unequivocally. In the leadership of walking robots with many legs, an optimal distribution of reaction forces of the supports can be taken into account in determining the stepping strategy. Walking robots moving on a terrain, strewn with many obstacles, which can be convex or concave, present a danger that the position of these robots becames unstable.

One of conditions imposed on the motion of walking machines is the stability. The movement of the legged robots can be divided in two modes:

- under condition of the static stability;

- under condition of the dynamical stability.

The main difference between robots which walked under the static stability and under the dynamical stability conditions originates from the fact that during statically walking, the vertical projection of the gravity center of the robot must lay into the supporting polygon, where as during the dynamical walking, this condition can be unsatisfied.

The problem of quasi-statical stability analysis in condition of arbitrary step when the accelerations of points of component elements are much smaller than gravity acceleration is identical with the problem of stability analysis when the robot does not walk. The inertial forces are neglected and the walking can be controlled in a kinematics way.

The investigation of statically stability is based on the notion namely of *hardening configuration*. Hardening configuration is a term used to indicate the rigidly structure of the robot, obtained through the shutting off of the driving motors. The position of the walking robot is stable if the hardening configuration is in posture of stable equilibrium under the action of gravity forces.

The hardening configuration is statically stable if this accomplished the following conditions are accomplished:

The vertical projection of the gravity center of the robot in its entirety (platforms, connecting elements and leg chains, control system and driving, pregnancy, etc.) must be inside the *supporting polygon*. The supporting polygon is the minimum convex area which is obtained by connecting all support points. A body at rest in a gravitational field, subject to ideal connections, has the differential of the gravity center elevation equal with zero.

Vertical projection of the care center of gravity G of the robot is inside the support polygon if all the distances d_i , measured from this point to each P_iP_{i+1} , sides are positive, assuming that the supports on the perimeter of the polygon are numbered clockwise (Fig.4).

Magnitudes of these distances are calculated with the following relationship:



Fig.4 Polygon support of a walking robots

$$d_{i} = \frac{(\xi_{G} - \xi P^{i})(\eta P^{i+1} - \eta P^{i}) - (\eta_{G} - \eta P^{i})(\xi P^{i+1} - \xi P^{i})}{\sqrt{(\xi P^{i+1} - \xi P^{i})^{2} + (\eta P^{i+1} - \eta P^{i})^{2}}}$$

For establishing the stable positions of a walking robot it is necessary to determine the forces distribution in the shifting mechanisms.

Determination of the real forces distribution in the shifting mechanisms of a walking locomotion system which moves in rugged land at low speed is necessary for the analysis of stability. The position of a walking system depends on the following factors:

- the configuration of walking mechanisms;

- the masses of component elements and their position of gravity centers;

- the values of friction coefficients between terrain and feet;

- the stiffness of terrain;

- the shape of terrain surface.

The active surface of the foot is relatively small and it is considered that the reaction force is applied in the gravity center of this surface. The reaction force represents the resultant of the elementary forces, uniformly distributed on the foot sole surface. The gravity center of foot active surface is called *theoretical contact point*.

The modelling of the reaction forces \overline{R}_i corresponding to the support points in the walking robot's movement mechanisms

generally represents the solution of the following static equilibrium equation system:

$$\overline{N}_{i} + \overline{T}_{i} + \overline{R} = 0, \ i = \overline{1, 3},$$

$$\sum_{i=1}^{3} \overline{r}_{i} \times (\overline{N}_{i} + \overline{T}_{i}) + \overline{M} = 0,$$
(4)

where: \overline{N}_i and \overline{T}_i are the normal and tangent components of the reaction force in support point i;

 $\overline{r_i}$ is the position vector of the application point of the reaction force $\overline{R_i} = \overline{N_i} + \overline{T_i}$, in relation to the coordinate axis system annexed to the platform;

 \overline{R} represents the resulting of the forces applied to the elements of the walking robot;

M is the resulting momentum of the forces applied to the elements of the walking robot, calculated in the origin of the coordinate axis system annexed to the platform.

In the case of a uniform straight-line movement of the walking robot on a plane horizontal surface, the reaction forces $\overline{R_i}$ in the support points only have the normal component $\overline{N_i}$ and represent the solutions of the following static equilibrium equation system:

$$\sum_{i} \overline{N}_{i} + \overline{R} = 0; \ \sum_{i} \overline{r}_{i} \times \overline{N}_{i} + \overline{M} = 0;$$
 (5)

The last three items are considered known, as the forces applied to the robot elements – within the mentioned working conditions - are solely the elements' own weights and the weight of the load. As the robot's movement is straight and uniform and there is no slipping, the friction forces between its legs and the ground are zero. These forces have values different from zero only if the robot movement is not straight and uniform or if the support surface is not plane and horizontal. The inertial forces and momentums of the leg elements can be neglected. It is of note that this situation representing an extremely idealised reality – is very rarely met in practice and, as such, the results obtained by solving the system (5) have only theoretical importance. Determining the real distribution of forces in the leg mechanisms of a walking locomotion system, with the classical notations in figure 5, which is moving over rough terrain, is necessary in order to calculate the robot position and to analyze its stability. At a given time, the position of a walking system depends on the following factors: the surface form of the terrain it moves on, the values of the friction coefficients between the terrain and the legs, the configuration of the legs, the rigidity of the terrain the robot is walking on.

As the active support surface area of the last leg element, on the terrain where the movement takes place, is relatively small, it is accepted that the reaction force is applied in the weight centre of this surface Fig. 5 Mathematical modeling of the contact of the modular walking robots walking with the ground

This is the resulting force of the elementary forces, evenly distributed on the entire active surface, with which the last leg element is supported on the terrain. The weight centres of the active surface, for the support of leg i, is called the theoretic point of support and has the notation P^i . In order to determine the robot's position, it is necessary to determine, in each support point, the values of the components of the reaction

force between leg and ground, called: normal component ${\cal N}$, perpendicular on the terrain surface in the contact point and

the tangent component \overline{T} , or the Coulomb friction force – contained in the plane tangent to the terrain surface in the contact point. The value of this vector cannot be greater than the product between the module of the normal component and the respective friction coefficient, between leg and ground. If the value of the tangent component of the reaction force goes beyond this limit, the leg will slip on the support surface until it reaches a stable position, in which the value of this component is equal or below the mentioned limit. As such, the problem of determining the stable position of a walking robot on a arbitrary terrain does not have an unique solution. For each leg there is an interval, which includes all support points in which the equation $\mu N \ge T$ is verified. 'Equal' corresponds to the domain limit.

The studies and theoretical research, seen as in a global approach, lead to a large volume of calculations in modelling and controlling the walking robots motion, which has allowed us to determine solutions only for particular structures, such as MERO modular walking robots [5],[6], with offline data processing. For real time control, albeit these concepts were taken into account, much simpler solutions, described below, were preferred, which avoid the inclusion of laborious calculations in modelling the control process of walking robots.





Fig.6. Computer simulation force distribution for modular walking robot MERO 1

4 Building of a modular walking robot MERO1

At the University "Politehnica" of Bucharest a walking modular robot has been developed to handle farming tools A vehicle like that is the Romanian Walking Robot MERO1 (Fig.4) [7],[5],[12] The modular walking robot MERO1 is a multi-functional mechatronic system designed to carry out planned movements aimed at accomplishing several Two analog-digital and one digital-analogue interfaces (Fig. 5) and a processing computer derived from a *PC*- to acquire the measured data. The values found as a result of the measurements, enabled us to draw the diagrams emphasizing scheduled targets. The walking robot operates and completes tasks by permanently inter-acting with the environment where there are known or unknown physical objects and obstacles. Its environmental interactions may be technological (by mechanical effort or contact) or contextual ones (route identification, obstacle avoidance, etc) The successful fulfillment of the mission depends both on the knowledge, the robot, through its control system has on the initial configuration of the working place, as well as by those



obtained during its movement.

Fig.7 The modular walking robot MERO 1-experimental model (Ion, I., Simionescu, I., 2006)

The modular walking robot MERO is made up of the following parts:

a) The mechanical system made up of three or many modules articulated and shaped according to the requirements of the movement on an uneven ground, the robot's shift system is thus built that it may accomplish many toes' trajectories, which can alter by each step.

b) The actuating systems of feet's have a hydraulic drive;

c) The distribution system is controlled by 18 servo-valves, according to the robot's configuration;

d) The energy feeding system;

e) The system of data acquisition on the shift the system's configuration and the environment;

f). The control panel processing signals received from the driving and the acquiring systems.

For the experimental determination of kinematic and kinestatic parameters of movement of the mechanism system, the robot's activation on and control were accomplished through an "umbilical cordon" by which the robot is tied to an activating group and to the control equipment, respectively. The movement system of the walking robot, as designed in Fig.4, has three degree of freedom and a quite simple mechanical structure, being made up of three serially connected bars and three linear hydraulic motors.

the variation of the kinematic and kinetostatic parameters of the MERO walking robot's movement system.



Fig. 8 Stance transducer



Fig. 9 Block scheme of acquiring data system

Analyzing these diagrams, one can draw the following conclusions:

- the driving forces in the active couplers of the movement system, calculated for the extremity corresponding to the negative abscissa of the working field, taken into account, do have values higher than those calculated in the remainder of the field, and this due to the rise in the values of the pressure angles;

- this value area is avoided by the walking program, so that the variation in the driving forces, remains within the admissible limits of an optimal power consumption.

In order to determine the walking robot's walking stance while moving (maintenance of horizontal position) we attached the platform a sensor made up of two incremental transducers aimed at perceiving the walking stance (Fig. 8). It determines the platform's leaning both at the sagital and frontal level. Fig.10

The sensor is made up of a body hanged by a rigid rod to a sphere leaning on two rollers.



Fig.10. Variation tilt platform modular walking robot MERO1

Restoring fixed platform height is achieved by using information provided by sensor height of the platform by ordering simultaneous vertical movement, for all legs in support phase. Roll over safety limit is determined by the value of the limit of stability required for a particular regime or walk and sizes depending on the reaction forces of the contact points of the feet with the supporting surface, these forces should not fall below limits of the condition to avoid tipping over. Bringing the projection center of gravity of the robot walked in the area of stability, in inside the support polygon is achieved by horizontal movements of the platform (block ordered by maintaining stability).

5 Experimental measurement of the reaction forces in the contact points where the robot's feet reach the ground

We carried out several experiments simulating several walking situations, in order to measure the values of the forces at the point where the robot touches the ground. both in the sagittal and the front plane.

Subsequently we found the values of the reaction forces in pre-set circumstances, when the four-legged robot is walking across an even horizontal ground

The force cells placed on the *MERO* walking robot's feet and the increment conduct sensor enable the robot to control its direction by adjusting its slant, during the tests the robot carries only its own weight namely 84.8 kg.

The tests were resumed with an additional load weighing 21.2 kg and thus the full load reached 106 kg,

Using the walking robot as transportation means we can change a few parameters defining its dynamic properties, at an enough large extent For instance, an additional load placed on the platform, changes the weight, the barycenter's position and the inertia moments of the robot's body. We can apply on the walking robot several forces on the walking robot such as, the resultant of the wind's action, whose influence can hardly be anticipated. The contact cells also ensure the protection of the force transducers. The force transducers we have made use of, turn the variation of a mechanical value (such as linear or angular movement achieved by distorting an elastic element) into the variation of an electric value such as voltage current by means of the electro-rezistive transducers. Each module of the *MERO* walking robot is equipped with two identical sensors placed at the ends of the robot's feet Each sensor has four electro-resistive transducers (TER) connected in full bridge (Wheatstone type). Elastic element on which TER is bound has the form of a plane framework. It was sized for an admissible material tension of $\sigma_a = 270$ MPa and a maximum load on the sensor of F = 600 N

6 Conclusions

The movement simulation of the walking robots may be idealized into a mathematical model for the purpose of kinematic analysis. The techniques of idealization can play the decisive role in easiness, precision and time of calculus for the problem solving. The Denavit – Hartenberg method is numerically robust, the solutions are either exact in the sense that is possible to refine them up to an arbitrary accuracy. A modular walking robot could have one or more modules. The motions of the legs must be coordinated so that the conditions of the gait stability of the system to be ensured.

The *MERO* modular walking robot is a multi-functional mechatronic system designed to carry out planned movements aimed at accomplishing several scheduled tasks. The walking robot operates and completes tasks by permanently inter-acting with the environment where there are known or unknown physical objects and obstacles. Its environmental interactions may be technological (by mechanical effort or contact) or contextual ones (route identification, obstacle avoidance, etc)

The successful fulfillment of the mission depends both on the knowledge the robot, through its control system has on the initial configuration of the working place, and by those obtained during its movement.

7 References

[1] A.P. Bessonov, N.V. Umnov, (1973) *The Analysis of* Gaits in six-legged Robots according to their Static Stability, Proc. Symp. Theory and Practice of Robots and Manipulators, Udine, Italy,

[2] Denavit J., Hartenberg R.S., (1955), A Kinematic Notation for Lower Pair Mechanisms Based on Matrices, Journal of Applied Mechanics, Tr. ASME, Vol. 77

[3] S. Hirose and Y. Umetani, (1980). "The basic motion regulation system for a quadruped walking machine", ASME Paper No. 80-DET-34, September

[4] V. Kumar and K.J. Waldron, "Gait analysis for walking machines for omnidirectional locomotion on uneven terrain", in Proceedings of the 7th CISM-IFTOMN Udine, Italy, 37-62 (1988).

[5] I.Ion, I. Simionescu, A. Curaj, (2005) MERO Modular Walking Robot Support of Technological equipments, The 8th International Conference on Climbing and Walking Robots, September 12-14, 2005 London, UK

[6] Ion I., Vladareanu L., Simionescu I., Vasile A.,(2008).The gait analysis for modular walking robot MERO walk on the slope Proceedings of the 9th WSEAS International Congress on Automation and Information (ICAI '08) pp 222-229, 21-24 June Bucharest Romania

[7] Ion I., Simionescu I., Curaj A. (2003). The displasement of Quatrupedal Walking Robots, Proceedings of the 11th World Congress in Mechanism and Machine Science, August 18-21, Tianjin – China.

[8] Ion I., Simionescu I., Curaj A. (2002). Mobil Mechatronic System with Applications in Agriculture and Sylviculture. The 2thIFAC International Conference on Mechatronic Systems, December 8-12, -Berkeley University – USA

[9] Ion, I., I Simionescu, A., Curaj, A. Vasile (2010), MERO Modular Walking Robots, Solution for Displacing Technological Equipments on Irregular Terrains Proceedings of the 13th World Congress in Mechanism and Machine Science, Guanajuato, México, 19-25 June, 2011

[10] R.B. McGhee, G.I. Iswandi, (1979) Adaptive Locomotion of a Multi-legged Robot Over Unarranged, IEEE Trans. On Systems, Man, and Cybernetics, Vol. SMC-9, No. 4, April 1979, pp. 176-182

[11] R.B. McGhee and A.A. Frank, 1968). "On the stability properties of quadruped creeping gaits", Mathematical Biosciences, **3**, 331-351 (

[12] Simionescu I., Ion I., Ciupitu L (2008). Mechanisms of Industrial Robots Ed.AGIR Bucharest 2008(in Romanian)

[13] S.M. Song, and J.K. Waldron, (1987) An Analytical Approach for Gait Study and its Application on Wave Gait, International Journal of Robotics Research, Vol. 6, No. 2, 1987, pp. 60-71

[14] Waldron, K.J., "Force and Motion Management in Legged Locomotion", Proceedings of 24 th IEEE Conference on Decision and Control, Fort Lauderdale, 1985, pp.12-17.

NSET - A Computational Fluid Dynamics Educational Tool

B. Chernyavsky¹

¹Hydrogen Research Institute, Université du Québec à Trois-Rivières, C.P. 500, Trois-Rivières, Quebec, Canada, G9A 5H7

Abstract - This paper introduces open-source Numerical Simulation Educational Tool(s) (NSET) software package designed to facilitate teaching of numerical simulation techniques for fluid mechanics and other scientific and engineering applications requiring numerical modeling. It consists of the three main modules — case library, a set of numerical solvers, and graphical interface integrating all components into a user friendly package. Unlike classical numerical simulation packages, the main purpose of NSET is to provide assistance in teaching of variety of numerical methods, the criteria, reasons and guiding principles of their selection for specific problems and comparison of their strengths and limitations. NSET package contains tools designed to facilitate comparison of the results obtained using various numerical methods, to illustrate the logic of solvers design, and to enable exploration of stability and accuracy of numerical methods. NSET package is designed with emphasis on simplicity, transparency and modularity of its components to increase its versatility by making it suitable for students with a diverse range of educational backgrounds and level of training in CFD.

Keywords: Educational tool, CFD, Modeling, GUI.

1 Introduction

Computational fluid dynamics (CFD) has become a powerful tool used to solve a wide range of scientific and engineering problems involving fluid flow and heat transfer. The success of its utilization depends, however, on the user's skills in the numerical techniques and their ability to select numerical approach most appropriate for a physical problem under consideration. As CFD gains wider and wider application in the research and engineering community, the number of cases where it is being misapplied or used less than optimally is, unfortunately, also growing, often being caused by the lack of understanding of areas of applicability of specific numerical methods or software packages. Modern software packages provide increasingly powerful tools incorporating wide variety of advanced numerical solvers. This growth of capability - and complexity - often unavoidably results in a non-transparent interface options and implied, but not explicitly stated compatibility issues between component solvers, made further complicated by the lack of transparency in underlying solver mechanics, which is often

proprietary and unavailable for user inspection. In the result, proper utilization of advanced modeling software - or developing research grade problem specific numerical simulation code - requires detailed knowledge of numerical methods strengths and limitations.

This situation emphasizes the importance of teaching undergraduate and graduate students the proper understanding of the requirements of numerical simulation approach and principles and reasoning behind the process of selection of the numerical method, generation of the grid and the choice of optimal computational parameters for a specific physical problem. Teaching of this principles should constitute an integral part of CFD (or other numerical simulation-related) courses, and is stand to particularly benefit from a hands-on approach, where students are encouraged to apply various methods to solution of a standard problem and discuss the results, and to explore the areas of applicability and limitations of specific methods. Unfortunately, both of approaches commonly used to introduce students to numerical simulation methods exhibit unfortunate weakness in this regard.

One approach, particularly popular in teaching undergraduate courses and courses aimed at a more applied utilization of CFD is to use commercial codes, such as, for example, contained in ANSYS package [1], to simulate a variety of flows. The advantage of using commercial codes lies in the possibility to illustrate the power of CFD modeling due to their ability to simulate and to perform post processing, analysis and visualization of complex 3-D transient turbulent flows around realistic geometries, thus illustrating representative application of CFD to a real life engineering and research. It also allows students to familiarize themselves with production CFD codes used in industry and academia and give them a powerful tool which can be used for their subsequent research projects.

On the other hand, these very capabilities make commercial CFD software less than optimal for teaching CFD methods themselves. The disadvantage of the commercial codes is their complexity, which is an unavoidable consequence of their applicability to a maximally broad range of research and engineering purposes. The interface is usually not intuitive and presuppose familiarity with both numerical simulation concepts and specific product features, resulting in a steep

learning curve. Significant time had therefore to be spent on learning the functioning of a specific software package rather than on understanding of underlying numerical methods. Underlying mechanics of commercial simulation software is also not obvious, usually being proprietary and unavailable for users inspection, which often makes simulation results difficult to interpret. The lack of transparency of solver algorithms creates difficulty in performing tests aimed at illustration of relative strengths and disadvantages of specific models, while versatility makes for a very complicated interface.

The attempt to circumvent this lack of solver transparency had been made in a number of codes specifically targeting educational filed, both commercial and open-source (e.g. [2] -[6] just to give a few examples), but majority of these codes are still aimed at simulation of relatively complex realistic 2and 3-D flows, resulting in both necessarily complicated underlying mechanics, which become difficult to understand without extensive prior experience with CFD development, and often also in restriction on number and implementation of numerical methods, favoring the most versatile ones at the expense of diversity of more illustrative methods applicable to a narrower range of problems. This category can also include the increasingly popular OpenFoam [7] project, which has advantage of being community developed, allowing potential user to easily integrate user-written components, and provides its users with both source files and documentation allowing understanding of code working. Its primary purpose as a research grade software, however, results in the necessarily complicated structure and interface, and a rather steep learning curve. While it can be successfully used in a CFD course for a term project, it is less useful as a flexible concept exploration tool which would have benefited from a greater simplicity of both algorithms and interface and a shorter learning time. Learning relative strengths and limitations of various numerical methods and exploring fundamental concepts such as stability and accuracy often does not requires or significantly benefits from complex multi-dimensional cases, and can be more instructively accomplished on an example of a simplified 1-D problem, which would also facilitate a comparative study of the results obtained by various methods.

Yet another alternative approach, particularly popular within dedicated CFD courses, is to introduce a variety of numerical schemes in class, with writing of simple simulation codes implementing these schemes and usually requiring solution of a simplified but representative problem becoming part of an individual or group assignments. This approach both provides students with a hands-on experience and allows them to compare advantages and limitations of different numerical methods. While this approach is undoubtedly very efficient for a graduate students of appropriate specialty, it presupposes good familiarity with programming techniques, and the limited number of course hours and assignments often prevents in depth investigation of subtler details of individual methods, with most of time and effort being spent on code development With NSET, we propose an intermediate approach by providing universal software package which can be used for a different purposes depending on student body background. It can be used to illustrate and compare numerical methods, allowing students to individually investigate their advantages and limitation for standardized cases using simple interface designed for this purpose and open source solvers. Alternatively, for a more advanced user group it can be used to provide a framework into which student-developed subroutines and solvers can be easily integrated, reducing purely programming workload and eliminating the need to develop service subroutines (such as control interfaces, input, output, visualization and standard mathematical routines such as matrix inversion and LU decomposition), which are necessary for numerical simulator functioning but not directly related to the numerical methods per se. It can also provide students a valuable experience of integrating user-developed components into a pre-existing packages, which is as common (or, arguably, even more common) task in CFD research work as the development of completely new simulation software.

The objective of NSET development was, therefore, to assemble a computational package with a minimal learning curve which can be quickly understood by majority of students, transparent user interface and solvers, which source files can be inspected by advanced students, and to provide students with a flexible tool focused at facilitation of learning numerical simulation methods, rather than mastering a specific software package utilization. NSET is also designed to have enough versatility to be useful for a variety of student groups from different fields and with different background and knowledge, as discussed in the next section.

In addition to simplicity and transparency, the third characteristic feature of NSET package is modularity. All package components, both computational and interface related, are written as self-contained as possible and communicate with each other using a standard interface. That allows student developed components (both problem cases and solver routines) to be integrated with the package, resulting in a gradually expanding library with increasing capability, which can potentially make it into a community development project. It should be noted that a concept of community expandable package is far from being new and a number of open source community development projects are active, including the aforementioned OpenFoam project. NSET, nevertheless, retains its advantages of simplicity, transparency and fast learning time, stemming from it focused educational purpose.

The advantage of purpose-designed educational simulation software had been realized by academic community, with several software packages specifically designed to facilitate teaching of numerical simulation approaches had been presented in the recent publications, including [8], which allowed students to perform hands-on exploration of heat conduction in fins, and even design simple fins, and [9], which facilitate teaching of numerical solution of differential equation by using a package composed of both commercial (MATLAB, FLUENT, etc) and in-house developed solvers, which allows student to compare the results obtained from different approaches and learn their respective strengths. NSET package takes this approach further by supplementing the capability to study the simulation methods and their operation requirements with the additional capability aimed at teaching the logic and the structure of these methods implementation, for which it is provided with additional graphical interface tools, illustrating solver's flow-charts, modules association logic and numerical stencil.

NSET package consists of user interface providing control, exploration and visualization tools, a set of numerical solvers including a number of representative numerical methods, which operations are controlled by parameters set in user interface, and a set of case files, containing data describing a problem being modeled, including initial and boundary conditions, mesh, etc. An alternative implementation which is currently being considered, can utilize an internet browser based interface, with computational core located on the centralized server. This approach had been recently grown in popularity (see, e.g., [10]), and it might become the preferred approach to a group-shared educational software. This approach would eliminate potential problems with platform dependencies and installation issues at the expense of complicating community based development. This version of the package is currently being tested to evaluate its practicality compared to baseline.

This paper had been written primarily within a context of teaching Computational Fluid Dynamics (CFD), but NSET design can be easily applied to a variety of other problems involving numerical solution of partial differential equations (e.g., heat transfer, etc).

2 NSET description

2.1 Target audience

In order to be accepted in the educational community, the software should have as wide area of applicability as possible, as it is not realistic to expect teachers to embrace numerous unrelated software products, each with its own unique features and interfaces, intended for a narrow niches or even individual courses. This lack of universality is one of the reasons of relatively limited acceptance of dedicated educational CFD software, and predominance of commercial CFD packages in teaching environment. Accordingly, we attempted to expand the appeal of presented tool to include several different student categories. The first category encompass the students from outside of the usual fields associated with CFD (physics or engineering). Last couple of decades have seen the increased penetration of numerical methods, including methods usually associated with CFD, into a wide range of previously unrelated fields, including biology, chemistry, geology and geophysics, etc. Utilization of numerical methods often lead to important, or even breakthrough development, but the process of incorporation of numerical methods had been often slowed down by the absence of the courses providing necessary background knowledge in the curriculum of students specializing in these fields. The attempts to include the mainstream engineering/physics CFD courses in the curriculum of students of these disciplines resulted in a limited success, due to a significant differences in both goals and the level of knowledge of numerical methods and programming concepts and practices.

For this category of students the objective of NSET is to help familiarize the students without previous exposure to numerical simulations with the basic concepts of numerical modeling, and to provide hands-on experience of numerical solution of partial differential equations. NSET would be used to introduce basic concepts such as numerical accuracy and solution stability as well as parameters determining these concepts. Students should become able to run simulation of well defined physical problem using a variety of applicable models built into NSET. While this task can be achieved using a wide variety of other CFD packages, both commercial and open source, the advantage of NSET for this student category lies in its simplified interface, reducing distractions caused by availability of numerous (and often difficult to understand without previous experience) options usually encountered in research/engineering-grade CFD software. A number of reports (as well as authors personal experience) indicate that successful introduction of numerical approach in the field not traditionally associated with CFD (or numerical methods at large) is often facilitated by utilization of simplified and straightforward interface. While current prototype of NSET is limited to solution of 1-D problems, it had been demonstrated that 1-D simulation can be used for legitimate research purposes in a wide range of research fields varying from marine geochemistry (e.g., [11]) to fuel cell design and optimization. While NSET itself is not intended to be used as a research code, it can be used by students in the term projects illustrating application of numerical methods in their primary disciplines.

The second student category is represented by those engineering and physics students (undergraduate and part of graduate) whose focus lies in utilization of CFD for solution of applied problems rather than fundamental research. These students are the primary users of commercial CFD software packages, and typically have a solid background in fundamentals of numerical methods. Their training is often entirely conducted using commercial packages, which allows them to perform simulation of complex physical and/or engineering problems from the early stages of their training. The downside of this approach is that the students are often initially focused on learning details of software interface in order to be able to set conditions for various cases, after which their focus shifts on visualization and interpretation of the results. Since the exact mechanics of commercial solvers is unavailable to users, and manuals and user guides often provide examples and recommendations on which methods and parameters should be used in common (or representative) cases, but not the detailed reasoning and analysis behind such choice, the skills necessary for the choice of the proper method and parameters for specific application is often left out. This leads to an unfortunately not uncommon cases when numerical modeling is performed using methods less than optimal for specific applications, or when significant efforts are spent on resolution of problems stemming from the wrong choice of parameters (e.g., violation of stability requirements at certain conditions).

The focus of using NSET while teaching that student category is to familiarize students with a variety of factors affecting the results of numerical simulations. Concepts of stability and accuracy can be investigated in depth, and simple but representative numerical problem can be solved using a variety of numerical schemes available with NSET solvers, allowing students to observe and explore relative advantages and disadvantages of various numerical schemes. This study should be conducted after students had already been introduced to several well-known numerical schemes and fundamental concepts of numerical simulation, which would allow them to concentrate on the in-depth investigation of limitations and controlling parameters of simulation using simple interface optimized for this purpose. While similar study can be performed using simple (typically 1-D) numerical solvers written by students themselves, utilization of NSET allows students to avoid many mundane programming requirements unrelated to numerical methods per se, which unavoidable arise in the course of development even a simple simulation code, and to concentrate on the main objective of understanding rationale of numerical method and conditions selection.

The third student category encompasses the graduate students specializing in CFD or related disciplines who will be expected to be able to develop (or, more commonly, substantially modify existing) research grade CFD software. These students can use NSET to both explore various methods and to take advantage of NSET modular nature, simplicity and open access to the underlying solver routines to practice implementation of various numerical methods by adding user developed cases and solver routines. Structure/Flow-chart tool is specifically added to assist students in understanding a connection between a logical structure of the solver and its implementation on programming language level. A typical learning assignment can be calling for implementation of new method which would require development a new solver subroutine, its integration into a package and testing and validation. Students are also encouraged to participate in the expansion of NSET libraries, with the best user-developed cases and routines added to publicly available integrated package

2.2 Software description

The choice of programming language for educational software is often made on the basis of portability and straightforward support of modularity and graphical interface, rather than outright performance. Many educational packages use MATLAB and its derivatives, Visual Basic (see, e.g. [5]), or Mathematica ([6]). NSET package is written in JAVA language. While JAVA is relatively seldom used for scientific simulation codes, the choice was made on the basis of platform independence, support of modularity, streamlined support of graphical interface features and the existence of numerous publicly available libraries, which can be used by students in development their own modules. An additional advantage of JAVA is its relative popularity among students of a wide variety of majors, who, unlike engineering/physics majors, are often not familiar with "classical" simulation languages, such as C or FORTRAN.

NSET package consists of the three main modules, including interface module providing control graphical and incorporating method exploration and visualization tools, case library containing data describing a problem being modeled, initial and boundary conditions and mesh, and computational module containing several numerical solvers (both explicit and implicit) and incorporating generalized computational stencils which by altering appropriate coefficient can be set to represent a wide variety of methods. Numerical solver It includes several standard mathematical functions, such as LU decomposition and fast matrix inversion, which are often used in numerical solvers. The choice of specific solver and setting of stencil coefficients to represent desirable numerical method is done from the user control panel of graphical interface. Solvers used in NSET prototype allows user to select from several well known methods, including Crank-Nicholson, Godunov, Roe, Lax-Wendroff, McCormack, Beam-Warming and Jameson methods.

Since NSET is primarily an educational tool aimed at in-depth exploration of the numerical methods features and areas of applicability, rather than generic simulation software aimed at obtaining solutions for a variety of problems, user interface does not provide the facilities for setting arbitrary user defined problems (which would require a complex interface and preprocessing software of its own). Instead, it allows to load a predefined case from a case library, which includes a set of initial and boundary conditions associated with the test problem, and can also contain a pre-generated non-uniform fixed grid (uniform grid is generated automatically based on user supplied grid resolution Δx , NSET also include an option to use dynamically updated grid which is computed within solvers during the simulation). Teachers or advanced students can create their own case files using standard subroutine interface, increasing a range of problems which can be used with NSET for demonstration of numerical methods properties.

The test case used during prototype development is Riemann shock tube problem, frequently used for demonstration of variety of numerical methods (see, e.g. [12]) and often used for illustration of concepts of numerical diffusivity and oscillatory behavior associated with odd and even orders of accuracy, respectively. In this problem, two halves of the tube separated by the diaphragm are filled with initially stationary gas under different initial conditions (designated "left" and "right" in Fig. 1). At t = 0 the diaphragm is instantaneously removed and two domains come into contact with each other. The advantage of using this problem as a test case is the existence of the exact solution, against which numerical simulation results can be compared. For the purpose of numerical simulation, Riemann shock tube problem is usually formulated in terms of 1-D Euler equations which in the absence of external forces can be written in differential form as

$$\frac{\partial \rho}{\partial t} + \frac{\partial \rho u}{\partial x} = 0 \tag{1}$$

$$\frac{\partial \rho u}{\partial t} + \frac{\partial \rho u^2 + p}{\partial x} = 0 \tag{2}$$

$$\frac{\partial \rho E}{\partial t} + \frac{\partial \rho u H}{\partial x} = 0 \tag{3}$$

where ρ is density, u is velocity, E is total energy and H is stagnation enthalpy. Other test cases can be developed and integrated into NSET package by users.



Figure 1. Riemann Shock tube problem schematic.

User interface (see Fig. 2) provides controls for loading a case from the library, choosing numerical method and controlling parameters, and allows to call stencil, stability and flow-chart tools, export results in ASCII output file and launch the simulation. The results of the simulations appear as a plot in the main graphical panel, along with time step and physical time. User can either explicitly select desirable method from a drop-down menu on the method selection panel, or specify desirable characteristic(s) (i.e. "central difference", "second order", "two-step" or "implicit") first, using method panel to narrow down the selection. User can subsequently select the main parameters controlling the simulation, e.g. grid resolution Δx , time step (fixed or dynamic) or CFL number, number of time steps or physical time at the end of simulation, etc.

Stencil tool is used to graphically represent the discretization used in computational stencil of selected numerical method. When user selects numerical method to be used in the simulation, stencil panel shows the computational stencil associated with selected method (Fig. 2). Alternatively, user can create custom method by creating computational stencil in this panel, adding or removing nodes used in method by clicking on appropriate (j,n) node. Such experimentation with custom computational stencils can improve student's understanding of numerical methods operations, stability requirements (since randomly created stencil can prove to be unconditionally unstable), etc.

More advanced users can also use Stability and Solver Structure/Flow-Chart tools. Stability tool provides information on numerical method stability requirements and allows student to explore their origin using graphical representation of Von Neumann analysis technique. Solver Structure/Flow-Chart tool is used to illustrate the logical scheme of the chosen numerical method. It uses graphical representation of solver algorithm logic which can significantly facilitates its understanding by students [13]. In the Flow Chart mode it can help student to understand the logic of the solver operation and facilitate making the connection between mathematical description of numerical method approach and its software the implementation without distracting their attention by peculiarities of specific programming language or design implementation. Solver Structure mode is intended for more advanced users, and serves to illustrate the details of solvers implementation on a more detailed level. It should be primarily used by students wishing to integrate user-designed components into numerical solvers package, for which purpose it illustrates details such as interface parameters used to connect individual routines and hierarchy of solver elements.



Figure 2. Illustration of NSET user interface showing open Stencil panel, illustrating computational stencil for Crank-Nicholson method, and UML components diagrams

Users can extend NSET package, either by creating additional case files defining new problems using standard format, or by creating incorporating additional solver methods by adding user-written subroutines. User-specified subroutines can be essentially self-contained, with the only modification to the base package being addition of subroutine call in appropriate location. The interface between subroutines is provided by a set of public classes, described in software manual, containing definition of shared variables. Shared variables include three main types, including primary flow variables, auxiliary variables (e.g., global iteration counters, stability relaxation coefficients, etc) and service variables (used for output and dynamic plotting control, etc.), incorporation of which within user-created classes ensure smooth integration of additional numerical methods into the package with minimal alteration of its base components.

NSET package is currently in the final stage of individual components integration within unified interface and, after a final verification test run by UQTR students, the source code will be accessible to public from UQTR website.

3 Conclusions and further development

An educational tool aimed at assisting teachers to illustrate the considerations involved in the choice of the appropriate computational method for numerical solution of a specific physical problem is developed. The tool allows students to solve the problem using a variety of numerical methods and compare and contrast their advantages and disadvantages, and to investigate limitations of each approach. Graphical interface is used to integrate various solvers and tools, to provide user-friendly control options and to visualize the results. Additional tools aimed at advanced students allows examination of code structure, computational stencil configuration and stability conditions. Open source object oriented architecture allows advanced students to develop their own routines and test cases, which can be subsequently integrated in the package, further increasing its area of applicability.

Plans for further developments include field testing it in actual numerical methods course environment and investigation of the alternative web-based interface design.

4 References

[1] http://www.ansys.com/

[2] http://www.easycfd.net/

[3] Romeu Andre Pieritz et al. "CFD Studio: An Educational Software Package for CFD Analysis and Design"; Computer Applications in Engineering Education, Vol. 12, issue 1, pp 20-30, 2004.

[4] C. R. Maliska et al. "Heat Transfer 1.0 - An educational software for heat conduction teaching"; ASME Proceedings of the 32nd National Heat Transfer Conference, Baltimore, Vol. 6, 1997, pp. 53-59.

[5] http://www.openfoam.org/index.php

[6] Mustafa Günal and Atilla Özcan. "Open channel design using Visual Basic"; Computer Applications in Engineering Education, Vol. 16, Issue 2, pp 127–136, 2008.

[7], J. T. Chen et al. "One-dimensional wave animation using Mathematica"; Computer Applications in Engineering Education, Vol. 17, Issue 3, pp 323–339, September 2009.

[8] F. Del Cerro Velázquez et al. "A powerful and versatile educational software to simulate transient heat transfer processes in simple fins"; Computer Applications in Engineering Education, Vol. 16, Issue 1, pp 72–82, 2008.

[9] Filipa Carneiro, et al. "Teaching differential equations in different environments: A first approach"; Computer Applications in Engineering Education, Vol. 18, Issue 3, pp 555–562, September 2010.

[10] Miladin Stefanovic, et al. "Web-based laboratory for engineering education"; Computer Applications in Engineering Education, Vol. 18, Issue 3, pp 526–536, September 2010.

[11] Boris M. Chernyavsky and Ulrich G. Wortmann. "REMAP: A reaction transport model for isotope ratio calculations in porous media"; Geochemistry, Geophysics, Geosystems, Vol. 8, Number 2, 20 February 2007, Q02009, doi:10.1029/2006GC001442.

[12] C.Hirsch. "Numerical Computation of Internal and External Flows"; John Wiley & Sons, 1997.

[13] Reyes Juárez-Ramírez et al. "Teaching undergraduate students to model use cases using tree diagram concepts"; Computer Applications in Engineering Education, Vol. 18, Issue 1, pp 77–86, March 2010.
Experiences with process interaction based simulation in education and research

H.P.M. Veeke, J.A. Ottjes, G. Lodewijks

Dep. of Marine and Transport Technology Faculty of Mechanical, Maritime and Materials Engineering, 3mE Delft University of Technology Mekelweg 2, 2628 CD Delft, the Netherlands E-mail: H.P.M.Veeke@tudelft.nl, J.A.Ottjes@tudelft.nl, G.Lodewijks@tudelft.nl

Abstract - This paper describes the development of simulation education and research in an academic environment. From the start – mid 70's - the real process interaction method has been used. It soon became clear that explaining the principles of simulation is a precarious task. Although the process oriented approach is the most natural way to describe system behavior, it is difficult to explain to students as well as to team members of a design project. Using commercial packages does not help to explain how the model is constructed internally, it merely makes a student an experienced user of the software. We decided to use a general programming platform (Delpi®) and implemented a selfdeveloped simulation tool TOMAS. For communication about the model we developed a gradual translation method: the Process Description Language. This method also enables a student to verify his model and creates greater trust in the quality of a model

Keywords: simulation, process interaction, education, research

1 Introduction

In about 40 years of educating simulation experience, the way of explaining simulation principles dramatically changed. When we received the first simulation lectures ourselves, the lecturer gave the next introductory explanation of the basics of the process oriented approach. "Suppose we have some industrial system, then we can divide the system into dead and living elements. Only living elements have a process. Let me give you the next example". And then we started programming examples. The course focused more on programming skills than simulation skills. Only half of the students understood immediately how to make models, the other half did not and would never do.

Most simulation projects report on the possibilities of simulation software and successes in applying it to complex business cases. As a consequence of this, many packages are built around a limited number of cases and contain predefined components or objects in order to speed up the development of models in the same application field. Often beautiful visualization tools are emphasized to convince management of the necessity to buy a package.

In a university environment however, the requirements on simulation software are quite different.

For education purposes it is necessary to know in detail what is inside the package in order to explain the building elements of a model, the timing mechanism and the use of

queues and distributions. Commercial packages tend to hide these internal constructs.

For research, the requirements are even more different. A researcher requires that he/she is not restricted by the possibilities of a package. The essence of research is to investigate something new, a new conglomerate, which has not been constructed anywhere before. The used platform should offer all possibilities to support this kind of research.

Nowadays most students have to master modelling, before they even start to think about simulation. It appears still very difficult to make a connection between the abstract view in modelling and the concrete process descriptions in simulation. The missing link appeared to be a tool that provides stepwise translation.

This appeared to be also true during simulation projects in practice; it was often difficult to explain to other project team members how the simulation model is constructed. This is especially important when the system to be designed does not exist yet. The project management and team members are usually asked to trust the model based on animations and (positive) results. Apparently there is a gap between the programming language and natural language, in which we can perfectly explain what we mean with "the behaviour".

Therefore the simulation approach on which the software is based, should closely connect to the way students and project team members are used to build models of technological or logistical systems. It makes the threshold to use and understand simulation lower. And for educational reasons it is required that the teacher can easily verify the conceptual model that students have used for their simulation experiments. Experimenting with graphics and animations does not give the required results. Graphics and animations are in fact the results of the simulation activities, and how good they may look, animations can be made with artificial data, as can graphics.

Communication is needed about everything to be designed in a project in order to agree about the result. Communication is required between student and teacher in order to let the student prove he/she understands the theory learned. It appeared that there is a big difference between procedural thinking (program) and process thinking (simulate), which could only be solved by a phased approach before programming.

The next paragraphs explain a new approach, supported by a process interaction simulation platform, to teach simulation and communicate about a simulation model, primarily to support the verification of it and to increase the confidence that it operates correctly.

2 The Platform

In [1] the construction of a tool "Tool for Object oriented Modelling And Simulation" (TOMAS) has been presented. It is implemented as a toolbox in the general programming platform Delphi®, so the complete functionality can be used too. Using Delphi® is not essential, but many students in Delft learn to program Pascal, so keeping the simulation in this language makes it familiar for them.

Delphi offers all possibilities and the flexibility of a general programming language, so there will be no restrictions other than the creativity of the student or researcher. The way of modelling closely matches the modelling as defined in the Delft Systems Approach (DSA) [2,3]. It differs widely from the approaches used in well-known packages. DSA uses as its main modelling element the concept of a "function" that expresses why a particular process is executed and what its contribution is to the environment. Within this notion, a process is described from the company's viewpoint (as a repetitive series of activities of a department /group that handles orders, materials, or resources). Many packages use the viewpoint of the customer or the flowing element itself (a visitor's view). Some call the company/visitor views different; for example, Zeigler et al. [4] calls it flow oriented vs. real process oriented approach.

The role of a visitor will not be denied in this article; it is often an eye-opener for problems, but it is only one possible view. Above that a company view is required to solve not only the problem right, but also the right problem, because it can only be achieved by keeping an eye on the environment continuously. Mostly one starts a research when something is not working satisfactory and then functionality needs to be added, deleted or restructured. Process thinking is generally accepted nowadays; in order to change or design processes correctly, systems thinking is required.

Now that the principal viewpoint has been explained, the question is how to teach students to select a viewpoint in simulation.

3 Education

First of all students have to learn how to describe the time dependant behaviour of a system from the department / machine / employee viewpoint. The model developer (student) should first learn to look at these elements as individual entities, each with its own, single process description. In order to describe the behaviour he/she should take the position of each department / machine / employee element itself, he/she should identify with it. So each process is separately described as if the student plays the element itself.

In describing a process one should ask oneself 2 questions:

What to do when one gets control (becomes "current")
 What to do when one releases control

Question 1 is simply answered by "start programming what the process should do at this moment"; in simulation terms it concerns the description of an "event". Question 2 is more complicated. Does the process have to regain control autonomously or should it depend on some other element in the simulation world? Does the process know the moment in time when it wants control again, or is only the condition known?

Based on the preceding discussion, there are only 2 states possible for each element (in the simulation world):

- (planned to be) active, which is referred to as "Scheduled"
- Inactive for an indeterminate period, which is called "Suspended"

When an element is created, it enters automatically the Suspended state, which it can only leave by a command from another element (eg. Start or Resume). It then becomes Scheduled, where it can decide about it's own state change as shown in fig. 1.



The term "advance" means that the element release control

- for an indeterminate period
- for a defined period of length t time-units
- for a period ending with a defined condition change.

"Interrupt" is forced by some other element.

The next difficulty in simulation programming concerns the communication between processes. It should be arranged in a practical way. In many cases a process has to react when the contents of its queues change. Mostly it signifies that some other process puts an element in it, or an element itself leaves a queue etc. All these cases can be elegantly described by making use of "Queues". Queues should not only be considered as items to collect statistical data about waiting times, but also as sets of elements expressing some state of the element they belong to (idle, busy, interrupted etc) or to transfer elements between processes (e.g orders to be executed, items to be selected from). For this reason, elements should be allowed to be present in more than one queue at any moment.

It is this basic structure that is also used for communication between teachers and students. Starting with a description in natural language one can discuss the structure of the model and which elements should have processes. After that the processes of the elements concerned can be described in a socalled "pseudo language", readable and interpretable by everyone. In this language only the notion of queues and the timing mechanism of fig, 1 is strictly defined.

The example below shows this procedure with the example of a job shop.

Within a process description, concrete activities like selecting a next flow element to be handled, handling an element and delivering it to another function, are modelled. All these activities are combined into one process because they belong to the same actor element. It is this exchange between abstract (functional) views and concrete (process) descriptions that makes the approach distinct from coincidental successes. The functional view guarantees the correct structure and (sub)goals, while the process descriptions fill in the material realization of each separate function with respect to their successors and predecessors. Together with statistical distributions (and a random generator of course), the timing mechanism (fig.1) and the queue concept, TOMAS offers all the basic simulation requirements. The process oriented approach for simulation, described above, was first introduced in Delft by Sierenberg & de Gans [5], later adopted by ECT [6], and finally implemented in an object-oriented way in Delphi® by Veeke & Ottjes [1]. During the last decade it was extended with a server mechanism for distributed simulation [7], and it was made combined discrete continuous by implementing a variable step 1st - 8th order Runge-Kutta

• Machine Group with a job queue, consisting of similar machines and a method to assign jobs

• Job Generator generates jobs at a random time (*Process*)

• Job has a due date and a list of tasks

• Task is an activity with predefined setup time and process duration for a specific Machine Group

Only two element types have a process: a Machine and the Job Generator. For illustration purposes the Machine's process is elaborated further. It can be described with one sentence as a first step:

"Do Continuously: Wait for a job of MyGroup to be assigned, and process it".

This is the first step of describing the behaviour in a way that everybody understands. But how does the Machine notice the availability of a job? For this the reference MyJob of a Machine is introduced. As long as no Job has been assigned, MyJob has value "Nil". The process looks then like:

Repeat

Advance(while MyJob = Nil) MyTask = First Task in Tasklist of MyJob Advance(SetUpTime + Duration of MyTask) Remove MyTask from Tasklist of MyJob If Tasklist of MyJob is empty Register Job Data

Else

MyTask = First Task in TaskList of MyJob Put MyJob In JobQ of Group of Task

Based on this description the modeller can discuss the correctness e.g. is there no operator required during setup time? The modeller could also ask for more information, e.g. How is a job selected if there are more jobs waiting?

Element classes in a jobshop:

[•] Machine belonging to a group MyGroup (Process)

method. Currently TOMAS is being transported to the Lazarus platform (www.lazarus.freepascal.org/). Delphi is distributed bv Embarcadero (www.embarcadero.com/products/delphi). In order to serve a population of about 1000 students that actually use Delphi, easily, we decided to change to the "Lazarus" platform. It is an initiative to construct a Delphi - like environment, which is completely free source (just as TOMAS already was). First tests have shown that TOMAS almost works as-is under Lazarus. This would be a breakthrough for education, because any number of students can use a professional simulation tool then, which requires some basic modelling and programming skills and offers all the flexibility of a general programming language (including the visual part).

4 Research

Often simulation is only applied in the very last stage of a design process, so when the design is almost ready. It only has to be simulated to verify the feasibility of the design and to validate the dimensions and capacities. The construction of a model is a rush job, and if there appear to be problems with the design, it is too late or too expensive to change something.

We advocate another simulation approach.

Firstly, simulation should be used during the design process from scratch. Construct a "growing" model, and extend it gradually together with detailing the design. For example at the start of the design process of the new Rotterdam port area "Maasvlakte II" [8,9] first the total quay length, stacking area and intensity of traffic flows have been determined by means of simulation. Moreover, by using simulation the quantified specifications of the required dimensions support an improved insight for the rest of the project duration. By zooming into this rough model, subsystems like individual terminals can be revealed and detailed. At the end of the day, there is a detailed model that accurately reproduces the behavior of the system and can also be used during the operational stage for shadow operations, planning and 'what if' questions.

Secondly, simulation can perfectly be used to develop and test real time control. An example is the development and implementation of a new routing method for Automated Guided Vehicles (AGV). The method was tested with simulated AGV's (fig. 2a). After that, these were replaced with real AGV's (fig. 2b), but the core of the control still was the original simulation-developed control. Especially in this type of applications the use of a general programming platform pays off. Communicating with hardware equipment is far from standardized.

At this very moment we are designing the control for robots in a food packaging line. The control of one robot is quite easy, and already proven technology. But the cooperation of robots in a packaging line requires extra intelligence that doesn't belong to the robots separately. 332.80



Fig. 2a Simulation of AGV-system



Fig. 2b. Prototyping an AGV-system

By using an open and general programming platform it is possible to design cooperating robots. First simulated robots, but now we're in the process of using real robots. Calibration should be added to the (simulated) control software to ensure accurate positioning. Currently shadow operations in a laboratory environment are prepared to solve all kinds of real disturbances. It is our intention to continue the research in this kind of hardware-in-the-loop developments.

5 Conclusions and future research

Teaching process interaction simulation to (presumably) highly intelligent students seems a straightforward task. However, in an academic environment, the scientific attitude requires proof of both students and teacher that the presented theory has been understood well. Therefore the teacher needs a way to explain his thoughts in a rational way, but the student needs a way to present his way of thoughts and proof the correctness of his model.

Now after years of experience we are convinced to have found a correct way for teaching process interaction. The difficulty was in process thinking instead of programmatic thinking, a difficulty we encountered also in commercial projects where simulation was being used. For education purposes the environment to program should also be as open as possible. This last issue was very important to show the principles of the simulation software itself.

It was even more important because in an academic environment, research is being done to far from standard cases. It has the objective to develop new standards and applications. For this purpose, it should be possible to mould the used simulation software in a way that matches the requirements of the researchers. In the field of hardware-inthe-loop much work has still to be done.

6 References

[1] Veeke, H.P.M., Ottjes, J.A., "TOMAS: Tool for Objectoriented Modeling And Simulation", Proc. Of Business and Industry Simulation Symposium, Washington, Ed. Maurice Ades, pp. 76 – 81, 2000

[2] in 't Veld, Prof. J., "Analysis of organisation problems", Wolters-Noordhoff bv, Groningen, 8th edition, ISBN 90-207-3065-7, 2002

[3] Veeke, H.P.M., Ottjes, J. A., Lodewijks, G., "*The Delft Systems Approach*", Springer,

ISBN 978-1-84800-176-3, 2008

[4] Zeigler, B.P., H. Praehofer, and T.G. Kim. 2000. "Theory of Modeling and Simulation", Academic Press, San Diego

[5] Sierenberg, R.W., de Gans, O.B., "*PROSIM text book*", lecture notes Delft University of technology, Delft, 1982

[6] Ham, R.Th. van der, "*Must: Simulation Software User and Reference Manual*", version 5.50, Upward Systems, Rijswijk, The Netherlands, 1992

[7] Boer, C.A., Verbraeck, A., Veeke, H.P.M., "Distributed simulation of complex systems: application in container handling", SISO European Interoperability Workshop, Harrow, Middlesex UK, june 24-27, 2002

[8] Veeke, H.P.M., Ottjes, J.A., "A Generic Simulation model for Systems of Container Terminals", Proc. Of the 16th European Simulation Multiconference (ESM 2002), Darmstadt, pp. 581-587, ISBN 90-77039-07-4, 2002

[9] FAMAS.MV2 2000-2002, "Towards a new generation of automated terminals on Maasvlakte 2", FAMAS Report, 2000

Design of a Real-Time Simulator for an Electric Vehicle

T. Silloway, Y. Jung, D. Mackellar, F. Mak

Electrical and Computer Engineering, Gannon University, Erie, PA, USA {silloway001, jung002, mackellar001, mak001}@gannon.edu

Abstract - A cost-effective Real-time Simulator (RT-Sim) for academic research and education has been developed to simulate the subsystems of an electric vehicle with both Software-/Hardware-in-the-loop (SIL/HIL). The RT-Sim including the associated software/hardware components was developed for intuitively and swiftly managing different levels of scale and complexity of design while providing sufficient accuracy to satisfy the real-time constrains. The RT-Sim suite consists of a multi-core computer, a hardware platform including an FPGA/microcontroller board, and a console. A methodology was also developed for organizing simulation procedures, managing simulation modules, and achieving simulation accuracy. The RT-Sim suite was evaluated with the Simulink/VHDL models of the electric vehicle after configuring the software/hardware models. The RT-Sim was evaluated by performing SIL/HIL simulations and by analyzing the simulation results. This RT-Sim suite is expected to bridge the various design tasks that span different disciplines by providing intuitive model configurations and management, and precise model simulation and rapid analysis of the results.

Keywords: real time, hardware-in-the-loop simulation, software-in-the-loop simulation, electric vehicle

1 Introduction

Contemporary systems are necessary to meet the swiftly increasing complexity and demanding requirements of applications. Engineering society has been challenged to deliver such rapidly evolving systems within the tight time-tomarket constrain. Simulation-based rapid prototyping is one of the viable solutions for system developers. In particular, developing the systems in automotives, avionics, and power electronics is necessary for employing real-time simulation capability. Therefore, a real-time simulator (RT-Sim) must be capable of addressing the complexity of system integration, performing accurate simulation with real-time ability, and efficiently managing design resources as a competent development tool while continuously supporting developers with an intuitive and seamless design and evaluation environment.

An efficient and perceptive simulation environment is required to allow quickly discovering and evaluating alternative design and control strategies [1]. In particular, an efficient management including model integration, input/output (I/O) allocation, accurate execution, and systematic test creation is one of the key parts for the successful utilization of the RT-Sim. Professional RT-Sims [2] have been utilized in order to provide adequate interfaces for various aspects of software and hardware models, I/O allocations, and test analysis. Other application-specific RT-Sims were developed for different applications, such as electric control unit in electric vehicles (EV) [1], fuel cell in hybrid EVs [3], electric and hydraulic systems in avionics [4], and power electronic system [5].

In general, hardware-in-the-loop (HIL) simulation becomes crucial for the systems that required comprehensive evaluations under realistic conditions with real-time constrains. The HIL simulation [6, 7, 8] not only offers costefficient, fast verifications, but also avoids or at least mitigates risks of contingency situations for microcontrollers in the systems. The HIL simulation permits models of a part of the system to be simulated in real time with the actual hardware of the remainder of the system. In order to integrate models and real hardware together, various I/O interfaces generally provide analog-to-digital and digital-to-analog conversions (ADC/DAC), voltage conversions for different digital signaling, and signal conditioning equipment. More specific HIL simulations, such as controller HIL simulation for rapid controller prototyping and power HIL simulation with transferring real power to the hardware system, are discussed in [9]. Thus, the HIL simulation capability in RT-Sim is instrumental in the development of EV. Developers promptly test the EV platform, evaluate the control strategy, and figure out a breakthrough with the HIL simulation.

As the demand for RT-Sims increases in industry, significant growth in the number of RT-Sims has been seen during the last decade in academia [10]. Unlike industry, academic version of RT-Sims [11, 12] are expected to embrace specific features including meaningful and relevant experience without being limited by laboratory equipment, user-friendly interfaces with increasing sophistication, the flexibility for continuous expansion, and preferred cost-effectiveness. A cost-effective academic RT-Sim with SIL/HIL capability has been developed in order to successfully utilize the RT-Sim in academia, especially for different disciplines including Power Electronics, Electric Drives, Embedded Systems, and Communications.

Section 2 describes the architecture and operation of the RT-Sim. Section 3 expresses evaluation of the RT-Sim with the simulation results and analysis of an EV model. The conclusions and future work are depicted in Section 4.

2 Architecture & Operations of RT-Sim

Architecture of the developed RT-Sim suite is illustrated in Figure 1 (a). The RT-Sim suite consists of three modules—



Figure 1. Block Diagrams of the RT-Sim Suite (a) Architecture and (b) Operation

a console module, real-time simulator module, and hardware module. The console module is aimed for an intuitive management of model creation and configuration by providing seamless interface to develop, convert, and simulate a software model such as Simulink [13] model. In addition, the console provides a means to develop and configure a hardware model, such as Hardware Description Language (HDL)-VHDL or Verilog HDL [14]-model, via a Field Programmable Gate Array (FPGA) or a microcontroller as a hardware module. The console communicates with the RT-Sim through a 100M bps Ethernet channel. The hardware model programmed in the hardware module is integrated to the RT-Sim via a Peripheral Component Interconnect (PCI) card installed in RT-Sim. The developed RT-Sim contains multi-cores, running the QNX real-time operating system with an analog/digital input and/or output PCI card.

Figure 1 (b) shows a brief model creation and operation procedure. A software and/or hardware model is created through the console by invoking a built-in modeling procedure. This modeling procedure initiates to launch underlying modeling tool suites including Simulink for software model and Xilinx ISE [15] for FPGA-based hardware model. A model configuration step must be completed to properly configuring the models created in different environments, such as software and hardware before transmitting the entire models to the RT-Sim for simulation. In particular, hardware models must be programmed via the underlying hardware development tools before forwarding

| ile Help | |
|--|--|
| and mine dimulator | |
| ceal-Time Simulator | b Desl-Time Simulator |
| usign, configure and similate reacts are | |
| Model Model | 1 Simulation |
| Composition Mode | 1 Name: Speed_Controller |
| Mode | 1 Location: Z:\RT_Eval_2010\: Browse |
| Design | |
| | Start |
| HW Model 3 | mulation |
| Design | |
| Output | |
| Configuration | |
| | |
| Shared W Memory | |
| Paremeter | |
| | Close |
| Electrical and | d Computer Engineering |
| | |
| | |
| | (a) |
| ○ O Source Block Parameters: rtsimReadValu | (a) |
| O Source Block Parameters: rtsimReadValu T-Sim Read Value (mask) (link) | (a) |
| O Source Block Parameters: rtsimReadValu T-Sim Read Value (mask) (ink) eads value from the RT-Sim Shared Memory Region. | (a) |
| O Source Block Parameters: rtsimReadValu T-Sim Read Value (mask) (link) aads value from the RT-Sim Shared Memory Region. treasure | (a) |
| Source Block Parameters: rtsimReadValu T-Sim Read Value (msk) (ink) kads value from the RT-Sim Shared Memory Region. traineter: infable Name - variable name to be read | (a) |
| O O Source Block Parameters: rtsimReadValu T-Sim Read Value (mask) (ink) Reads value from the RT-Sim Shared Memory Region. arameter: Infable Name - variable name to be read | (a) |
| Source Block Parameters: rtsimReadValu T-Sim Read Value (mask) (link) eads value from the RT-Sim Shared Memory Region. arameter: arable Name - variable name to be read arameters | (a) e Sink Block Parameters: rtsimWriteValue RT-Sim Ware Valar (max) (Ink) Write valar (max) (Ink) Write valar to RT-Sim Shared Amony region. Parameters Variable Type - type used for saving variable Parameters Variable Kame |
| Source Block Parameters: rtsimReadValu (T-Sim Read Value (msk) (ink) eads value from the RT-Sim Shared Memory Region. arameter: fariable Name - variable name to be read arameters Arable Name | (a) |
| Source Block Parameters: rtsimReadValu T-Sim Read Value (msk) (ink) kads value from the RT-Sim Shared Memory Region. tarable Name variable Name Motor Speed .RPM | (a) Performance of the second secon |
| Source Block Parameters: rtsimReadValu Tr-Sim Read Value (mask) (link) etads value from the RT-Sim Shared Memory Region. arameter: arameter: arameters Arable Name Motor, Speed, RPM | (a) e Sink Block Parameters: rtsimWriteValue RT-Sim Wree Value (masku) (Ink) Writes value for Raving Variable Variable Type - type used for saving variable Parameters Variable Kame Motor: Speed, RPM Variable Type globie |
| Source Block Parameters: rtsimReadValu T-Sim Read Value (mask) (link) teads value from the RT-Sim Shared Memory Region. arameter: frable Name Motor_Speed_RPM | (a) e Sink Block Parameters: rtsimWriteValue RT-Sim Write Value (mask) (link) Writes value to the RT-Sim Shared Memory region. Parameters: Variable Name - Variable Name to be used Variable Name - Works Speed, RM Variable Name Works Speed, RM Variable Type double |
| Source Block Parameters: rtsimReadValu T- Sim Red Value (max) (ink) keads value from the RT-Sim Shared Memory Region. Varameter: anable Name - variable name to be read arameters Ariable Name Motor Speed_RPM | (a) (a) (b) (c) (c) (c) (c) (c) (c) (c) (c |
| Source Block Parameters; rtsimReadValu T-Sim Read Value (mask) (link) tatads value from the RT-Sim Shared Memory Region. Airameter: Airable Name - variable name to be read arameters (Ariable Name Motor_Speed_RPM OK Cancel Help | (a) (a) (b) Sink Block Parameters: rtsimWriteValue RT-Sim Write Value (mask) (link) Writes value to the RT-Sim Sared Memory region. Parameters: Variable Type - type used for saving variable Parameters Variable Type - type used for saving variable Parameters Variable Type Jourbe Type double |
| O Source Block Parameters: rtsimReadValu Tr-Sim Read Value (mask) (link) leads value from the RT-Sim Shared Memory Region. arameter: arable Name Motor_Speed_RPM OK Cancel Heip | (a) Re RT-Sim KBlock Parameters: rtsimWriteValue RT-Sim Write Value (mask) (link) Writes value to the RT-Sim Shared Memory region. Parameters: Variable Name - Variable name to be used Variable Name - Variable Name to be used Variable Name - Wander Name - Variable Name Variable Type double OK Cancel 1940 Apply |
| Source Block Parameters: rtsimReadValu Tr-Sim Read Value (mask) (link) Read Value (mask) (link) Reads Value (mask) (link) Arameters Arable Name Motor Speed, RPM OK Cancel Help | (a) e Sink Block Parameters: rtsimWriteValue RT-Sim Wre Valae (mass) (Ink) Writes valate in RT-Sim Shared Memory region. Parameters Variable Type - type used for saving variable Parameters Variable Kame Motor: Speed, BPM Variable Type double OK Cancel Help Appr |
| Source Block Parameters: rtsimReadValu Tr-Sim Read Value (mask) (ink) teads value from the RT-Sim Shared Memory Region. arameter: arameters arameters (ariable Name Motor Speed_RPM | (a) e Sink Block Parameters: rtsimWriteValue RT-Sim Write Value (mask) (link) Writes value to the RT-Sim Skared Memory region. Parameters Variable Time - write Value Parameters Variable Time - type und Variable Tipe double CK Cancel Help Apply BT-Sim Write Value |
| Source Block Parameters: rtsimReadValu Tr-Sim Read Value (mask) (link) Eads value from the RT-Sim Shared Memory Region. Tainneter: Tariable Name Motor, Speed, RPM OK Cancel Help RT-Sim Read Value Notor Speed RPM | (a) e Sink Block Parameters: rtsimWriteValue rt-Sim Write Value (max) timk Wrate to the rt-Sim Shared Memory region. Parameters Variable to the rt-Sim Shared Memory region. Parameters Variable Name Name Name Name Name Name Name Name |
| Source Block Parameters: rtsimReadValu Tr-Sim Read Value (mask) (link) tarameter: tarameter: fariable Name variable name to be read tarameters fariable Name Motor_Speed_RPM | (a) (a) (b) Sink Block Parameters: rtsimWriteValue RT-Sim Write Value (mask) (link) Writes value to the RT-Sim Skared Memory region. Parameters: Variable Type - type used for saving variable Parameters Variable Type - type used for saving variable Parameters Variable Type double (c) Cancel Help Apply Apply Cancel Melp Apply Motor_Speed_RPM Motor_Speed_RPM Motor_Speed_RPM |
| Conce Block Parameters: rtsimReadValu T-sim Read Value (mask) (ink) Reads value (mask) (ink) Reads value from the RT-Sim Shared Memory Region. Tarameter Tarable Name Motor_Speed_RPM OK Cancel Help RT-Sim Read Value Motor_Speed_RPM rtaineBaadValue | (a) e Sink Block Parameters: rtsimWriteValue RT-Sim Write Value (mask) (link) Writes value to the RT-Sim Shared Memory region. Parameter: Variable Name - Variable Name to be used Variable Name - Variable Name to be used Variable Type double OK Cancel Help Apply RT-Sim Write Value Motor_Speed_RPM Type: double -trimWrite Value Type: double -trimWriteValue -t |

Figure 2. Console of the RT-Sim Client (a) Frontend and (b) Model Configuration

both the software model and software/hardware configuration to the RT-Sim. The RT-Sim compiles software model according to the model configuration after receiving a signal for launching the simulation from the console. Then, the RT-Sim partitions and schedules the executable and distributes them to the appropriate cores according to the system partitioning and scheduling outcomes established during the model compilation. The RT-Sim performs SIL/HIL simulations if necessary. Consequently, the RT-Sim generates and transmits outputs of the simulation to the console for displaying and storing the simulation results.

2.1 Console Module for Model Management

The console provides an environment for modeling both software and hardware components by launching underlying modeling tool suites, i.e., Simulink, Xilinx ISE, and so on. The console in general has the responsibility to transmit the source code of the software model and configuration code of both software/hardware models if necessary. The console also controls simulation including starting and terminating simulation while displaying and saving the simulation outputs receiving periodically if needed.

Figure 2 (a) illustrates a frontend of the console developed. Four menu buttons initiates aforementioned operations: (1) a software model design button is for invoking a software modeling tool, such as MATLAB/Simulink. The MATLAB/Simulink code must be converted to C code via Real Time Workshop (RTW) before transmission to the RT-

Sim; (2) a hardware model design button is for launching a hardware modeling tool, such as Xilinx ISE. A HDL model is programmed to the target FPGA via the FPGA programming tool; (3) an output configuration button is for starting configuration of each model's interface. An in-house configuration module was developed for this purpose. I/O configuration of a Simulink model and/or a HDL model provides information to establish SIL/HIL for real-time simulation: and (4) a simulation start button establishes communication to the RT-Sim and initiates transmission of necessary codes. A socket-based communication between the console and the RT-Sim continues until the simulation is terminated. Similar to other GUI-based tool, the console supports users creating and accessing projects to encapsulate the model preparation procedure and to swiftly retrieve existing simulation.

Figures 2 (b) illustrates blocks of the custom Simulink system function, S-Function. A primary goal of the S-Function blocks is for seamless model integration regardless of the models designed in software or in hardware. Another important role of the S-Function blocks is for efficient interoperation of simulation models. Thus, the S-Function blocks permit accessing to the shared memory region allocated on the RT-Sim. The shared memory contains values of I/Os of software/hardware models being simulated. The S-Function blocks are composed during the model configuration process in the console. The composed S-Function blocks are dynamically used for creating real-time simulation in SIL/HIL during the model simulation.

Two S-Function blocks, named as "rtsimReadValue" and "rtsimWriteValue," are available for model composition as seen in Figure 2 (b). The "rtsimReadValue" block with a parameter, such as a variable identification, grants accessing from the associated location of the shared memory region. The "rtsimWriteValue" block with the associated parameters, such as variable identification and type, allows accessing to the specified location of the shared memory region. With the S-Function blocks, RTW generates header information during the model conversion process. RTW utilizes the Target Language Compiler (TLC) file to generate the inline macro call required to access the shared memory region. In particular, the RT-Sim verifies I/O configuration of all models for simulation according to the parameters given



2.2 Hardware Module for HIL Simulation

Hardware module contains hardware models, which provide time-critical functionalities of the hardware components utilized in the target system for the HIL simulation. This module is designed for expandability and interchangeability so that any software model can be swiftly and intuitively replaced with the associated hardware models whenever necessary. As seen in Figure 3 (a), an FPGA (i.e., Xilinx Spartan 3E)/microcontroller (i.e., 32-bit MicroBlaze) board is used as the hardware module. In particular, the FPGA board is connected to a PCI card in the RT-Sim via an interface board seen in Figure 3 (b), which passes analog/digital I/O signals after converting levels of the signals if necessary. For instance, 3.3 Volt LVTTL signals are converted to 5 Volt TTL signals and transmitted to the PCI input ports in the RT-Sim.

2.3 Real-Time Simulator (RT-Sim)

The RT-Sim is a multi-core (i.e., Intel Core i7 Quad Core Processor 870 with VT (2.93GHz, 8M)) computer executing a real-time operating system (i.e., QNX 6.5). The RT-Sim also contains a PCI card for the HIL simulation. Figure 4 (a) shows an organization of the RT-Sim. The four cores with double thread per core are scheduled to execute the distributed tasks in parallel on the real-time operating system. The RT-Sim performs both SIL and HIL simulations. In particular, the RT-Sim provides a hardware-interface task to process and establish the hardware I/O for the HIL simulation via the installed PCI card which is capable of processing analog/digital I/Os. The rate of the interface task is dependent upon the rate of the fastest software model executing on the simulator. Threads 1 to 4 in core 1 and 2 are allocated for the SIL simulation, core 3 interfaces to the PCI for the HIL simulation, and core 4 processes the socket server and communication with the console.

Figure 4 (b) illustrates the operation flow of the simulator. The RT-Sim behaves as a server and waits for a console module to request a simulation to begin. After the initial request, the RT-Sim launches a thread to service the console. The RT-Sim receives all software simulation models written in C programming language and associated configurations. The RT-Sim automatically compiles the models and configurations upon receiving. The RT-Sim transmits an acknowledgement to the console whether or not file transmission and compilation are successful. The RT-Sim holds any further operation if any error is detected during the model compilation. The RT-Sim, then, resumes the aforementioned simulation model and configuration receiving and compiling operations. The RT-Sim initiates the simulation of the models according to the configuration upon receiving a command of the simulation start. During the simulation, the



Figure 3. (a) FPGA/Microcontroller and (b) Interface Boards in the Hardware Module



Figure 4. RT-Sim (a) Organization and (b) Operation

RT-Sim continuously transmits output data of the models to the console at a predefined rate (i.e., a default rate is set as output data transmission per second) for displaying and/or saving the simulation results. The default rate is determined to provide the most processing power on the simulation and so that the console is not overwhelmed with too much output data at once. The RT-Sim continues to simulate the models either reaching to the simulation time pre-set or end of the simulation. The RT-Sim is also interrupted to terminate current simulation by the console.

The RT-Sim employees a client-server socket-based 100 M bps Ethernet connection as the communication layer to the console. FTP was initially installed for the communication. However, the FPT-based communication was not sufficient for the real-time data transmission because of causing unnecessary overheads including the constant authentication for a file transmission. For instance, one of the FTP overheads was maintaining the IP addresses of clients. This overhead is disappeared with the socket-based architecture. Therefore, the RT-Sim adopted the client-server architecture in that the RT-Sim and the console, respectively, act as the server and the client. The RT-Sim and the console transmit custom XML messages to provide flexibility in any model to be executed on the simulator without any modification to the communication layer.

Quard-core in the RT-Sim is partitioned for concurrent operations in order to support the real-time simulation. A new simulation generates a set of partitions. The RT-Sim reserves a core for the console communication for any simulation. The remaining threads in the cores are assigned for the different tasks if necessary. This core partition is performed according to a set of partitioning scenarios built-in the RT-Sim. For instance, a real-time simulation involves both of the SIL and HIL simulations. RT-Sim assigns a core for the HIL simulation, another core for managing current simulation, and the remaining cores for the SIL simulation.

The core assigned for the HIL simulation accesses the PCI card for exchanging I/O signals to the hardware models running on the hardware module (i.e., VHDL/Verilog models in FPGA and/or a microcontroller) via the interface board. Up to four independent software models run concurrently on the four threads in two cores. More than 5 software models can be virtually executed in parallel by sharing the same core in different scheduled time. The remaining core is responsible of monitoring core usage, task distribution and scheduling, and output management. For instance, this core distributes multiple software models to the available cores according to the dynamic scheduling scheme implemented in the RT-Sim.

In order to efficiently deliver the output from a model to an input of the same and/or other model, the RT-Sim is designed with a shared memory scheme. In addition, this flexible shared memory scheme supports seamless configuration for various simulations. The shared memory provides a virtual interface between all the models in simulation. Each model has the responsibility to update its outputs to the dedicated location of the shared memory during the simulation. The contents of the shared memory are transmitted to the console periodically for processing and analysis of the simulation outcomes. These values are also saved on the simulator as a precaution. Figure 5 illustrates an example of the shared memory configuration. Simulink models 1 to N connect to the shared memory region during the model initialization and access the region as required throughout the simulation. The analog and digital I/Os of the hardware model for the HIL simulation are updated to the shared memory locations at the specified frequency in the simulation configuration. The diagnostics screen seen in Figure 5 (a) displays values of the I/O signals transmitted and received via the PCI for verification of the HIL simulation.

3 Evaluation of the RT-Sim

A schematic of the RT-Sim suite shown in Figure 1 (a) has been developed for evaluation of the RT-Sim. The developed RT-Sim suite is illustrated in Figure 6. Figure 6 (a) is the console client. Figure 6 (b) and (c) are the hardware module and the RT-Sim respectively. The hardware module consisting of an FPGA and an interface boards is also shown in Figure 6 (c). The console and the RT-Sim are connected with 100 M bps socket-based Ethernet. The hardware module for the HIL simulation is connected to the RT-Sim via the PCI card for 8 analog outputs (i.e., +/- 12 Volts), 3 digital outputs, and 4 digital inputs (5 Volts). The interface board converts

LVTTL and/or PCI 33/66 MHz and AGP-2X single-ended type signals (i.e., 3.3 Volts) received from the FPGA board to 5 Volts TTL signals vice versa although the FPGA board supports 19 different standards including LVTTL, LVCMOS2/18, Bus LVDS, and so forth. Two 14-bit output ADCs with 1.5 M Hz sampling rate are used for analog inputs to the hardware models programmed in the FPGA. A 12-bit resolution DAC with 6 analog outputs including a VCC and a ground signals is used for analog outputs to the PCI card.

3.1 A Case Study: An Electric Vehicle (EV)

A subsystem of an EV was used for evaluation of the RT-Sim suite. As seen in Figure 7 (a), the complete Simulink closed-loop model of an EV subsystem was modeled. The EV subsystem consists of a battery for supplying DC power, a DC motor for driving, and components for sensing speed of EV and charging level of the battery, and a control unit for controlling speed and torque of the DC motor. The EV subsystem includes a 25 horsepower, four quadrant operation wound DC motor. The model is designed to execute in



Figure 7. (a) A Simulink Model, (b) Simulation Results of an Electric Vehicle, and (c) Simulink Model with RT-Sim Inputs/Outputs

discrete time with a sampling rate of 1 micro-second. The control unit controls speed of the armature. Inputs of the control unit are the desired speed, in RPM, and the armature speed, in RPM. Outputs of the control unit are the speed change and armature current. The control unit controls the armature current and prevents the current from surpassing the rated armature current. Inputs of the control unit are the armature current, armature speed, speed change, and change in armature current. The control unit outputs the PWM pulses to set the armature voltage to the desired voltage to achieve the desired armature speed. The control unit also outputs the control values used to generate the PWM pulses. The DC Motor is connected to a 30 volt battery with a linear load torque. The DC Motor contains a DC-DC converter connected to the PWM pulses to provide the desired armature voltage. The motor inputs are the torque load, PWM pulses, and battery voltage. The motor outputs are the armature current, armature speed, armature voltage, and the field voltage.

Figure 7 (b) illustrates the simulation results with the reference speed and engine speed. The results, with a reference speed of 300 RPM, illustrate that the motor is able to achieve the desired speed within 0.3 seconds with a 10 RPM overshoot. These results are used as the baseline for the SIL and HIL evaluation of the RT-Sim.

Figure 7 (c) illustrates the Simulink model after connecting the RT-Sim s-functions to each sub-model. Three models—(1) DC Motor with a Battery, (2) Current Controller, and (3) Speed Controller—are separated into distinct models based upon the intended processing threads. The blocks in the left are the "rtsimReadValue" s-functions for specifying the values to read from the shared memory. The blocks in the middle are the three separated models of the EV. The blocks in the right are the "rtsimWriteValue" sfunctions for configuring the values to write to the shared memory.

3.2 The SIL and HIL Simulations of an Electric Vehicle

As described in Section 3.1, three sub-models are identified before compilation. The executable codes of the associated sub-models are distributed to three threads in two cores (e.g., thread 1 in core 1 for the speed controller, thread 2 in core 1 for the current controller, and thread 1 in core 2 for the DC motor and battery unit) for the SIL simulation.

The same EV model was modified for the HIL simulation. Figure 8 (a) shows details of the current controller sub-model for the HIL simulation. Since the PWM used in the current controller requires 0.625 MHz frequency, the PWM generator in the software model was substituted with a PWM generator implemented on the hardware module (e.g., FPGA). Since the current controller requires four PWM inputs, a 0.625 MHz PWM (i.e., PWM-1) was generated and the other PWMs were also generated. For instance, the PWM-2 is the inverse of the PWM-1. The PWM-3 and -4 are based on the negative value of the given input. Four 3.3 V LVTTL signals of the PWM outputs are generated by the FPGA board and converted to 5V TTL signals by the interface board. Finally,

Current Controller/Current Controller



Figure 8. (a) Modified Current Controller Sub-model for the SIL and PWM Model for the HIL simulations and (b)/(c) the SIL and (d)/(e) the HIL simulation results.

the PCI in the RT-Sim receives the signals for the HIL simulation.

The results of the SIL simulations are shown in Figure 8 (b) and (c). Figure 8 (b) illustrates the comparison between the Simulink and SIL simulation results obtained on the RT-Sim. The results confirm that the RT-Sim accurately simulated the EV models and that the RT-Sim s-Functions operated properly in providing access to the shared memory region. Because of asynchronization at the beginning of the model execution between the shared memory region and models, the RT-Sim results are different within 0.1% of the Simulink simulation results seen in Figure 8 (c).

The HIL simulation results are illustrated in Figure 8 (d) and (e). According to the HIL simulation results, the EV model reacted two and half times faster than those of the SIL simulation. This result demonstrates the execution differences between SIL and HIL. In particular, HIL proves that the hardware is capable of generating faster signals for satisfying real time constrains than those generated by software.

| Table 1. Comparison of Simulink, SIL, and HIL |
|---|
| Simulation Results* |

| | Quarahoat | Overshoot | Simulation |
|----------|-----------|------------|------------|
| | Overshoot | Percentage | Difference |
| Simulink | 311 | 3.67% | |
| SIL | 312 | 4% | +0.33% |
| Simulink | 205.8 | 2.9% | |
| HIL | 206.4 | 3.2% | +0.3% |
| | (2.2.2.0) | | |

*: Speed references (RPM) for SIL and HIL simulations are chosen as 300 and 200 RPMs respectively.

Accuracies of Simulink, SIL, and HIL simulations are compared for the further evaluation since the simulation models on the RT-Sim are executed in parallel as opposed to operating in serial as the Simulink simulation. Table 1 summarizes the overshoots of the motor speed measured from three different simulations with ideal references (i.e., 300 rpm for SIL and 200 rpm for HIL). According to the results collected, both the SIL and HIL simulations, respectively, have 4 and 3.2% overshoots while Simulink simulation has 3.67% for 300 rpm and 2.9% overshoots for 200 rpm. Differences of the simulation results are about 0.3%. This results in the developed real-time simulation suite performs close enough for the EV evaluation.

3.3 Evaluation of the Real-Time Operation

The results of determining the access time of the shared memory region, PCI card, and PCI ADC are shown in Figure 9 (a). The required amount of time to access the RT-Sim shared memory region is an average of 0.1486 μ s. The time required access the hardware PCI card is an average of 1.7889 μ s, due to the bus used to access the PCI card. The ADC in the PCI card was identified as a critical path with an average of 13.2093 μ s delay, which limits the faster model execution for the HIL simulation. Figure 9 (b) shows that the resolution of the real-time clock on the RT-Sim is set to 10 μ s. Since the most reliable time to consistently operate at the same timing interval was identified as 50 μ s, the model execution.



Figure 9. RT-Sim (a) Operation Timing Analysis and (b) Clock Operation

4 Conclusions

An academic real-time simulation platform has been developed for rapid and intuitive management, accurate simulation, and cost-efficient real-time environment for research and education. The RT-Sim suite consists of a console module, a hardware module, and a real-time simulator. The console implements software/hardware model development with existing design tools integrated to in-house model configuration expansion, and simulation result management. The hardware module contains FPGA/microcontroller integrated with an interface board. The RT-Sim comprising of a quard-core with double thread computer and a PCI card equipped with Ethernet communication performs model compilation, real-time task scheduling/distribution/interface/ simulation/result collection, and console/hardware module communication. Subsystems of an electric vehicle are modeled and performed the SIL/HIL simulations for extensive evaluation. A Simulink model of an EV subsystem was successfully simulated and analyzed via the SIL simulation. The same EV model was modified for the HIL simulation. The real-time simulation methodology developed was evaluated for organizing real-time simulation procedures, managing simulation modules, and achieving simulation accuracy. The developed RT-Sim suite supports seamlessly expansion of various design tasks that span different disciplines in the academic environment.

5 References

[1] L. Cheng and Z. Lipeng, "Hardware-in-the-Loop Simulation and Its Application in Electric Vehicle Development," IEEE Vehicle Power and Propulsion Conference, 2008.

[2] C. Culianu and D.Christini, "Real-time Linux experiment interface system: RTLab," IEEE Bioengineering Conference, pp. 51 – 52, 2003.

[3] C. Dufour, J. Bélanger, T. Ishikawa, and K. Uemura, "Advances in Real-Time Simulation of Fuel Cell Hybrid Electric Vehicles," Proceedings of the 21st Electric Vehicle Symposium (EVS-21), April 2-6 2005, Monte Carlo, Monaco.

[4] J. Casteres and T. Ramaherirariny, "Aircraft integration real-time simulator Modeling with AADL for architecture tradeoffs," DATE09, 2009.

[5] G. Parma, and V. Dinahavi, "Real-Time Digital Hardware Simulation of Power Electronics and Drives," IEEE Transactions on Power Delivery, pp. 1235-1246, 2007.

[6] C. Dufour, S. Abourida, and J. Belanger, "Hardware-inthe-Loop Simulation of Power Drives with RT-LAB," IEEE PEDS pp. 1646-1651, 2005. [7] R. McNeal and M. Belkhayat, "Standard Tools for Hardware-in-the-Loop (HIL) Modeling and Simulation," IEEE pp. 130-137, 2007.

[8] O. Mohammed, N. Abed, and S. Ganu, "Real-Time Simulations of Electrical Machine Drives with Hardware-in-the-Loop," IEEE pp. 1-6, 2007.

[9] M. Faruque and V. Dinavahi, "Hardware-in-the-Loop Simulation of Power Electronic Systems Using Adaptive Discretization," IEEE Transactions on Industrial Electronics, Vol. 57, No. 4, pp. 1146 - 1158, 2010.

[10]P. Menghal and A. Jaya laxmi, "Real Time Simulation: A Novel Approach in Engineering Education," IEEE pp. 215-219, 2011.

[11]C. Dufour, C. Andrade and J. Bélanger, "Real-Time Simulation Technologies in Education: a Link to Modern Engineering Methods and Practices," the 11th Int Conf on Engineering and Technology Education, 2010.

[12]S. Abourida, C. Dufour, J. Bélanger, and V. Lapointe, "Real-Time, PC-Based Simulator of Electric Systems and Drives," Int. Conf. on Power Systems Transients, pp. 1–6, 2003.

[13]MathWorks, Simulink® Getting Started Guide R2012a, 1012.

[14]B. Cohen, Real Chip Design and Verification Using Verilog and VHDL, 2002.

[15] Xilinx, ISE 10.1 Quick Start Tutorial, 2008.

Efficient Pseudo-Random Numbers Generated from Any Probability Distribution

Clarence Lehman¹ and Adrienne Keen²

¹University of Minnesota, 123 Snyder Hall, 1474 Gortner Avenue, Saint Paul, MN 55108, USA ²London School of Hygiene and Tropical Medicine, Keppel Street, London WC1E 7HT, UK

"Truth is much too complicated to allow anything but approximations." —John von Neumann, 1947

Abstract—Microscale simulations and other applications in science, engineering, and commerce need an abundance of pseudo-random numbers drawn from non-classical probability distributions, including empirical distributions that may be incompletely known. Discrete-event simulations that assign random times to future events have further requirements, including numbers drawn from subsets of distributions to help establish initial conditions, or to deal with events that are partially complete. Fast methods are known for generating pseudo-random numbers accurately from arbitrary probability distributions, but those methods do not combine the full range of necessary algorithms outlined here. In this paper we provide techniques and computer code for practical high-speed generation of pseudo-random numbers from any continuous, discontinuous, or discrete probability distribution, reducing the need for approximation by standard probability functions. The techniques are designed for the kinds of scientific simulations presently emerging.

Keywords: non-uniform random numbers, uniform random numbers, pseudo-random numbers, probability distributions, numerical simulations

1. Introduction

Computer generated pseudo-random numbers are needed at every step in stochastic simulations—as well as to establish representative sets of initial conditions in deterministic simulations, to draw samples for statistical bootstrapping and other operations, to identify uncertainty in models by varying parameters, to randomize experimental designs, to develop test cases for commercial software, and for many other applications in science, industry, and art. They can arise in such quantities as to become a significant part of the total time for the computation itself.

For example, an emerging application arises in discreteevent simulation [1], where stochastically assigned times of future events must be determined in advance. When a simulated individual is born, the future time of death may be assigned from empirical probabilistic "life tables" for the year, geographic location, and other conditions of the individual being simulated. Initial conditions for the simulation can start with an empirically or hypothetically derived "age distribution." Subsets of the life tables (sub-distributions) are sampled to determine how long each individual will live, based on initial conditions. Moreover, when new individuals may enter the population as immigrants, sub-distribution sampling is also needed to assign future times of death. Examples and code that follows will illustrate these points.

In what follows we shall omit the prefix "pseudo" in "pseudo-random", it being understood that repeatable sequences of numbers generated by deterministic computer algorithms can be at best apparently random. The starting point for random number generation from a desired distribution is random number generation from the standard uniform distribution. That is, random numbers greater than equal to zero and less than one, all drawn with equal likelihood and with no correlation between any prescribed pair of numbers. This is difficult to achieve in practice, but much theory and effort have been dedicated to the problem, and a number of acceptable algorithms are known (e.g. [2]).

The question is, given a standard uniform random number, how can that be accurately converted to a random number from an arbitrary distribution? Answers were found in the earliest days of computers, as early as John von Neumann in 1945, and some of those methods are still in use [3]. The early methods, of necessity, used essentially no computer memory. However, with the vast computer memories now available, not using available memory is wasting it, and timefor-space tradeoffs that support arbitrary distributions have been published recently [3].

In this paper we (1) provide further background on the kinds of probability distributions needed in random number generation, (2) introduce a new aspect we call "sub-distribution sampling," (3) provide detailed algorithms for generating numbers from the kinds of distributions that arise in practice, and (4) compare processor time required of selected methods.

2. Density and cumulative functions

The "probability density function" is the most familiar representation, but it is the corresponding "cumulative prob-

ability distribution" that is employed for generating random numbers. In what follows, for brevity we shall write "density function" for the probability density function, "cumulative function" for the cumulative probability distribution. Also, we shall write "inverse function" or "inverse cumulative function" for the inverse of the cumulative probability distribution.

The vertical axis of a density function is the "probability density" that a corresponding value on the horizontal axis will be drawn by chance. This is not the actual probability of it being drawn, since in many distributions, the probability of any precise value being drawn is 0, amongst the infinite set of possible values. The probability density can take on any value from zero to infinity. If one considers the probability that a single random number drawn from the distribution will fall between two specified values, such as between 1.49 and 1.51, that probability is equal to the average value of the density function over that interval, multiplied by the width of the interval. Or what is equivalent, it is the area beneath the density function from the left to the right endpoint of the interval. Thus the area under the entire density function becomes 1. Modes of the distribution correspond to peaks in the density function, while the median and the mean are not immediately visible.

The cumulative function carries the same information as the corresponding density function, but in a different form. The vertical axis is the probability that a number less than or equal to the corresponding value on the horizontal axis will be drawn from the distribution. The vertical axis of a cumulative function is thus constrained to a range of 0 to 1. While the density function can have peaks and valleys, the cumulative function is either level or increasing. The cumulative function is one degree smoother than the corresponding density function, since it is its integral. The median of the distribution corresponds to the half-way point on the vertical axis (y = 0.5) and modes of the distribution correspond to places of maximum slope. The mean is not immediately visible in the cumulative function.

See Figures 1A through 1D for examples of four cumulative functions shown above their corresponding density functions.

3. The inverse cumulative technique

Using the inverse cumulative function to generate random numbers from any desired distribution is a long-recognized approach [4], and the idea is straightforward. Start with a uniform random number U between 0 and 1, locate that number on the vertical axis of the cumulative function, and find which value on the horizontal axis corresponds. That value is the desired random number. See the arrows in Figures 1A through 1D, where P = 0.25 on the vertical axis maps to -1, 4.25, 1.0638, and 2, respectively, on the horizontal axes of the various cumulative distributions. Note that domains from which a random number generator may

never select random values correspond to perfectly level stretches in the cumulative function (Figure 1B and 1D), and that discrete random numbers correspond to vertical jumps (Figure 1D). Also, as Devroye points out [4], using the same value of U like this for multiple distributions, or using correlated values $f(U\pm\rho_i)$, with ρ_i being a small random variate, draws correlated random numbers from multiple distributions. Likewise, using uniform random numbers equal to U and $f(1-U\pm\rho_i)$ will generate negatively correlated numbers from any distribution, or pair of distributions. Both kinds of correlation can be useful in applications.

The inverse technique is simple graphically but not necessarily numerically, for it involves computing the inverse function. A few classical distributions, such as the exponential, Cauchy, and Pareto, have inverse functions that can be written in terms of elementary functions [4] and therefore computed directly. In general, however, inverting an arbitrary cumulative distribution is computationally difficult, in which case generating the corresponding random numbers is slow.

Hörmann and Leydold [5] explained how to improve the speed by computing the inverse function only once when the simulation begins, then approximating it by interpolation as the simulation proceeds. This can be done by computing a series of x, y pairs from the cumulative function, then exchanging x and y and fitting the y, x pairs to obtain the inverse. Approximating the inverse in this way is reasonably fast and can be made as accurate as desired for many distributions.

In this paper we exhibit a variation that is simpler yet still fast, and that supports the full set of functions required. We approximate not the inverse cumulative function, but the cumulative function itself, with quadratic pieces that join smoothly, without sharp corners, where one piece ends and the next begins, or which approximate smoothness as accurately as we need. We can then invert the function piecewise as the simulation runs, since inverting quadratic functions involves only the quadratic formula. We use this because it directly supports "sub-distribution sampling," which is useful in general discrete-event simulations [1] and other micro-scale computations. It also supports efficient random numbers from discrete density functions (e.g. Poisson), step functions (e.g. empirical histograms), or piecewise linear (e.g. empirical function estimates). See Figure 1 for examples. The corresponding algorithms (appendix) are short and relatively simple.

The algorithms we present can be made to reproduce known distributions as accurately as desired, although in many cases such accuracy is superfluous. Classical distributions may be used because they are known and available, even though they may not closely approximate the distributions of interest. For instance, waiting times in an individual-based computer model may be selected from an exponential distribution that is used largely as a convenience, for correspondence with differential equation models. More-



Figure 1. Random number generation at four levels of continuity. Horizontal axes represent values of random variables. The vertical axis for cumulative functions represents the probability that the random variable is less than or equal to the corresponding value *x* on the horizontal-axis, and for density functions represents the average probability of a random number being in an arbitrarily small surrounding interval. (**A**) A standard Cauchy distribution, which does not converge to a mean. That renders it difficult to replicate purely with finite approximations. It represents the distribution of slopes that are associated with random angles. Other classical distributions, such as normal, lognormal, exponential, and chi-square, which do converge to means, are also higher order like this. (**B**) A hypothetical empirical distribution of two lobes developed from a piecewise linear density function, with each lobe equally likely. See Figure 2 for data structures of this example. Multi-modal distributions like this arise, for example, in carbon-14 dating, though those functions are typically more complicated than this illustration. (**C**) A hypothetical empirical distribution of failure rates of machine parts. Probability of failure is relatively high for new parts, drops to a minimum at intermediate ages, then rises again when parts get older. (**D**) A Poisson distribution with mean of 3. This discrete distribution represents frequencies of co-occurrence of random events with delta-function spikes.

| i: | 0 | 1 | 2 | 3 | 4 | 5 | 6 | 7 | 8 | 9 |
|-------|---|---|-------|-------|-------|-------|-------|-------|----|----|
| X[i]: | 0 | 1 | 3 | 5 | 8 | 10 | 11.75 | 15.75 | 18 | 20 |
| Y[i]: | 0 | 0 | .1111 | .3333 | .5000 | .5000 | .5729 | .9063 | 1 | 1 |
| Q[i]: | 0 | 0 | 1/9 | 1/9 | 0 | 0 | 1/12 | 1/12 | 0 | 0 |

Figure 2. Data structures. Three one-dimensional arrays, X[i], Y[i], and Q[i] carry the x values, the cumulative y values, and optionally the density y values, respectively, for a probability function. This example corresponds to the distribution of Figure 1B. In practice, arrays are often much larger than this illustration.

over, a simulation may benefit from using data directly, such as lifetime survivorship data available for multiple years and geographic areas, rather than approximating the data with standard distributions, then handling those distributions exactly. On top of all this, many empirical distributions are poorly known. Important distributions, such as the duration of infectiousness for certain asymptomatic diseases, may not be known to more than one digit of accuracy. Therefore, we are not recommending these algorithms as general solutions for all cases, but they can be particularly useful in practical cases where even the density function may be imperfectly known.

4. Data structures

Each cumulative function and, where needed, each corresponding density function, is defined by a set of two or three matched one-dimensional arrays that define the functions at selected points in their domains (Figure 2). These are constructed by the user and supplied to the algorithms, either to approximate classical distributions but more typically to represent distributions derived empirically.

The first array is X[i], which contains the x values that cover the range of random values to be generated. With entries in one-to-one correspondence, Y[i] contains the y values for the cumulative distribution, with Y[0] = 0 and $Y[i_{max}] = 1$. Optionally, Q[i] can be supplied, which carries the derivatives of Y[i]. That is, it carries the probability density function. When Q[i] is used, it is piecewise linear, making Y[i] piecewise quadratic. With a sufficient number of points, this can model any distribution. When Q[i] is not used, Y[i] is taken to be piecewise linear or piecewise constant and the implied Q[i] is a piecewise constant histogram. That corresponds to many empirical distributions, such as life tables.

5. Algorithms

The appendix defines four algorithms sufficient to accomplish the goals of this paper:

- 1. Cinverse: Evaluates inverse cumulative functions.
- 2. Cforward: Evaluates general functions.
- 3. Cdiscrete: Evaluates discrete functions.
- 4. Cintegral: Prepares cumulative functions.

Algorithm 1, *Cinverse*, accepts a probability, typically as a uniform random number, plus a starting point g, and returns a corresponding random number for the cumulative

distribution defined by X, Y, and optionally Q. It calls *Cforward* and *Cdiscrete* to accomplish its work. The starting point g is a minimum value for the random number returned. It is typically set to the minimum value of the distribution, X[0], but may be a greater value. In that case numbers are selected from the remainder of the distribution, starting at g, transformed as Z(x) = (Y(x+g) - Y(g))/(1 - Y(g)). This we call "sub-distribution sampling," which has such uses as initializing a population with individuals of various ages before the simulation begins, or handling immigrants to a population of random ages and projecting their remaining lifetimes. Note that "memoryless" distributions (the exponentials) are invariant under this transformation.

Algorithm 2, *Cforward*, determines the value y of a function at a specific point x, here used to obtain the value of the cumulative function at a random value x. It calls *Cdiscrete* to accomplish its work. The function is defined by tables X and Y of values for corresponding points, and optionally a table Q of derivatives of Y at each point in X. Values in tables X and Y are increasing. Passing Y[i] as a function of X[i] processes forward functions, while passing X[i] as a function of Y[i] processes their inverses. (Y need not be increasing if Q is not supplied, though processing non-cumulative functions is not a purpose of this subroutine.)

Algorithm 3, *Cdiscrete*, is where the process begins. For cumulative distributions of discrete density functions that have precisely one entry per non-negative integer, as in Figure 1D, the process is complete and this routine may be called directly. For non-discrete cumulative functions, *Cinverse* is called instead. That calls this routine to start and then handles any necessary inverse linear or quadratic interpolation. This routine is nothing more than a recursive binary search that processes ordered tables of n entries in time proportional to $\log_2 n$.

Algorithm 4, *Cintegral*, creates a cumulative function on demand, given an array of points X[i] and a density function Q[i]. It integrates using quadratic interpolation. It is not needed for linear interpolation or discrete distributions, where the cumulative distribution is simply the sum of the density function values, as in Figure 1D.

6. Timing

In timing tests, drawing 10^7 random values from a Poisson distribution with a mean of 3, as in Figure 1D, averaged 0.89 seconds on a 2.4 GHz processor for both a standard

iterative method [6] and for this method. For means greater than 3, the standard iterative method was slower. Generating numbers from discrete distributions like the Poisson is exact and is particularly fast because no interpolation is needed. In a continuous case, drawing 10^7 random values from a lognormal distribution required 1.14 seconds overall with the standard Box–Muller method and 1.07 seconds with this method, on the same processor. The 6% improvement in speed is not significant, but it is significant that a general method like this is competitive with the speed of classical custom methods. Timing tests showed that the recursive binary-search of Algorithm 3 could be speeded by about 10% by using an equivalent iterative algorithm, at the cost of a little greater complexity.

7. Discussion

Devroye [4] listed six factors for assessing general random number generation, (1) speed, (2) initialization time, (3) memory requirements, (4) portability, (5) generality, and (6) simplicity/readability. He pointed out that the sixth factor is the most neglected.

We find that his first factor, speed, is as important now as then. Despite the enormous increase in computer speeds, computers now labor under proportionally longer simulations. The second, initialization time, is less important, since it typically vanishes into the time for the simulation itself. Moreover, with large memories spaces that can be allocated, cumulative distributions may be precomputed and read from files, so initialization times become essentially zero. The third factor, memory, is largely irrelevant now except insofar as it increases complexity—thanks to the thenincomprehensible rise in computer memory sizes. The fourth factor, portability, is now easily had with careful coding, for which countless examples exist. Therefore, speed, generality, and simplicity remain as important factors.

The algorithms we present satisfy the simplicity factor well. They require under 60 lines of computer code altogether and are completely exhibited here, with code and accompanying meta-code. They are also fast, slightly outperforming even well-established methods like the Box–Muller algorithm for drawing numbers for standard distributions like the lognormal. Finally they are general, written to handle any continuous or discrete distributions for which the density function or cumulative function is known.

Modest increases in speed could be obtained in these algorithms at the cost of complexity, trading space for time by storing the inverse cumulative function as a direct-access table, with one entry per lattice point in the probability space Y. This would eliminate the binary search, though that is not slow. It would work if the slope of the cumulative function never gets very close to zero. A more modest increase in speed, whenever sub-distribution sampling applies, could result from rescaling the random number differently so that the entire table need not be searched, but only the part

Generality could be increased even further. As written, the algorithms handle all the kinds of distributions shown in Figure 1, but not mixtures of those types—for example, density functions that are discrete in some parts of the domain and continuous in other parts. Such functions are compatible with the algorithms detailed here and the algorithms could be extended to accommodate them, but at the cost of a little complexity, should the need for such hybrid distributions ever arise.

8. Conclusions

The algorithms presented here can be incorporated wherever efficient random numbers drawn from arbitrary distributions are needed. These algorithms have been successfully used in a large-scale simulation model developed by one of us (A.K.) for tuberculosis in the UK. Compilable copies of the code described here and related simulation algorithms are available free from the authors upon request.

9. Acknowledgements

We are grateful to Todd Lehman and Lori Thomson for discussions and help with the presentation, and to Todd Lehman for timing comparisons of binary-searching iteratively and recursively. This project was supported in part by a resident fellowship grant to C.L. from the UMN Institute on the Environment, by grants of computer time from the Minnesota Supercomputer Institute, and by doctoral research funding to A.K. from the Modelling and Economics Unit at the Health Protection Agency, London.

10. Contributions

A.K. conceived the approach to sub-distribution sampling, which was necessary in her discrete-event simulations, and which inspired this effort. C.L. extended the ideas to piecewise quadratic functions and coded the resulting algorithms. A.K. tested and applied the techniques in a large-scale individual-based model for tuberculosis. Both authors contributed to the manuscript.

References

- J. Banks, J. S. C. II, B. L. Nelson, and D. M. Nicol, "Discrete-event system simulation, fourth edition," *Pearson Prentice Hall, New Jersey*, 2005.
- [2] W. H. Press, S. A. Teukolsky, W. T. Vetterling, and B. P. Flannery, "Numerical recipes: The art of scientific computing, third edition," *Cambridge University Press, New York*, 2007.
- [3] W. Hörmann, J. Leydold, and G. Derflinger, "Automatic nonuniform random variate generation," Springer-Verlag, Berlin, 2004.
- [4] L. Devroye, "Non-uniform random variate generation," Springer-Verlag, New York, 1986.
- [5] W. Hörmann and J. Leydold, "Continuous random variate generation by fast numerical inversion," ACM Transactions on Modeling and Computer Simulation, vol. 13, pp. 347–362, 2003.
- [6] D. Knuth, "Seminumerical algorithms, volume 2, third edition," Addison-Wesley, Reading, MA, 1997.

11. Appendix

To use the algorithms described in this paper, it is only necessary to understand the entry and exit conditions that appear at the beginning of each, not the code itself. Nonetheless, to allow complete evaluation of the algorithms, and to encourage further development of them, we present them as pseudo-code inspired by and simplified from the programming languages C, R, and Python.

The algorithms are defined with sufficient precision that they can be run, tested, timed, modified, or translated to other languages. Familiarity with a relatively few operators* and

with the syntax of flow control (if, for, while, etc.), is sufficient to follow the algorithms. Text copies of this pseudo-code translated into operational C are available from the authors upon request, or from the associated website www.cbs.umn.edu/modeling.

The algorithms assume a uniform random number generator Rand, which returns values in the range of 0 to 1, including 0 but not including 1, as is typical for uniform random number generators. WarnMsg and ExitMsg display error messages and the latter terminates the program.

Algorithm 1. Evaluate inverse cumulative function, with sub-distribution sampling.

Upon entry to the algorithm, (1) k describes the piecewise order of the function: 0=constant, 1=linear, 2=quadratic. (2) y contains a value between 0 and 1, representing a probability. (3) g is given value in the range X[0] to X[n-1], inclusive. (4) n is the number of entries in tables X, Y, and Q. (5) X is a table of strictly increasing values in the set of numbers to be generated. (6) Y is a table of probabilities, each being the probability that a value will be less than or equal to the corresponding value in X. (7) Q is a table of probability densities, in effect the derivative of Y at every point in X. At exit, *Cinverse* returns the value from the given distribution corresponding to probability y, starting at value g. Note that if g > 0, this is the value from the rescaled distribution.

real Cinverse(k, y, g, n, X, Y, Q) integer k, n; real y, g, X[], Y[], Q[]; integer i; real r, s, d, h, a, b, c, p, w; if X[0] > g or X[n-1] < g: *ExitMsg*(1); if $Y[0] \neq 0$ or $Y[n-1] \neq 1$: ExitMsg(2); if y < Y[0]: return X[0]; if y > Y[n-1]: return X[n-1]; $v \rightarrow r; g - X[0] \rightarrow d;$ if d: Cforward(k, g, 0, n-1, X, Y, Q) $\rightarrow p$, $p + r * (1 - p) \rightarrow r;$ $Cdiscrete(Y, 0, n, r) \rightarrow i;$ if k = 0: return X[i] - d; $\begin{array}{ll} X[i+1] - X[i] & \rightarrow w; \\ \text{if } k = 2 \text{ and } Q: \ (Q[i+1] - Q[i])/(2*w) & \rightarrow a, \\ Q[i] & \rightarrow b, \ Y[i] - r & \rightarrow c; \end{array}$ else $0 \rightarrow a$; if a: $b*b-4*a*c \rightarrow s$; if s < 0: ExitMsg(3); $sqrt(s) \rightarrow s, \ (-b+s)/(2*a) \rightarrow h;$ if h < 0 or h > w: ExitMsg(4); else $Y[i+1] - Y[i] \rightarrow s;$ if s: $(r - Y[i])/s \rightarrow s$; else $1 \rightarrow s$; $s * w \rightarrow h$; return X[i] + h - d;

* The pseudo-code given here is two-dimensional, as in the language Python, so that indentation completely defines the nested structure, with no need for bracketing characters such as '{' and '}'. Variables and function names are italicized and flow control and reserved words are bolded.

The assignment operator is represented either as ' \leftarrow ' or ' \rightarrow ', similar to assignments in R. The compound assignments ' $a + 1 \rightarrow$ $a \rightarrow b \rightarrow W[i][j]$ and $W[i][j] \leftarrow b \leftarrow a \leftarrow a+1$ are equivalent, first incrementing a and placing the results back in a, then in b, and then in the i, jth element of the array W.

The expression structure 'c? u : v', where c is a condition, u is

- 1. Check the bounds of both tables.
- 2. Handle variables outside the normal range.
- 3. Rescale the probability value if only part of the distribution is to be sampled.
- 4. Bracket the probability value and return if it is piecewise constant.
- 5. If this is piecewise quadratic, generate the coefficients of the quadratic equation, $ax^2 + bx + c$.
- 6. If the equation actually has a quadratic term, invert it using the positive root of the quadratic formula.
- 7. If it is only linear, invert it with linear interpolation.

an if-expression, and v is an else-expression, follows that of C. Using up-tick and down-tick operators to write ' $\uparrow a$ ', ' $\downarrow a$ ', ' $a\uparrow$ ', and ' $a\downarrow$ ' form pre- and post-increments by one, as in '++a', '--a', 'a++', and 'a--' of C.

Arrays are indexed as in the language C, starting with 0. Data types are 'integer' and 'real', with the latter specifying floating point. Operator precedence is that of C, with assignments having lowest precedence. Logical operators such as 'and' and 'or' are preemptive, terminating a chain of logical operations as soon as the result is known. Permanent global assignments, as would be represented '#define $\alpha \beta$ ' in C, are rendered as ' $\alpha \equiv \beta$ '.

^{8.} Return the result.

Algorithm 2. Evaluate general function.

Upon entry to the algorithm, (1) k describes the piecewise order of the function: 0=constant, 1=linear, 2=quadratic. (2) x specifies the independent variable. (3) i0 and i1 define the first and last entries, respectively, in tables X, Y, and Q. (4) X and Y define the independent and dependent variables, respectively. (5) Q defines the derivative of the function, if k is 2. Otherwise Q is null. At exit, *Cforward* returns the value of the function at point x. If x is below or above the range defined in table X, the minimum or maximum value, respectively, in table Y is returned.

| real Cforward(k, x, i0, i1, X, Y, Q) integer k, i0, in integer i; real h, s, u, w; | I; real x, X[], Y[], Q[]; |
|---|---|
| if $k > 1$ and $Q = 0$: <i>ExitMsg</i> (5); | 1. Check for certain invalid calls. |
| if $x < X[i0]$: return $Y[i0]$; if $x > X[i1]$: return $Y[i1]$; | 2. Handle variables outside of the normal range. |
| $Cdiscrete(X, i0, i1 - i0 + 1, x) \rightarrow i;$ if $k = 0$: return $Y[i + 1];$ | 3. Bracket the independent variable and return if piecewise constant. |
| $X[i+1] - X[i] \rightarrow w, \ x - X[i] \rightarrow u;$ | 4. Compute <i>x</i> -width and displacement. |
| if $k = 2$: $(Q[i+1] - Q[i])/w \rightarrow s,$ $u*(Q[i] + u*s/2) \rightarrow h;$ | 5. If a derivative is supplied, compute the <i>y</i> -value with quadratic interpolation. |
| else if $w: u/w \to w$; else $1 \to w$; $w*(Y[i+1] - Y[i]) \to h$; | 6. Otherwise interpolate linearly. |
| return $Y[i] + h;$ | 7. Return the computed y -value. |

Algorithm 3. Evaluate discrete function.

Upon entry to the algorithm, (1) T addresses a strictly increasing table of two or more values. (2) b indexes the beginning entry to be examined in T. (3) n defines the number of entries to be examined in T, at least 2. (4) v specifies the value to be located in T, with $T[b] \le v \le T[b+n-1]$. At exit, Cdiscrete indexes the local pair of table entries containing v, such that $T[loc] \leq v \leq \overline{T}[loc + 1]$.

integer Cdiscrete(T, b, n, v) integer b, n; real T[], v; integer m; $(n+1)/2 - 1 \to m;$ **return** $m \le 0$? b: v < T[b+m]? Cdiscrete(T, b, m+1, v): Cdiscrete(T, b+m, n-m, v);

Algorithm 4. Prepare cumulative function.

Upon entry to the algorithm, (1) n is the number of entries in tables X, Y, and Q. (2) X is a table of strictly increasing values in the set of random numbers to be generated. (3) Y is a table to receive the cumulative function associated with the corresponding values in X. (4) Q is a table representing the piecewise linear density function associated with the corresponding values in X. At exit, (1) Cintegral returns the number of entries in Y up to and including the first entry that saturates at 1. (2) Y contains the piecewise quadratic cumulative distribution function associated with the corresponding values in X.

| integer Cintegral(n, X, Y, Q) integer n; real X[], Y[], y integer i, m; real w; | Q[]; |
|---|--|
| for <i>i</i> from 1 to $n-1$: if $X[i-1] \ge X[i]$: <i>ExitMsg</i> (6); | 1. Make sure the domain increasing. |
| $\begin{array}{ll} 0 \ \rightarrow Y[0], \ 1 \ \rightarrow m; \\ \text{for } i \ \text{from } 1 \ \text{to } n-1; \\ Y[i-1] \ \rightarrow Y[i], \ X[i] - X[i-1] \ \rightarrow w; \\ (Q[i-1] + (Q[i] - Q[i-1])/2) * w + Y[i] \ \rightarrow Y[i]; \\ \text{if } Y[i] > 1; \ WarnMsg(7), \ 1 \ \rightarrow Y[i]; \\ \text{if } Y[i] < 1; \ i+1 \ \rightarrow m; \end{array}$ | 2. Integrate the probability function to obtain the c distribution function. |
| if $Y[n-1] \neq 1$: WarnMsg(8); return $m + 1$; | 3. Make sure it adds and number of operational of |

- is strictly
- density umulative
- eturn the entries.

Public Concerns and Projected Repercussions of Health Care Reform Using Simulation

Ahmed YoussefAgha, PhD Dept. of Applied Health Science Indiana University Bloomington, Indiana, USA E-mail: ahmyouss@indiana.edu

Wasantha Jayawardene, MD Dept. of Applied Health Science Indiana University Bloomington, Indiana, USA E-mail: wajayawa@indiana.edu

Abstract: Health Care Reform (HCR) is currently a major concern of many Americans. The White House outlined as three salient issues which are reflective of the fears of the American public. The first dimension, stability and security, refers to the increased dependability of health care and the reduction of discrimination against individuals with health conditions. Dimension two affects individuals who do not have health insurance and pertains to the quality and affordable choices they will gain under the new reform. Dimension three refers to funding concerns and future fiscal management. Purpose: The purpose of this paper is to examine the interrelationships of the three dimensions of public concern regarding health care reform. The proposed research will compare the perspectives of a simulation and actual Media data regarding the Health Care Reform before it was enacted. Method: We used stochastic tree simulation to study the interrelationships of predicted repercussions of the three dimensions for those subjects with multiple diseases, a common condition among 65+ Age groups or obese groups with several chronic illnesses. Results: Analysis revealed that improved quality of coverage and greater security, dimensions of the Affordable Health Care Act, are associated with stable health insurance, leading to longer periods of QALY health. Conclusions: For those subjects on multi disease, this work simulated the expected utility related to the HCR dimensions interaction. At quality/ affordability of (20-50%), the simulation showed that adding more funding within the third dimension for (20%), given that the risk of no security/ stability is within (0-to-20%) is the more reasonable situation. Better health care policies will increase the number of Americans with health insurance and improve overall public health in the U.S.

I. INTRODUCTION

A. Healthcare Reform Significance

Since 1992, the American public has considered health care one of the four most important problems facing the country [1]. Analysis of public opinion polls

David Lohrmann, PhD Dept. of Applied Health Science Indiana University Bloomington, Indiana, USA E-mail: dlohrman@indiana.edu

Lesa Lorenzen-Huber, PhD Dept. of Applied Health Science Indiana University Bloomington, Indiana, USA E-mail: lehuber@indiana.edu

shows that the majority of Americans are satisfied with the quality of their health care, but are not satisfied with the cost [1, 2]. Although the majority of Americans agree that reforms are needed to control the costs of health care how these reforms should be realized is still controversial [3-5]. Blendon & Benson [2] suggest that the inability of the public to agree on a strategy of health care reform has significantly inhibited the political progress of the current reform bill.

Health care expenditures in the United States are four times greater than for national defense and the average health insurance premiums and individual contributions for family coverage have increased approximately 120% during last 10 years [6]. The U.S. health care system has been blamed for inefficiencies, excessive administrative expenses, inflated prices, inappropriate waste, fraud and abuse. While many people lack health insurance, others who have insurance allegedly receive care ranging from inexcusably poor to very high quality. President Obama and congressional Democrats focused on creating a national health care system to address some of these problems with the goal of improving efficiency, restraining expenses, and increasing quality in healthcare, prompting application of business practices to medicine.

In their analysis of health care in the United States with a focus on saving money, many methodologists, policy makers, and the public seemed to dismiss the major disadvantages of other national health care systems and the previous experiences of health care reform in the United States. A primary reason is that many of those without health insurance do not receive preventative care, and are afraid to seek medical attention when sick or injured. As a result, they often suffer and wait until it is absolutely vital they seek help at a hospital emergency room which may cost them an exorbitant about of money. The Affordable Care Act will stop this pattern by making routine care affordable.

Media and publications by stakeholder organizations are trying to convince the general public to take their side in support or opposition to the Affordable Care Act [7] and a heated debate has ensued at all levels. Most of the assessments of public opinion, including opinions of older adults, are conducted by media companies, such as television channels, magazines and newspapers, which nave not necessarily conveyed accurate information in a statically unbiased format. Often these companies publish information that reflects their biases. Scientifically designed research to assess and analyze the perspective of older adults is imperative to gain a more comprehensive understanding of the views of this population.

B. Age 65+ Groups Situation

In 1993-94, opposition to the Clinton healthcare reform centered on middle-aged voters. In contrast, opposition to the ObamaCare legislation centers on the elderly who suspect that it will bring a significant restriction of their medical care.

Oddly, though this population that now receives health services through government-run Medicare, many would prefer a privately run system (67%) to one managed by the government (7%). Most resistant to change, the elderly voters cited fears that they will have to change existing healthcare arrangements as the greatest reason they opposed ObamaCare. The political implications of this reaction are enormous because such dissatisfaction has sent many Americans into opposition, similar to the Social Security reforms in 2005. Many older adults may strongly oppose the proposed health care reform in the 2012 elections.

The concerns of older adults are personal concerns related to their access to health care and whether they would have to ration visits to health care providers, rather than the macro perspectives that account for the deficit and taxes that influence the implementation of this new plan. Open-ended questions on a number of surveys find the elderly very worried that they will not be able to get quality-of-life treatments, such as hip or knee replacements, under the Obama program. Others worry that the program will encourage them to give up when facing serious illness to minimize government costs (Table 1).

Table 1: Kaiser Health Tracking Poll (January 2010)

| | | Age | |
|--|-------|-------|-----|
| | 18-39 | 40-64 | 65+ |
| Support/Oppose Current Reform Proposals | | | |
| Support proposals being discussed in Congress | 47% | 41% | 37% |
| Oppose proposals being discussed | 33% | 46% | 48% |
| No opinion | 20% | 13% | 15% |
| Awareness of Medicare provisions Aware proposals would limit future increases in Medicare | | | |
| provider payments | 54% | 46% | 37% |
| Aware proposals would help close the Medicare doughnut hole | 53% | 38% | 37% |
| Limiting future increases in Medicare provider payments | | | |
| Make you more likely to support proposals | 53% | 43% | 24% |
| Make you less likely to support | 29% | 37% | 46% |

Source: Kaiser Health Tracking Poll (January 2010)

Elderly Americans seemingly turned out to fight the Obama health plan, and President Obama told them they have nothing to worry about. While claims about euthanasia and "death panels" were over the top, senior fears exposed a fundamental truth about the Obama approach: namely, once health care is implemented rationing of care is inevitable, and those who have lived the longest will find their care the most restricted. Far from being a scare tactic, this is a logical conclusion based on experience and common-sense. Once health care is a "free good" that government pays for, demand will rise as well. When the public finally reaches its taxing limit, health care providers will have to limit care in order to restrain spending and break even. In other words, care will be rationed. Obama's response is that private insurance companies already ration services by limiting which treatments are covered and which are not.

For knowledge of the efficacy, strengths, and limitations of national health care, much can be learned from the systems implemented in other countries. Virtually every European government with "universal" health care restricts access in one way or another to control costs. The British system is most restrictive, using a formula known as "quality-adjusted life years" (QALYs), that determines who can receive what care. If a treatment is not deemed to be cost-effective for specific populations, particularly the elderly, the National Health Service simply does not pay for it. Even France, which has a mix of public and private medicine, strictly controls the use of specialists and the introduction of new medical technologies such as CT scans and MRIs. Medicare already rations care by refusing, for example, to pay for virtual colonoscopies, and has payment policies or directives to restrict the use of certain cancer drugs, diagnostic tools, asthma medications and many others. Seniors routinely buy supplemental insurance (Medigap) to patch Medicare's holes. This is not trivial as a substantial portion of Medicare spending is incurred in the last six months of life [8].

This study aimed to explore the interaction of the 3Ds in the proposed reform for uninsured adults in a multi-illness condition and those aged 65+. The research question is: What is the best interaction situation for uninsured individuals with multi-illness and those age 65+ regarding the three dimensions of security/stability, quality/affordability, and funding.

II. METHODOLOGY

A. Decision Trees

A decision tree is an essential element of decision analysis under uncertainty [9-11]. Conventional decision trees do not completely represent "the real world since they cannot investigate problems that are cyclic in nature" [10]. As a result and for enhancing decision making, researchers have developed methodologies, via Markov processes, allowing for the examination of cyclic diseases. Recently, stochastic trees as a technique has been developed for solving continuous-time Markov cyclic trees by Hazen [12-14]. Hazen's approach enables "state-time to be modeled as a range of stay where subject-or-patient's health state transitions can occur at any time; it can in addition contain quality of patient's life and patients' preferences regarding intervention risk"; Hazen [13]. Hazen's Stochastic Tree measures treatment outcome, such as Cost or Quality-Adjusted Life Years (QALYs) of health gain to an individual facing recurrent disease when he/she is under a treatment/medication. Hazen's model was developed as a technique for solving continuous-time Markov cycle trees. This stochastic tree uses recursive (rollback) evaluation of outcome function; this is a computation of an expected outcome measure use iterative technique similar to the technique of value iteration in the Dynamic Programming [13, 14]. Rollback computation has defined as: a predetermined sequence of health status y_i and duration t_i where it is understood that status y_1 is first occupied for duration t_1 , followed by a status y_2 for duration t_2 , and so on. The recursive equation, defines a utility function u(y) over the subject health status transition and history [14]:

B. Computing Expected Utility

If there is a net of nodes (health status) on which a subject will move through until he/she has death; such a net can be represented by stochastic tree model equivalence to Hazen's model (1998-2002) [14].

Effectively, the model in this study, utilized simulation; it describes the interaction among the HCR plan's dimensions (The three dimensions have been displayed by the White House website to explain the Healthcare Reform Plan: www.whitehouse.gov/assets/documents/obama_plan_card.PDF

The three dimensions (3Ds) had been used by the researchers as in the following:

- Stability/Security renamed as d1 to represent the Risk of No Health-Insurance, 0 < d₁ ≤ 1.
- Quality-Affordability renamed as d1 to represent "If you don't have insurance: Quality, Affordable Choices for All Americans", $0 < d_2 \le 1$
- Percent Added Funding (d3); which represents funding arguments, $1 \le d_3 \le 2$

This study model is analogues to Hazen model (1998-2003) [14]. Hazen model, computes the expected quality adjusted life years (QALY); this analogy is to compute the expected QALY within which an individual will have a value adjusted to the healthcare d1, d2, and d3. The adjustment will be attuned to three interacting healthcare dimensions. The QALY values will be greater or equal to zero. The higher value of QALY the better outcome and the longer life within which the subject gains the better 3Ds interaction; his/her life is computed under a situation of the HCR dimensions interrelations.

$$u(node_{y}) = \int_{0}^{\infty} (\int_{0}^{t} q(y)e^{-\alpha^{*}k}dk) * \lambda e^{-\lambda^{*}t}dt$$

(y) is subject's health quality to each small interval *dk* spent in status *y*, $1 \le q \le 0$. Alpha (α) is subject's risk assessment

factor (under the simulation), $-1 \le \alpha \le 1$. The expression $e^{-\alpha^*k}$ represents the subject tendency to face the risk of intervention at specific time: if alpha is positive, then the subject is risk averse at status y, whereas if alpha is negative, the subject is risk seeking. If alpha equals to zero, the subject is risk neutral; and also "alpha may be interpreted as a discount rate of the associated cost of staying in status y." [13, 14]. The λ is the time rate to be in health status y.

$$u(node_{y}) = \int_{0}^{\infty} (\int_{0}^{t} q(y)e^{-(d1/d3)^{*}k}dk) * (\lambda * d2)e^{-(\lambda * d2)^{*}t}dt$$

is the analogues of Hazen's. The internal part of the formula models the subjects' status under the interaction of d_1 and d_3 in the healthcare system:

$$(\int_{0}^{t} q(y)e^{-(d1/d3)^{*k}}dk)$$

Once this value reaches zero a death has occurred. The Alpha (α) parameter in Hazen's model was replaced with (d_1/d_3). The two factors: d_1 and d_3 interact together so that the funding constant, d_3 , controls the risk of no insurance, d_1 , for the subject who receives healthcare benefits, (e.g. d_1 may be 0.40 risk of no health insurance, and d_3 may be 1.20; which represents 0.20 of an added finance to the healthcare budget to reduce d_1 to 0.33).

The second part of the formula: $(\lambda * d2) * e^{-(\lambda * d2)*t} dt$

It represents the amount of time spent at health status (y) with rate λ , under specific (d₂) quality/affordability benefits. The λ in Hazen's model was replaced with (λ *d₂). λ and d₂ are the factors that give rise to greater length of stay in a health status y, $0 \le \lambda \le 1$ and $0 \le d_2 \le 1$. Higher values of d₂ indicate greater levels of "added quality/affordability". At a specific health status y, the greater value of λ indicates the higher rate of y under specific quality/affordability benefits generated by healthcare resources, such as hospital. The quality/affordability dimension, d₂, interacts with λ for the sake of the subject who receives healthcare services.



Figure 1: Stochastic Tree, to simulate the interaction among the three dimensions of the Healthcare Reform Plan

At each node, the tree computes QALY, which is *u*; the root node "Well" will have the summed QALY of all the tree nodes. Transition from node to node is by continuous rate $0 \le \lambda_i \le 1$; or by chance, p, $0 \le p_i \le 1$. The next three formulas are node-based formulas (a node means a health state); used to evaluate the expected value of QALY at each node until the QALY values stabilized (Semi-Monrovian chain); if stabilize after a number of rolling back computations; then the tree's root will contain the final QALY in the Hazen model.

$$E[u(node_{y})] = \frac{q(y) + d2 * \sum_{i=1}^{n} \lambda_{i} E[u(node_{i})]}{(d1/d3) + (\lambda * d2)}$$
$$E[u(node_{y})] = q(y) + \sum_{i=1}^{n} p_{i} E[u(node_{i})]$$
$$E[u(node_{y})] = \frac{q(y) + (\lambda * d2) \sum_{i=1}^{n} p_{i} E[u(node_{i})]}{(d1/d3) + (\lambda * d2)}$$

Hence: u(node y) means the QALY at node y; E[u(node g)] means the expected QALY; and q(y) means quality of life at health status y, $0 < q(y) \le 1$. $(d_1/d_3 + \lambda * d_2) > 0$; and $\lambda = sum(\lambda_i)$. The first formula used for node branching with rates (λ_i) ; the second for node branching with probabilities (p_i).

C. Simulation - Matlab Code

SubjectO = 1.0; % Normal Health Quality of a Subject ChronicIllnessQ = 1.0; % Health Status Chronic illness OtherillnessP = 1.0; % Health Status under Other illness % Subject Risk of No Insurance d1s = 0:.1:1;FUND = [1.2 1.4 1.8]; d2 = .50; % Quality/Affordability % death rate as illness deathR1 = .01; deathR2 = .03; %death rate as other p1 = .25; % Chronic illness probability p2 = .50; % probability of death lambda1 = 1.0; %Rate level lambda2 = deathR1 + deathR2;lambda3 = deathR1 + deathR2: lambda4 = 1.0;Well = 0; IllnessA = 0; IllnessB = 0; IllnessC = 0; IllnessD = 0; ChronicIllnessQ = 0; for iii = 1:length(FUND) d3 = 1/FUND(iii);for ii = 1:length(d1s) d1 = d1s(ii);for i = 2:100 IllnessD(i) = (otherillnessP +(lambda3*d2*0+lambda4*d2*ChronicIllnessQ(i-1)))/ ((d1*d3)+(d2*(lambda3+lambda4)));IllnessC(i) = (otherillnessP +(lambda2*d2*0+lambda4*d2*ChronicIllnessQ(i-1)))/ ((d1*d3)+(d2*(lambda2+lambda4)));ChronicIllnessQ(i) = (ChronicIllnessQ+(p2*0+(1-p2)*IllnessD(i-1))); IIInessB(i) = 1*IIInessC(i-1);IllnessA(i) = otherillnessP + (p1*ChronicIllnessQ + (1-p1)*IllnessB(i-1));(d2*(lambda2+lambda1))); end dWell = diff(Well): stablepoint = (find(dWell==0)); stablepoints = stablepoint(find(stablepoint>5)); QALY(ii) = Well(min(stablepoints)+1); end QALYfinal(:,iii) = QALY; end

III. RESULTS

Results show that the primary concern reflected pertains to funding for the health care reform; however, increased stability and security were discussed as favorable qualities.



Figure 2: More quality/affordability has no increase effect after the 20% limit.



Figure 3: Zero risk of "no security and stability maximize QALY.



Figure 4: The interaction among the 3Ds - The QALY value at (d_1 =0.01, d_2 =0.01) represents the highest QALY; at which the subject has about zero value risk of "no security and stability", and has about zero "quality /affordability". This is logic: if he/she is 100% covered/insured under (d_1); then he/she does not need any choices under (d_2) "Quality/Affordability".

Table 2: Averages of QALY on d_1 vs. d_3 , at d_2 set to 50%

| | Adap | tive Fund Ctr | :l (d3) |
|----------------------------------|-------------------|-------------------|-------------------|
| Risk of No Security/Stability | (20%) Increase | (50%) Increase | (80%) Increase |
| d1 = 0-to-20% | 3.9 | 4.0 | 4.2 |
| d1 > 20% | 1.4 | 1.5 | 1.8 |

Note: Each cell value includes the average of QALY in Years

IV. DISCUSSION

For those subjects on multi disease such as the situation among the 65+ groups or obese groups with several chronic illnesses, this work simulated the expected utility related to the Health Care Reform dimensions interaction.

Table 2 shows that adding funding for 20%, given that the risk of no security stability is within (0-to-20%) is the more reasonable situation; it has not significant difference (p<.05) comparing to the other d_3 status at the same level of d_1 "the same row cells".

The importance of Hazen's Stochastic Tree was mentioned in the technical report, by Professor Keefer et al [15]. Keefer stated that the work of Hazen, stochastic trees, "is one of applications that presented significant advancement in decision analysis methodological tools" [15].

This study will further the work of government in the community due to the focus on public health and strategic planning required to maintain affordable, accessible healthcare in USA.

The federal government needs to develop an effective plan to address reform funding to appease Americans concerns. Better health care policies will increase the number of Americans with health insurance and improve overall public health in the U.S.

V. CONCLUSION

For those subjects on multi disease such as the situation among the 65+ groups or obese groups with several chronic illnesses, this work simulated the expected utility related to the Health Care Reform dimensions interaction. At quality/ affordability of (20-50%), the simulation showed that adding more funding within the third dimension for (20%), given that the risk of no security/stability is within (0-to-20%) is the more reasonable situation. Better health care policies will increase the number of Americans with health insurance and improve overall public health in the U.S.

REFERENCES

- L. R. Jacobs, "1994 all over again? Public opinion and health care," *New England Journal of Medicine*, vol. 358, pp. 1881-1883, May 1 2008.
- [2] R. J. Blendon and J. M. Benson, "Understanding How Americans View Health Care Reform," *New England Journal of Medicine*, vol. 361, pp. E13-U18, Aug 27 2009.
- [3] A. S. Chen and M. Weir, "The Long Shadow of the Past: Risk Pooling and the Political Development of Health Care Reform in the States," *Journal of Health Politics Policy and Law*, vol. 34, pp. 679-716, Oct 2009.

- [4] A. Gelman, D. Lee, and Y. Ghitza, "Public Opinion on Health Care Reform," *Forum-a Journal of Applied Research in Contemporary Politics*, vol. 8, 2010 2010.
- [5] S. E. Gollust, P. M. Lantz, and P. A. Ubel, "The Polarizing Effect of News Media Messages About the Social Determinants of Health," *American Journal of Public Health*, vol. 99, pp. 2160-2167, Dec 2009.
- [6] L. Manchikanti and J. A. Hirsch, "Obama Health Care for All Americans: Practical Implications," *Pain Physician*, vol. 12, pp. 289-304, Mar-Apr 2009.
- [7] D. R. McCanne, "THE ORGANIZATION FOR ECONOMIC COOPERATION AND DEVELOPMENT AND HEALTH CARE REFORM IN THE UNITED STATES," *International Journal of Health Services*, vol. 39, pp. 699-704, 2009 2009.
- [8] J. M. LUCE and G. D. RUBENFELD, "Can Health Care Costs Be Reduced by Limiting Intensive Care at the End of Life?," *American Journal of Respiratory and Critical Care Medicine*, vol. 165, pp. 750-754, March 15, 2002 2002.
- [9] H. Raiffa, Decision Analysis: Introductory Lectures on Choice Under Uncertainty. New York: McGraw-Hill, 1997.
- [10] A. Jaafari and J. Yao, "Combining real options and decision tree: An integrative approach for project investment decisions and risk management," *Journal of Structured and Project Finance*, vol. 9, p. 53, 2003.
- [11] T. A. Press, "Internet users lured to oppose health bills," ed, 2009.
- [12] G. B. Hazen, "STOCHASTIC TREES A NEW TECHNIQUE FOR TEMPORAL MEDICAL DECISION MODELING," *Medical Decision Making*, vol. 12, pp. 163-178, Jul-Sep 1992.
- [13] G. B. Hazen and J. M. Pellissier, "Recursive utility for stochastic trees," *Operations Research*, vol. 44, pp. 788-809, Sep-Oct 1996.
- [14] G. B. Hazen, J. M. Pellissier, and J. Sounderpandian, "Stochatic-Tree Models in Medical Decision Making," *Interface*, vol. 28, pp. 64-80, 1998.
- [15] D. Keefer, "Summary of Decision Analysis Applications in the Operations Research Literature, 1990-2001," Arizona State University, Technical Report2002.

Using Serious Games to Teach Business Process Modeling and Simulation

Cláudia Ribeiro, João Fernandes, André Lourenço, José Borbinha and João Pereira INESC-ID, Rua Alves Redol 9, Lisbon Portugal Department of Information Systems and Computer Science, IST/UTL, Lisbon, Portugal

Abstract - Serious game and virtual-based environments have recently been used in business contexts to promote training of business related competences as is the case of business process modeling. Due to its characteristics modern BPM methods rely heavily on computers, since they are used as a tool for model, simulate and analyze a process and its effects on the organizations. Several authors argue that these approaches offer several limitations with respect to the modeling processes namely, they typically offer a low degree of collaboration between users, the simulations are simplistic and do not take into account the possible interaction of users during the simulation process and the interfaces are many time complex and abstract requiring an extensive training in order to use those tools. Serious games and virtual environments are mentioned as a promising approach for teaching and simulating BPM since due to their characteristics they can solve or minimize many of those limitations. In this paper we describe a serious game named ImPROVE for teaching business process modeling in a realworld context, the Manchester triage system. The simulation laver of the game was based on Time-based Activity-based Costing which also provides feedback information related to costing of activities in a business process.

Keywords: business process modeling, Time-based Activitybased Costing, Simulation, Serious Games

1 Introduction

Simulation is the process of designing a model of a real or imagined system and conducting experiments with that model. The purpose of simulation experiments is to understand the behavior of the system or evaluate strategies for operation of the system (Smith, 2009). Enterprises design and implement systems to fulfill certain functions. What often becomes a challenge is to identify the complex interactions and interfaces between different organization functions and the roles and responsibilities of the various stakeholders involved. Business Process Modeling is an important part of understanding and restructuring the activities and information a typical enterprise uses to achieve its business goals. It is also a powerful method used for better understanding business concerns and communication between stakeholders, allowing that every interested part is actively involved in these activities. Modeling and simulation are tools and methods that are widely used in enterprise engineering/organizational study where these are considered effective, efficient and economical for organizational analysis and design (Barjis & Verbraeck, 2010).

Games and game-like tools have been chosen as the means of information spreading and training for various reasons. Some of these include the visually impressive outlook, apparent easiness of the game user interface compared to conventional software, the way younger generations have been accustomed to the playing of games, and entertaining aspects of the gaming and the attractiveness of the gaming (Beck, 2004) (Prensky, 2001). The use of serious games in business contexts has grown dramatically over the last fourth decades and it has been reported that they provide three main benefits in promoting organizational learning, namely: (i) to orient and train new employees; (ii) to select current managers or future managers; and (iii) for ongoing management training (A.J. Faria, 2009).

In serious games human actors are the active decision makers whose actions will affect the future state of the simulation run. This active role of human players separates simulation games from pure simulations (Martin, 2000). Pure simulations use static rules to calculate the outcome of the situation. Human players are unpredictable therefore they might choose to do things that extend the possibilities taken into account by the original designer. This is one of the central problems with which modeling and simulation tools still debate nowadays. There is a considerable gap between the model of a sub-set of the reality, the model a current simulation tool support and the ability to view the results of simulating such a model in the real-world.

In this paper we argue that using serious games and their specific characteristics for modeling and simulating business processes could help bridge these existing gaps and therefore promote suitable tools both for promoting organizational learning as well as organizational change. For that purpose a serious game based on modeling and simulation Time-Driven Activity-Based Costing (TDABC) processes was developed and applied to a hospital emergency unit case study.

This paper begins with an overview of the current developments of serious games (section 2), namely their advantages and current use in educational contexts and some examples of how they have been applied in business process modeling contexts. This is followed (section 3) by a brief description of the TDABC methodology and the advantages that serious games could provide when applying this methodology. Next (section 4), a detailed description of the ImPROVE game is provided. This game was developed in the context of a master class lecture and inspired by a realworld case, namely the current triage system of a Portuguese hospital emergency unit. Finally (section 5) conclusions are drawn pointing to the future scope for development that lies ahead in the vast and interesting field.

2 Serious Games and M&S

Serious game and virtual-based environments are an important response from the education technologist to the "digital natives" (Squire, 2005), a generation of students who were raised on interactive games and expect the same kind of interactive experiences from their education media. Indeed, it may possibly be wrong to call the use of serious games in education a novelty, since by nature young children begin to learn through games at their earliest years (Rieber, 1996).

Due to their characteristics, games can introduce clear advantages in supporting complex learning processes and knowledge transfer. Through games it is possible to simulate environments and systems allowing learners to experience situations that are impossible in the real world for reasons of safety, cost and time (Corti, 2007) (Jenkins, 2004).

Games and virtual simulations are often referred in the literature as experiential exercises (Gredler, 1994) in which people "*learn by doing*" avoiding "*mimicry learning*" (Turkle, 1984). They in fact increase greatly the learning outcomes since they easily change our mood towards the learning of specific topics (McGonigal, 2012). As Savill-Smith argue (Savill-Smith, 2004), games can support the development of a number of different competences such as: analytical and spatial skills, strategic skills and insight learning, recollection capabilities, psychomotor skills, visual selective attention, showing promising results when compared to traditional methods (Szezurek, 1982) (VanSickle, 1986) (Randel, 1992) (Van Eck, 2006).

The use of games in an education context is not a distant concept, they were adopted long time ago by organizations in a wide range of sectors, nowadays no pilot will ever pilot an airplane without an intensive training using simulation and games, no power plant manager will ever run a nuclear plant without an intensive training using simulation and games (Aldrich, 2004). We can thus conclude that, the organization that cares most about the training use simulations and games to do their training and to support some of their day-to-day processes.

The field of business is not an exception in the permeation to this kind of approach, the first business game widely recognized, 'Top Management Decision Simulation', was created in 1956 by the American Management Association. Since then, a great number of different business games have been developed (Carroll, 1954) (Faria, 2004) and used in management training by different business schools, faculties and professional associations all over the world (Walters, 1997) (Chang, 2003) (Sánchez Franco, 2009). Nowadays these games are seen as a useful tool to learn how to manage firms and to explore new strategic opportunities (Jensen, 2003).

Several authors referred that the most important advantages of applying games in a business context, are the immediate feedback, active participation of students, learning from the experience, observation of the key factors in an on-the-job situation, preparation for the uncertainty of business, and the high motivation to learn created by the competitive environment (Fu, 2009) (Gilgeous, 1996) (Zantow, 2005).

An activity that would benefit substantially by the introduction of these approaches is the Business Process Modeling (BPM). Due to its characteristics modern BPM methods rely heavily on computers, since they are used as a tool for model, simulate and analyze a process and its effects on the organizations (Tapani, 2008). A wide range of process modeling grammars (also called notation and technics) were proposed during recent years (Rosemann, 2006), generally they are composed by two dimension representation making use of shapes like circles, squares, rectangles. The authoring of such schemes is generally supported by process modeling tool suites (Hill, 1990), that provide a graphical model editor and sometimes some complementary functionality regarding simulation, bug detection, reporting and analysis (Recker, 2006).

Several authors argue that these approaches offer several limitations with respect to the modeling processes, they typically offer a low degree of collaboration between users, the simulations are simplistic and do not take into account the possible interaction of users during the simulation process and the interfaces are many time complex and abstract requiring an extensive training in order to use those tools (West, 2010). Serious-games and virtual environments are mentioned as a promising approach for teaching and simulating BPM since due to their characteristics they can solve or minimize many of those limitations.

Towards meeting these objectives, a large number of three-dimensional approaches for the representation of the process models have recently been suggested (Brown, 2009) (Streit, 2005). More than merely providing a modulation and simulation tool, these environments put humans as active decision makers, conditioning their actions states of the simulation runs, thus differing from the traditional tools. Due to their superior collaboration capabilities, games and virtual environments can enhance the process of distance collaborative process modeling allowing both analyst and domain experts to collaborate in the some modulation process (Frederiks, 2006).

Perhaps the most well known reference in this field is IBM INNOV8, which as IBM states made a major impact on business games. INNOV8 is a tri-dimensional seriousgame that virtually simulates a business environment where the player assumes the role of a project manager with the goal of through his analysis over the several companies' internal business processes construct a more efficient company. The game focused on the introduction of BPM therefore the players do not need any prior knowledge of BPM methodologies to use INNOV8.

This game has been used for a few years in corporate and university environments allowing users to uncover process challenges and simulate real world bussing solution with clear benefits comparing to more traditional tools (Lapp, 2007).

3 Time-Driven Activity-Based Costing (TDABC)

Time-Driver Activity Based Costing is a costing methodology (Kaplan & Anderson, 2003). that aims at solving some of the problems of the traditional Activity-Based Costing, in which costs were calculated based on the assignment of resource expenses to activities that were verified through interviews and surveys. Although this is effective for small processes, it becomes inefficient and not at all accurate for more complex processes. To address this issue TDABC only requires two parameters: The unit cost of a resource (e.g. Logistics Department), that can be calculated from the total expenses related to that particular resource divided by its capacity (normally expressed in time), and the time needed to execute a particular task. With these two values we can know the cost of a particular activity simply by multiplying the time taken by the unit cost. Furthermore we define processes as a composition of activities.

Another ABC problem that TDABC aims to solve is that each time there is a small change or variation to an activity, there is the need for creating a new activity, for example standard packaging and hazardous packaging would result in two completely different activities, with TDABC we could have the same activity "packaging", but with different times whether it's a standard or a hazardous one. This greatly reduces the complexity of the processes since we can express all conditions in a single activity. TDABC also has the advantage of being able to give information about its own accuracy as well as help identifying waste, by comparing the calculated capacity of a given resource vs. the actual used capacity in a given period.

The TDABC methodology is used to analyze the costs resultant of a set of activities or of a process, but it can be used not only to see what was the cost of those processes but also to do what-if analysis, based on historical data. For instance it is possible with TDABC to see what would be the result of removing/adding an activity, or reducing/increasing the time it takes. On top of this TDABC also gives the ability to generate custom reports, to evaluate costs based on clients, departments, areas, etc. and to drill-down in order to identify the causes of waste.

Allying a game environment with TDABC is a breakthrough approach that has three major advantages:

- 1. What-if Scenarios Using a simulation environment it becomes easier to do an impact analysis based on historical data, since it becomes possible to do adjustments to the various components of the model, using visually appealing elements, and get real time feedback, instead of being mandatory to change the whole model and then reapplying it to the data. Also, it becomes easier to simulate local changes, without having to contemplate the whole model.
- 2. Definition of Processes\Activities\Time Equations A problem very common in the conventional implementation of an application that uses the TDABC methodology, is that the definition of all elements must be done using spreadsheets, tables or even by manually defining the time-equations. Obviously a process analyst is more used with Business Process Modeling Notations than with these methods. With a simulation it would be possible to ally a familiar concept like BPMN, and at the same time give some assistance to the user by asking for needed values or warning for inconsistency. On another perspective a simulation would help teaching the steps to define a model.
- 3.<u>Reaching all stakeholders</u> Another aspect of a simulation would be to help presenting results or to show the importance of certain data to both management and operational levels. For management it is always easier to make decisions based not only on ideas, but on real numbers as well. For operational level some decisions or changes are normally better accepted if their future impacts and advantages are shown. All these data would be possible to show using a simulation. Also regarding reaching stakeholders, this simulation/game could create a common ground on a company motivating every employee to the importance of an accurate definition of processes and activities and company's global objectives.

4 ImPROVE: A Serious Game Based on TDABC

ImPROVE is a 3D serious game developed using the Unity3D game engine and it was based on a real-world example, specifically the implementation of the Manchester triage system (Machester Triage Group, 2005) on a Portuguese Emergency unit hospital. Nowadays, triage is applied in various healthcare settings such as in mass casualty incidents, the intensive care unit, and emergency departments. Triage systems tend to rely on three different healthcare values. First, they intent to protect endangered human lives and human health. These systems therefore prioritize patients with urgent care need to treatment while less severely ill or injured patients can safely wait. However, in case several patients have to wait for life-saving interventions because one patient needs too many resources, the latter patient will not be treated first. This situation is related to the second healthcare value, efficient use of resources. Because healthcare resources are scarce, these resources will be allocated to the patients in greatest need and with the largest probability of survival. The third and final value on which triage systems rely is fairness and refers to the use of established guidelines for allocating resources to patients. With these guidelines, decisions are made on the basis of standards instead of personal preferences.

In this context the ImPROVE game provides a player with the ability to model the business process underlying the hospital emergency unit and check its impacts on healthcare values and hospital costs using the simulation features. The game main screen presents the player with a 3D representation of a set of swim lanes (Figure 1), each one visually distinguishing responsible for sub-processes of the emergency unit business process. The dimensions and number of swim lanes are setup either on a xml file or they can also be easily changed in the unity3D Editor. A simple graphical user interface (GUI) was also developed in order to manage the creation of business process primitives (Figure 1), namely activity, decision points, start and finish. The main goal of ImPROVE is to assist and enhance business process modeling and simulation activities in order to provide two main benefits, promote organizational learning and organizational change. In this sense, the set of activities that can be used are pre-defined which clearly described the range of possible activities involved or otherwise done by people in the real-world setting. Therefore, when the player is building a business process feedback is provided in order to guide the player and also to prevent designing impossible sequence of activities and or decision points. Providing real-time feedback while modeling and or simulating a business process represents an important improvement comparing to current tools. This not also motivates the player as it also gives information to the player important for decision-making. In this manner, this knowledge could then be more easily transform in explicit knowledge and used successfully in the real world when executing similar activities.







While the player is building the business process the layer responsible for checking its validity and simulation is automatically creating the time equations for calculating the cost according to the TDABC methodology. For this methodology considerations of time are very important, therefore the player has the ability to setup each activity duration and cost driver in real-time by accessing a context GUI. The relationships between business process primitives are manually created by the player by pressing a specific swim lane element and choosing the appropriate link.

Once the business process model is finished the player has the ability to test its impact by simulating it. This operation will provide the player with visual information regarding patient health, waiting time and associated costs. These values are the most important for the final score of the game, therefore it given the possibility for the player to test the final model and respective setup information three times before he/she submits the final model. In each of the intermediate simulations the player can also have access to relevant information related to the current business process being simulated.

The business process model was developed on a separate unity3D scene in order to be easily integrated in a different context. Ultimately, an adventure game can be built around this base context and modeling and simulating a business process can be just one of the many activities of the game. Therefore, the use of serious games in a company could, for example, serve as tool to create a better understanding of the prevailing organizational culture (represented on the adventure part of the game, with characters, etc), structure, and processes to access the risks, chances, and necessities of organizational change.





Figure 2 – Interaction with primitives [3]. Emergency unit business process [4]

5 Conclusions

The potential of Serious Games and Virtual-based environments to promote training has already been recognized and put into practice in several application areas, as is the case of a wide range of military sectors, aviation, healthcare, and the energy sector. This potential has also been recognized in business contexts where several simulations and games have been developed in the last decades. As described previously, several authors referred that the most important advantages of applying games in a business context, are the immediate feedback, active participation of students, learning from the experience, observation of the key factors in an on-the-job situation, preparation for the uncertainty of business, and the high motivation to learn created by the competitive environment. These particular advantages could greatly benefit the activity concerning business process modeling because although there are a lot of tools that support this activity there is still a considerable gap regarding the ability to do bug detection, reporting and analysis. Due to these constraints, business process modeling methods are still nowadays based on research design and best practices. Serious games could help bridge this gap. In this paper an example of how this could be accomplished was presented. The ImPROVE serious game allows the player to model a business process and visually receive real-time feedback of the impacts of implementing that particular business process in a business context, which in the presented example was the emergency unity of a Portuguese hospital. The ability to receive real-time feedback as well as visually witness the impacts of making certain decision (e.g. deciding on a particular sequence of activities or using certain resources) brings the activity of modeling business processes closer to the real world therefore, to the on-job situations. This also promotes the transformation of tacit knowledge into explicit knowledge representing a shorten period between learning a competence and being able to apply it on a concrete or similar situation. Finally, although several advantages have been pointed out, there is still a lot of space for improvements, namely adding multiplayer support and collaborative tools could greatly increase knowledge transfer through socialization and promote collaboration between employees.

6 References

A.J. Faria, D. H. (2009). Developments in business gaming. *Simulation & Gaming*, 40(4), 464-487.

- Aldrich, C. (2004). Simulations and the future of learning: an innovative (and perhaps revolutionary) approach to e-learning. Pfeiffer.
- Barjis, J., & Verbraeck, A. (2010). The relevance of modeling and simulation in enterprise and organizational study. *Enterprise and Organizational Modeling and Simulation*, 63, pp. 15–26.
- Beck, J. C. (2004). *Got Game: How the gamer generation is reshaping business forever*. Boston: MA:Harvard Business School Press.
- Brown, R. A. (2009). Improving the Traversal of Large Hierarchical Process Repositories. *Proceedings* of the 20th Australasian Conference on Information Systems.
- Carroll, T. H. (1954). Where is Business Education Going? *Business Quarterly*, 145-152.
- Chang, J. L. (2003). Business Simulation Games: The Hong Kong Experience. *Simulation and Gaming*, 367-376.
- Corti, K. (2007). Gamesbased Learning a serious business application.
- Faria, A. J. (2004). "A Survey of Simulation Game Users, FormerUsers, and Never-Users. *Simulation and Gaming*, 178-207.
- Frederiks, P. J. (2006). Information modeling: the process and the required competencies of its participants. *Data Knowl. Eng.*, 4-20.
- Fu, F. L. (2009). EGameFlow: A Scale to Measure Learners Enjoyment of E-Learning Games. *Computers and Education*, 101-112.
- Gilgeous, V. a. (1996). A Study of Business and Management Games. *Management Development Review*, 32-39.

- Gredler, M. (1994). *Designing and Evaluating Games and Simulations*. Gulf Professional Publishing.
- Hill, J. C. (1990). Magic Quadrant: Business Process Management Suites.
- Jenkins, K. S. (2004). Harnessing the Power of Games in Education.
- Jensen, K. O. (2003). "Business Games as Strategic Team-Learning Environments in Telecommunications. *BT Technology Journal*, 133-144.
- Kaplan, R., & Anderson, S. (2003). *Time-Driven Activity Based Costing*.
- Lapp, D. (2007). Innov8 Code: A BPM Simulator.
- Machester Triage Group. (2005). *Emergency Triage,* 2nd Edition. BMJ Books.
- Martin, a. (2000). The Design and Evolution of a Simulation/ Game for Teaching Information Systems Development. *Simulation & Gaming*, 31(4), 445–463.
- McGonigal, J. (2012). Reality Is Broken: Why Games Make Us Better and How They Can Change. Random House.
- Prensky, M. (2001). *Digital game-based learning*. New York: McGraw Hill.
- Randel, J. M. (1992). The effectiveness of games for educational purposes: a review of recent research. *Simul. Gaming*, 261-276.
- Recker, J. C. (2006). How Good is BPMN Really? Insights from Theory and Practice. 14th European Conference on Information Systems.
- Rieber, L. (1996). Seriously considering play: Designing interactive learning environments. *Educational Technology Research and Development*, 43-58.
- Rosemann, M. R. (2006). A Study of the Evolution of the Representational Capabilities of Process Modeling Grammars. Advanced Information Systems Engineering - CAiSE 2006, 447-461.

- Sánchez Franco, M. M. (2009). Exploring the Impact of Individualism and Uncertainty Avoidance in Web-Based Electronic Learning: An Empirical Analysis in European Higher Education. *Computers and Education*, 588-598.
- Savill-Smith, A. M. (2004). The use of computer and video games for learning a review of the. *Learning and Skills Development Agency*.
- Smith, R. D. (2009). *Military Simulation & Serious Games: Where we came from and where we are going*. Modelbenders Press.
- Squire, K. a. (2005). From Users to Designers Building a Self-Organizing Game-Based Learning. *Technology Trends*, 32-42.
- Streit, A. T. (2005). Visualization Support for Managing Large Business Process Specifications. 3rd International Conference, Business Process Management.
- Szezurek, M. (1982). Meta-analysis of simulation games effectiveness for cognitive learning.
- Tapani, N. (2008). Applying Serious Gaming to Business Process Management with Viprosa. *IADIS International Conference Gaming 2008*, 128-132.
- Turkle, S. (1984). *Video Games and Computer Holding Power*. The New Media Reader.
- Van Eck, R. (2006). Digital game-based learning: It's not just the digital natives who are restless. *Educause Review*, 1-16.
- VanSickle, R. L. (1986). A quantitative review of research on instructional simulation gaming: A twenty-year perspective. *Theory and Research in Social Education*, 245-264.
- Viknashvaran Narayanasamy, K. W. (2006). Distinguishing games and simulation games from simulators. Computers in Entertainment, 4(2), 9.

- Walters, B. A. (1997). Simulation Games in Business Policy Courses: Is There Value for Students? *Journal of Education for Business*, 170-177.
- West, S. B. (2010). Collaborative business process modeling using 3D virtual environments. Proceedings of the 16th Americas Conference on Information Systems : Sustainable IT Collaboration around the Globe, Association for Information Systems.
- Zantow, K. K. (2005). More than Fun and Games: Reconsidering the Virtues of Strategic Management Simulations. *Academy of Management Learning & Education*, 451-458.

Unit-Delay Simulation with the EVCF Algorithm

Peter M Maurer

Dept. of Computer Science, Baylor University, Waco, TX, 76798-7356, USA

Abstract - The EVCF (Event-Driven Condition Free) algorithm is a differential simulation algorithm that is capable of simulating virtually any type of circuit. In this paper we extend the EVCF algorithm to the Unit-Delay model, which allows for the detection of both static and dynamic hazards in a circuit. The Unit-Delay EVCF algorithm is fast, more than four times faster than conventional unit-delay simulation, and has an extremely small simulation core of around 2,000 bytes. It is suitable for use in debugging commercial-grade circuits.

Keywords: Digital Simulation, Logic Simulation, Unit-Delay Timing

1 Introduction

Event driven logic simulation has been the subject of much research[1]. A variety of logic models and delay models have been studied[2]. There have been several improvements in performance based on two different approaches. Internally, a logic simulator consists of scheduling algorithms and gate simulation algorithms. One approach to improving performance is to improve the scheduling algorithms [3,4,5]. The other approach is to improve gate-simulation speed [6,7]. Differential simulators do both. They focus on improving scheduling speed by improving the efficiency of event propagation and they eliminate much if not all of the gate-simulation code. [7].

The Event-Driven-Condition-Free (EVCF) logicsimulation algorithm [8] is intriguing in two respects. First, it operates in differential mode. Instead of computing the values of gate outputs, it computes differences between successive inputs to determine whether events propagate from gate inputs to gate outputs. In such a simulation, no net values are required, except at the primary inputs and primary outputs. This enables some peculiar optimizations. Because NOT and BUFFER gates always propagate events, they can be eliminated entirely from the simulation. AND, OR, NAND, and NOR gates all appear identical during the simulation, permitting gates to be combined in unusual ways.

The second intriguing aspect of the EVCF algorithm is that it relies entirely on metamorphic code for implementing internal states. Except for processing the primary inputs and testing for termination, the code contains no conditional branches and no loops. The code is peculiar looking, and extremely compact. The EVCF algorithm is capable of simulating virtually any circuit using a simulation core of around 1,500 bytes. In addition, the EVCF algorithm is many times faster than conventional simulation.

The major drawback of the EVCF algorithm is that it is a zero-delay algorithm. Run-time queues are complicated, because the circuit must be levelized during simulation to avoid simulating a gate before its inputs have been computed. Hazards can neither be detected nor reported because each gate is simulated only once. This means that potential instabilities in the circuit can go undetected. Cyclic circuits cannot be handled properly for the same reason.

In the following we show that the EVCF algorithm can be adapted to handle the unit-delay model in which each gate has an implicit delay of one. This is sufficient to detect hazards and oscillations in a circuit, and it provides a basis upon which an even more accurate timing simulation can be built.

2 Event Handling

The unit-delay EVCF algorithm has two types of structures, event structures which represent one fanout branch of a net, and gate structures which represent gates. A net is represented as a linked-list of events, one for each fanout branch. Primary outputs have an additional fanout branch that is used to compute the value of the output. Figure 1 shows the details of the event structure. The *Next* and *Prev* fields are used to create a doubly-linked list of events. (Double links are used to facilitate event cancellation.) The *Routine* field points to the processing routine for the event. The *Gate* field points to the gate-structure of the fanout branch or the value of the net for primary output branches.

Figure 2 shows the representation of a net. This linkedlist of structures is added and removed from the event queue as a unit. When an event occurs, the entire list is added to the event queue, and when an event is canceled, the entire list is removed from the queue.

The Routine field of the event structure points to one of four routines. These routines are accessed using a computed go-to. The four routines are shown in Figure 3. The code in Figure 3 uses the GCC extensions for handling pointers to labels and computed go-tos. EVUP and EVDN are used for AND, OR, NAND and NOR gates, while NOT is used for NOT, BUFFER, XOR, and XNOR gates. The MONITOR routine is used to compute the values of primary outputs.

In the routines of Figure 3, the variable *shp* points to the current event structure. The EVUP and EVDN routines alternate with one another by replacing the *Routine* pointer with each other's addresses. After doing this, these routines load the Gate pointer and then branch to a routine contained

in the gate structure for further processing. The NOT routine skips the additional processing and goes immediately to the scheduling routine.

| Next |
|---------|
| Prev |
| Routine |
| Gate |

Figure 1. An Event Structure.



Figure 2. A Net Structure.

The MONITOR routine computes the new value of the gate by inverting the existing value. The MONITOR event structure will never be scheduled unless the value of the output changes. No additional events are scheduled by a MONITOR event, so the routine simply goes to the next event structure and begins executing it.

| EVUP: | NOT: |
|---|-------------------------------|
| <i>shp->Routine</i> =&& <i>EVDN;</i> | shp2 = shp->Gate; |
| shp2 = shp->Gate; | goto * shp2->Schedule; |
| goto * shp2->Up; | |
| EVDN: | MONITOR: |
| <i>shp->Routine</i> =&& <i>EVUP;</i> | shp-> $Gate =$ |
| shp2 = shp->Gate; | (struct Gshadow |
| goto * shp2->Down; | *)((long) |
| | (<i>shp->Gate</i>) ^ 1); |
| | shp = shp -> Next; |
| | goto * shp->Routine; |

Figure 3. Event-Handling Routines.

3 Gate Handling

The gate structure is shown in Figure 4. This structure is used by EVUP, EVDN, and NOT events. This structure will be scheduled as if it were an event, so the three important scheduling fields, *Next*, *Prev*, and *Routine*, are present in this structure. The *Begin* and *End* fields point to the linked lists of events that represent the output of the gate. This list of events will be scheduled if the output of the gate changes. The *Up* and *Down* fields contain pointers to processing routines. These routines are used by EVUP and EVDN events. The processing routines are used to recompute the state of the gate and to determine whether the output of the gate changes. If the output of the gate changes, these routines will schedule the gate using the *Next* and *Prev* fields of the gate structure.

The *Schedule* field points to the current scheduler of the gate. This alternates between the QUEUE and the DEQUEUE routines. The QUEUE routine adds the gate structure to the

queue, while the DEQUEUE routine removes the gate structure from the queue, thus cancelling event propagation.

| Next |
|----------|
| Prev |
| Routine |
| Begin |
| End |
| Up |
| Down |
| Schedule |

Figure 4. The Gate Structure.

The state of an AND, OR, NAND or NOR gate is maintained by using several different UP and DOWN routines. Figure 5 gives a sample of these routines. As a practical matter, only a few routines are needed, because the number of inputs of an AND, OR, NAND, or NOR gate is limited by the underlying technology.

Note that each of the UP and DOWN routines schedules new UP and DOWN routines. (Except DN0 which should never be executed) Except for UP0 and DN1, the routines continue with the next event. Instead of going to the next event, the UP0 and DN1 routines execute the routine pointed to by the *Schedule* field of the gate. The pairs of routines, UPx/DNx, represent the number of inputs of an AND, OR, NAND, or NOR gate that have the dominant value (x). The output of the gate changes when the dominant count goes from 0 to 1 (UP0 executes), or when the count goes from 1 to 0 (DN1 executes). In all other cases, the output of the gate does not change.

| UP0: | UP1: |
|---|--|
| shp2->Up = &&UP1 | shp2->Up = &&UP2 |
| shp2-> $Down = &&DN1$ | $shp2 \rightarrow Down = \&\&DN2$ |
| goto * shp2->Schedule; | shp = shp -> Next; |
| | goto * shp->Routine; |
| DN0: | DN1: |
| shp = shp -> Next; | shp2->Up = &&UP0 |
| goto * shp->Routine; | shp2-> $Down = &&DN0$ |
| | goto * shp2->Schedule; |
| | |
| UP2: | UP3: |
| UP2: shp2->Up = &&UP3 | <i>UP3: shp2->Up</i> = && <i>UP4;</i> |
| UP2: shp2->Up = &&UP3 shp2->Down = &&DN3 | UP3: shp2->Up = &&UP4 shp2->Down = &&DN4 |
| UP2: shp2->Up = &&UP3 shp2->Down = &&DN3 shp = shp->Next; | UP3: shp2->Up = &&UP4 shp2->Down = &&DN4 shp = shp->Next; |
| UP2: shp2->Up = &&UP3 shp2->Down = &&DN3 shp = shp->Next; goto * shp->Routine; | UP3: shp2->Up = &&UP4 shp2->Down = &&DN4 shp = shp->Next; goto * shp->Routine; |
| UP2: shp2->Up = &&UP3 shp2->Down = &&DN3 shp = shp->Next; goto * shp->Routine; DN2: | UP3: shp2->Up = &&UP4 shp2->Down = &&DN4 shp = shp->Next; goto * shp->Routine; DN3: |
| UP2: shp2->Up = &&UP3 shp2->Down = &&DN3 shp = shp->Next; goto * shp->Routine; DN2: shp2->Up = &&UP1 | UP3: shp2->Up = &&UP4 shp2->Down = &&DN4 shp = shp->Next; goto * shp->Routine; DN3: shp2->Up = &&UP2 |
| UP2: shp2->Up = &&UP3 shp2->Down = &&DN3 shp = shp->Next; goto * shp->Routine; DN2: shp2->Up = &&UP1 shp2->Down = &&DN1 | UP3: shp2->Up = &&UP4 shp2->Down = &&DN4 shp = shp->Next; goto * shp->Routine; DN3: shp2->Up = &&UP2 shp2->Down = &&DN2 |
| UP2: $shp2 \rightarrow Up = \&\&UP3$ $shp2 \rightarrow Down = \&\&DN3$ $shp = shp \rightarrow Next;$ $goto * shp \rightarrow Routine;$ DN2: $shp2 \rightarrow Up = \&\&UP1$ $shp2 \rightarrow Down = \&\&DN1$ $shp = shp \rightarrow Next;$ | UP3: shp2 > Up = &&UP4 shp2 > Down = &&DN4 shp = shp > Next; goto * shp > Routine; DN3: shp2 -> Up = &&UP2 shp2 -> Down = &&DN2 shp = shp -> Next; |

Figure 5. Up and Down routines.

| QUEUE: | DEQUEUE: |
|---|----------------------|
| <i>shp2->Next</i> = & <i>Trailer</i> ; | shp2->Next->Prev = |
| shp2->Prev = Trailer.Prev; | shp2->Prev; |
| <i>Trailer</i> . <i>Prev</i> -> <i>Next</i> = | shp2->Prev->Next = |
| (NSHADOW *)shp2; | shp2->Next; |
| Trailer.Prev = | shp2->Schedule= |
| (NSHADOW *)shp2; | &&QUEUE |
| shp2- | shp = shp -> Next; |
| >Schedule=&&DEQUEUE | goto * shp->Routine; |
| shp = shp -> Next; | |
| goto * shp->Routine; | |

Figure 6. The QUEUE and DEQUEUE routines.

Figure 6 shows the QUEUE and DEQUEUE routines. The QUEUE routine takes the current gate structure, which is pointed to by the variable *shp2*, and adds it to the end of the queue. The variable *Trailer* points to the element that is currently at the end of the queue. This structure is a marker that must remain at the end of the queue, so gate structures must be inserted ahead of this element. The queue is initialized with a dummy element at the head of the queue, so the Trailer structure will always point back to a valid queue element. The first four statements of the QUEUE routine perform a standard double-link operation. Next, the QUEUE routine replaces its own address with the address of the DEQUEUE routine, and finally advances to the next event.

The DEQUEUE routine is simpler than the QUEUE routine, because only the fields of the events remaining in the queue must be changed. The routine delinks the gate structure, replaces its own address with the address of the QUEUE routine and advances to the next event.

4 Event Scheduling

Event scheduling is done by the gate data structures once they have been placed in the event queue. These structures are not placed in the queue unless the output of the gate changes. Trailer events are used to mark time intervals and to terminate the queue. Figure 7 shows the general structure of the queue at different times. Queue state A occurs right after a new input vector has been read, and new events have been scheduled. This state will recurr many times throughout the course of the simulation. The *Trailer* event marks the end of the queue. State B occurs after some events have been processed. These events will generally cause Gate structures to be scheduled. Since the gate processing belongs to the current time interval, the gates are inserted ahead of the *Trailer* event.

Queue state C shows the queue after all events have been processed and after some gates have been processed. Processing a gate will cause evens to be scheduled, but because these events belong to the next time interval, not to the current time interval, they are inserted after the *Trailer* event and ahead of the *Tail* event. Figure 8 shows the processing routine for the gate structures. All gate structures use the same processing routine.







Figure 8. The Event Scheduling Routine.

The FORWARD routine of Figure 8 first adds events unconditionally to the queue. If the gate structure is in the queue, then the output of the gate has changed, and it is necessary to schedule events. The list of events is pointed to by the *Begin* and *End* fields of the gate structure. The routine then resets the corresponding gate structure so that it will queue the next event.

| TRAILER: | TAIL: |
|--|---------|
| <pre>putvec(ifile);</pre> | return; |
| <i>if(shp->Next != &Tail) {</i> | |
| work = shp -> Next; | |
| shp-> $Next = &Tail$ | |
| <i>shp->Prev = Tail.Prev;</i> | |
| Tail.Prev->Next = shp; | |
| Tail.Prev = shp; | |
| shp = work; | |
| } else { | |
| shp = shp -> Next; | |
| } | |
| $goto * shp \rightarrow Rtn;$ | |

Figure 9. The Trailer and Tail Event Handlers

The trailer event is used to mark the boundaries between two consecutive units of simulated time. The first action of the *Trailer* event handler is to print an output vector containing the value of each primary output. Next, the event handler tests its *Next* pointer. If the next event is the *Tail* event, indicating that there are no events queued for the next time period, the trailer event processor simply continues with the next event, terminating the simulation of the current input vector. However, if there are events queued for the next time period, the trailer event handler inserts itself at the end of these events. Figure 9 shows the code for the *Trailer* and *Tail* event handlers.

Apart from the code used to read input vectors and print output vectors, all code for the simulator appears in this paper. Needless to say, the simulation core is tiny compared to that of conventional simulators: around 2,000 bytes.

5 Experimental Data

Figure 10 shows the experimental results for the unit delay EVCF algorithm. These results use the ISCAS-85 benchmarks [9] which have become a de-facto standard for measuring the performance of simulation algorithms The simulations were performed on a dedicated 3.06 Ghz Xeon processor with 2GB of 233 Mhz memory. The results are in seconds of execution time for 5,000 input vectors.

The unit-delay EVCF algorithm outperforms conventional unit-delay simulation by a factor of between 4 and 5 for most circuits.

| Circuit | Conventional | EVCF | Improvement |
|--------------|--------------|--------|-------------|
| <i>C432</i> | .528 | .138 | 3.826 |
| C499 | 1.672 | .344 | 4.860 |
| C880 | 1.752 | .464 | 3.776 |
| C1355 | 3.988 | .760 | 5.247 |
| <i>C1908</i> | 4.092 | .944 | 4.335 |
| C2670 | 13.226 | 3.180 | 4.159 |
| <i>C3540</i> | 5.448 | 1.172 | 4.683 |
| C5315 | 24.302 | 5.076 | 4.788 |
| C6288 | 30.928 | 13.692 | 2.759 |
| C7552 | 25.440 | 5.108 | 4.980 |

Figure 10. Experimental Results.

6 Conclusions

The unit-delay EVCF algorithm is an unconventional algorithm that is extremely efficient both in terms of execution time and computer memory. By operating in differential mode it is able to realize an improvement of 4-5x over conventional unit-delay simulation. Because it uses

metamorphic coding which is virtually devoid of loops and conditions, the simulation core is tiny, around 2,000 bytes. Despite the size of the simulation core, the unit-delay EVCF algorithm is able to simulate virtually any circuit, including the standard ISCAS-85 benchmark circuits. Unlike the zerodelay EVCF algorithm, the unit-delay EVCF algorithm is able to detect both static and dynamic hazards, making it a useful tool for debugging commercial-grade designs.

7 References

- E. G. Ulrich, "Event Manipulation for Discrete Simulations Requiring Large Numbers of Events," JACM, V.21, N.9, Sep. 1978, pp. 777-85.
- [2] Szygenda, S., D. Rouse, E. Thompson, "A Model and Implementation of a Universal Time-Delay Simulator for Large Digital Nets," Spring Joint Computer Conference, 1970, pp. 491-496.
- [3] D. M. Lewis, "A Hierarchical Compiled Code Event-Driven Logic Simulator," IEEE Transactions on Computer Aided Design, Vol 10, No. 6, pp.726-737, June 1991.
- [4] D. M. Lewis, "Hierarchical Compiled Event-Driven Logic Simulation," Proceedings of ICCAD-89, pp.498-501.
- [5] Z. Wang and P. M. Maurer, "LECSIM: A Levelized Event Driven Compiled Logic Simulator," Proceedings of the 27th Design Automation Conference, 1990, pp. 491-496.
- [6] M. Heydemann, D. Dure, "The Logic Automation Approach to Accurate Gate and Functional Level Simulation," Proceedings of ICCAD-88, pp. 250-253.
- [7] P. M. Maurer, "The Inversion Algorithm for Digital Simulation," Proceedings of ICCAD-94, pp. 259-61.
- [8] P. M. Maurer, "Event Driven Simulation Without Loops or Conditionals," ICCAD 2000, Nov. 2000, pp.
- [9] F. Brglez, Pownall, Hum, "Accelerated ATPG and Fault Grading via Testability Analysis," ISCAS-85, pp. 695-698.

Review of A Programming language Analysis in Simulation modleing in health care

Maher Amer Department of Industrial & Engineering Technology Southeast Missouri State University Cape Girardeau, MO

Abstract—this publication will review a paper discussing programming language Analysis in Simulation modeling in health care. The paper reviewed showed that simulation modeling can become more attractive if done correctly. The use of alternate programming languages or techniques such as variance reduction enhanced the simulation model greatly. With further research and improvement, flexible and accurate models that consume less and less time will be possible.

Keywords- Simulation Languages, language analysis, simulation modeling, health care

I. Introduction

Simulation modeling and programming has become a vital tool in solving problems and in aiding in decision making. Also realistic graphic simulation allows students to observe scientific, industrial, and brings reality into the classroom [1][2]. Due to the increased dependence on simulation modeling it is very important to validate these models. Validating a model is defined by "substantiation that a computerized model within its domain of applicability possesses a satisfactory range of accuracy consistent with the intended application of the model" [3]. In order to do so, the model is evaluated until confidence is established in the models validity [4][5]. Also the programming of the simulation needs to be accurate and done in a timely fashion [6]. This, today, is more essential than ever in the world of simulation modeling as the need for a more flexible modeling is rising to adapt to the new complex problems in science which led to new techniques such as Bayesian simulation modeling using Win [7] and stochastic, Monte Carlo simulation modeling [8]. The research paper that is reviewed here addresses the topic and shows a study made to assess the Relationship between Computational Speed and Precision.

II. Review of A Programming language Analysis in Simulation modleing in health care

In this section we will review the paper of McEwan et al.[9] The research paper at hand is concerned with the use of simulation techniques in disease modeling. Many models that were developed in this field are quite adequate and provide Suhair H. Amer, Ph.D. Department of Computer Science Southeast Missouri State University Cape Girardeau, MO

great results but some issues still offer challenges that need to be addressed. Mainly the biggest concerns with these models are time consumption and precision. For example, Markov models are adequate in simulation but require great computational effort. On the other hand Monte Carlo simulation models have got issues in their accuracy. To improve on the precision the paper recommends using variance reduction and evaluation of the number of simulation replications required. In doing so, the research paper will study two common programming languages to point out prespecified differences in simulation output.

For the purpose of this study, a model was made to evaluate the consequences of the introduction of new medications or therapies on a sample population of type 2 diabetes. The model uses fixed time increment stochastic simulation. For its input, the model takes in certain information such as: age, ethnicity, sex, whether the patient smokes or not, how long has the patient been diagnosed with diabetes. It also takes into consideration modifiable risk factors such as: total cholesterol, weight and blood pressure. The model will output risk factors of micro vascular (retinopathy, neuropathy or nephropathy) and macro vascular (congestive heart failure) and costs associated with those. In this study the model was used in two ways: mean input values that were applied through simulation replications and sampled input values. In the sampled input values the risk factors were taken from means and maximum and minimum values. The method in this research is basically to replace Visual basic applications program that was the original programming language for the model with C++ programming language and then test the differences. A scenario was used in which a treatment is introduced and then the effect on it on a sample of 1000 patients is studied along with running QALYs analysis to test effectiveness and cost effectiveness. A span of 40 years was used for the model. The quality-adjusted life year (QALY)s cumulative is going to be the estimator X. Studying half the width of the confidence interval in the QALYs cumulative will be compared to a certain preset value of the number of simulation replications required. The time required to achieve the calculation is measured by the use of a call to the system clock. The time will be measured from start to tend of computations. Also variance reduction was used to test its effect on calculation times.
The results in this research were as follows, a sample population of cohort 1000 patients was used in the model which requires 21 and 5653 replications for the mean and sampled input values respectively. For the Visual Basic applications the time to process these replications was 3.0 and 810.2 minutes. For C++ the time was 0.04 and 8.9minutes. Also analyzing the 25 input parameters using standard oneway sensitivity was done in 2 minutes using C++ while it took 155 minutes to do so in Visual Basic Applications. Introducing antithetic varieties minimized the number of simulation replications required by 53%. In general if applying 50 scenarios it would take C++ to finish them in 3.8 hours while it would take Visual Basic Applications 14.5 days to perform the same calculations.

The conclusion reached in this research paper is that Health models need not be time consuming and that calculation need not be a burden. This was shown in the huge difference in time needed by two different programs addressing the same problem.

The comparison made in this research paper is basically between two programming languages used to process data in the same model. The results have shown that the same model can perform better and consume much less time if its core programming is changed. This shows the superiority of one programming language over the other. The question is, does it really validate the use of simulations in medical studies. In the method used in this paper a comparison was made between 2 programming languages in a single model. It might have been good if also a comparison between 2 different models was made. It might be the case that there exists better models in whole that have both a good programming language and provide accurate data. It would broaden the view of the issue if a listing of different models used and a comparison of the times consumed by each is shown in the research. Also the research concludes that medial simulation models can be much faster and less burdensome. That can't be fully achieved by speeding up models. A cost to benefit analysis must also accompany the research after faster models are made to show the justification of these models. Also since these simulation programs will predict risk factors, statistics will need to show over time that the predictions were validated.

Factors such as age, sex, ethnicity and smoking status were used in the model since it's the relevant factors to diabetes 2 patients. What if the illness was different and the number of factors was higher or smaller? Would that affect the simulation model? Also the sample population could vary and might have an effect on the results. It would be of great benefit if other models with different factors and sample populations are also studied.

III. Conclusions

In simulation modeling it is crucial that the validity of the model is verified. It is also very important to critically analyze the programming language used and study the time spent in calculating results and the accuracy of those results. A continuous effort must be made to continually improve all these factors. This will give simulation modeling a leading edge in comparison to other methods used in scientific research or decision making. The paper reviewed showed that simulation modeling can become much more attractive if done correctly. The use of alternate programming languages or techniques such as variance reduction enhanced the simulation model greatly. With further research and improvement, flexible and accurate models that consume less and less time will be possible.

References

- [1] Wang, S., & Reeves, T. C. (2007). The effects of a webbased learning environment on students' motivation in a high school earth science course. *Journal of Education Technology Research and Development*, 55 (2), 169-192.
- [2] Zhang, J. (2007). A cultural look at information and communication technologies in eastern education. *Journal of Education Technology Research and Development*, 55 (3), 301-314.
- [3] Schlesinger, S. 1979 el al. Terminology for model credibility. *Simulation* 32 (3): 103-104.
- [4] Sargent, R. G. 1982. Verification and validation of simulation Models. Chapter IX in *Progress in Modeling* and *Simulation*, ed. F. E. Cellier, 159-169. London: Academic Press.
- [5] Sargent, R. G. 1984b. A tutorial on verification and validation of simulation models. In *Proceedings of the 1984 Win-ter Simulation Conference*, ed. S. Sheppard, U. W. Pooch, and C. D. Pegden, 114-121.Piscataway, NJ: IEEE.
- [6] Karnon J. Alternative decision modeling techniques for the evaluation of health care technologies: Markov processes versus discrete event simulation. Health Econ 2003 Oct; 12 (10): 837-48
- [7] Armero C, Garci'a-Donato G, Lo' pez-Qui'lez A. Bayesian methods in cost-effectiveness studies: objectivity, computation and other relevant aspects. Health Econ 2010; 19 (6): 629-43
- [8] Caro JJ. Pharmacoeconomic analyses using discrete event simulation. Pharmacoeconomics 2005 Apr; 23 (4): 323-32
- [9] McEwan, Phil; Bergenheim, Klas; Yong Yuan; Tetlow, Anthony P.; Gordon, Jason P. Assessing the Relationship between Computational Speed and Precision. PharmacoEconomics, 2010, Vol. 28 Issue 8, p665-674,

SESSION MODELING

Chair(s)

TBA

Modeling Discrete Distributed Heterogeneous Systems

Irene Pestov¹, Hiroki Sayama², and Chun Wong²

¹Defence R&D Canada, Centre for Operational Research and Analysis, Ottawa, Ontario, Canada ²Collective Dynamics of Complex Systems, Binghamton University, Binghamton, New York, USA

Abstract—In this paper we consider discrete spatiallydistributed systems that are designed to perform two main functions: detection of a significant event and execution of a response action. Under our approach, the system is modeled as a dynamical network with embedded heterogeneous agents. Each agent is encoded as a multidimensional string to represent its role in a network, the environmental domain in which it operates, and other case-specific attributes. Agents are linked in the standby network that represents the standby posture of the system. In response to a significant event, the system dynamically creates an operational network among agents of the standby network. Agents are chosen according to their specifications, as encoded in corresponding strings. The resulting dynamical network model provides visualization and simulation capabilities. The visualization capabilities are invoked to display the architectural make-up of the system and to explore its multidimensional agents. The simulation capabilities are used to simulate the development of the operational network. The utility of the modeling framework is demonstrated on an example of the Canadian Arctic Search and Rescue system.

Keywords: dynamical networks, heterogeneous agents, generative network automata

1. Introduction

Networks provide a natural way to represent discrete distributed systems composed of many interconnected components [1]. In general, each network node denotes a system component and the links represent some relationships between those components. Depending on the nature of dynamics, all network models can be subdivided into three broad categories:

- 1) Networks with rich functional/state dynamics, but fixed structure (e.g. Random Boolean Networks [2]).
- Networks with varying topology that can modify their structure according to specified rules of attachment, not coupled to functional/state dynamics (e.g. scalefree networks [1].
- Networks with coupled structural and functional dynamics. In these network models, structural changes cause changes in network states, and vice versa. Most real networks fall into this category [3].

The latter network models are of particular significance to this study. In [3], these network models are referred to as adaptive networks. Another commonly used term is 'complex dynamical networks' [4] or, simply, 'complex networks' [5]. A distinctive feature of these network models is the existence of a dynamical feedback loop between network topologies and functional states of network nodes.

In this paper, we are concerned with systems that are designed to perform two main functions: detection of a significant event and execution of an appropriate response action. Agents (or active players) in such systems are heterogeneous entities. They differ by their roles (specialization), by environmental domains (realms) of their operation, and also by other case-specific attributes. Pestov and Pilat argue in [6] that such systems are better represented by dynamical networks with embedded heterogeneous agents. As shown in [6], the agent heterogeneity gives rise to an operational network – a step-by-step representation of a response action. The operational network develops dynamically, as new links/shortcuts appear and disappear between agents that become actively involved in the prosecution of the response.

Let us consider the Canadian Arctic Search and Rescue (SAR) system [7]. In the SAR system, agents represent SAR assets: Canadian Coast Guard (CCG) officers, teams of SAR techs, Joint Rescue Coordination Centres (JRCCs), aircraft, ships, and various information and communication systems. These are highly specialized entities, trained or designed to provide specific services, often in a specified environmental domain. Unlike agents in social networks, heterogeneous agents of the SAR system cannot easily re-train and replace other agents. The agent specialization results in distinctive patterns of system dynamics.

When a distress alert is detected and transmitted to a responsible JRCC, the system mounts a SAR response. The nature and size of the incident determine the choice of SAR assets being called upon. One of the CCG officers on duty is appointed as the Search Master. From now on, the Search Master is responsible for the SAR response in question until closure of the case. The Search Master takes necessary steps to prosecute the response. These may include:

- verification of information about aircraft/vessel in distress;
- issuing notifications and updates;
- selecting and tasking agents to perform a certain task;
- coordination of rescue effort.

As a result, a new network develops between agents actively involved in the response: new nodes join the network (e.g. a vessel is charted to conduct rescue) and new links develop (e.g. between the Search Master and a search aircraft). In [6], [7], this new network is termed the Operational network. According to [6], [7], the Operational network is the most important network in the system, as its measures serve as performance characteristics for the entire system.

In this paper, the SARnet dynamical network model of the Canadian Arctic SAR system [7] is used to demonstrate visualization and simulation capabilities of the modeling framework proposed by Pestov and Pilat in [6].

2. Modeling Framework

This section describes a modeling framework for representation and analysis of discrete distributed heterogeneous systems. We have developed this framework with two purposes in mind: (1) to provide visualization capabilities that can be used to represent system architecture and to explore its heterogeneous agents; and (2) to support the development of a simulation tool for automated generation of operational networks in 'what-if' scenarios.

2.1 Agent Representation

Here, we follow Pestov and Pilat [6] in introducing five classes of agents, depending on their role (specialization) in the network, and six environmental domains (realms) of their operation. Five agent classes are: sensor, router, actor, database, and controller. Six realms are: maritime, land, air, space, cyber, and cognitive. Sensor agents represent the system surveillance capabilities. They sense, detect, and gathered information. The role of routers consists of distributing and directing the flow of information within the system. Actors accomplish the bulk of tasks, associated with a response action. Databases store information and make it available to other agents. Controllers coordinate a response and task other agents.

The agent classes differ with respect to how they access, process and exchange information, and also with respect to what degree of knowledge about the system they have. For instance, a sensor agent may only have information about its next neighbor – a router, whereas a controller may have knowledge about a significant segment of the network and agent specs within this segment.

For a sake of clarity, we assume that an agent cannot hold more than one role simultaneously. However, it could change roles during various phases of a response. For example, an actor could assume the controller functions, if required. Some agents cannot change the domain of their operation, whereas others can operate across several environmental domains. The latter is generally an exception.

Let us return to the SAR example. The Search and Rescue Satellite Tracking (SARSAT) system, which provides distress alert and location information on SAR incidents to SAR authorities, is a sensor agent in the SARnet model. SARSAT passes information through a Local Users Terminal (LUT), which is a router agent, to the Canadian Mission Control Centre (MCC), which is a controller agent. The realm of SARSAT is space. Some types of sensors can operate in air, maritime, or land realms. Routers typically operate in the cyber realm.

JRCC Trenton located in Trenton ON, which is a JRCC responsible for Arctic SAR, and the Search Master, who is appointed by JRCC to oversee SAR response, are examples of controller agents. The realm of the SARnet controller agents is cognitive. (Note that some systems may include non-human controllers operating in the cyber realm.)

In SARnet, actors are the teams of SAR techs, Royal Canadian Air Force (RCAF) SAR aircraft, CCG vessels, ground search teams, and volunteers. In general, most of these actors operate in one of three designated environmental domains: air, maritime, or land.

The Information System on Marine Navigation (INNAV) and SAR Mission Management System (SMMS) are examples of database agents. Both systems are dynamic information sources that are able to acquire and process new information and exchange that information through knowledge. The realm of database agents is cyber. The agent classification is summarized in Table 1.

In addition to agent classes and environmental domains (realms), we introduce operational domains according to traditional subdivision of system components. Thus, the SAR system is traditionally partitioned into four SAR operational domains: Maritime, Air, Ground, and Joint. (In this context, Joint refers to operations across more than one environmental domains.) The CCG oversee the Maritime SAR operations. Whereas the RCAF and Canadian Army oversee Air and Ground SAR operations, respectively. Joint SAR operations are jointly overseen by the RCAF and CCG. Note that each operational domain may include agents operating in different realms. For example, the database agent INNAV belongs to Maritime SAR, as it provides marine traffic information, but operates in the cyber domain.

Depending on the specifics of the system under consideration, one may also include such agent attributes as mobility and availability factors, knowledge areas and resources that are required for the prosecution of a response action, home organizations, standby times, geographic location of an agent, and more. Then in a machine-readable form, an agent can be represented by a string of dimension N of the following form:

$$\sigma = [\sigma_1, \sigma_2, \dots, \sigma_N]. \tag{1}$$

In Eq. 1, σ_i is either a binary variable that assumes values of [0, 1], or categorical variable that assumes values [0, 1, *] (* stands for N/A), or continuous variable that assumes any real value.

Dimension N equals to the sum of the following addends: 5 agent classes, 6 realms, 4 operational domains, 2 mobility factors (mobile or portable), 1 availability factor, # knowledge ares, # resources, # home organizations, 2 standby times

Agent Class Function Realm Sensor senses, detects, and passes gathered info space, air, maritime, land Router distributes the flow of info and enables com links cyber, cognitive Actor executes a response action as tasked air, maritime, land Database stores and provides access to info cyber Controller coordinates a response and tasks agents cognitive, cyber

Table 1: Summary of agent classes.

(working and quiet hours), and latitude and longitude of the location of a home unit (see Table 2).

Table 2: Multidimensional representation of agents.

| σ_i | Designation | Variable Type |
|------------|------------------------|---------------|
| 1–5 | agent class | binary |
| 6-11 | realm | binary |
| 12-15 | Op. domain | binary |
| 16 | subclass | categorical |
| 17-18 | mobility | binary |
| 19 | availability | binary |
| 20-48 | expertise area | binary |
| 49-60 | resource | binary |
| 61-86 | organization | binary |
| 87-88 | standby time | continuous |
| 89–90 | latitude and longitude | continuous |

The last 4 variables – standby times and geographic coordinates – are continuous variables, whereas others are either binary or categorical variables. Standby times are specified times within which an agent must be able to respond when called to duty. Different standby times are set for working and quiet hours. (Quiet hours cover hours outside a working period from 8:00-16:00 and weekends.)

An agent is deemed stationary, when both mobility factors are set to 0 (nether mobile nor portable). The availability factor is set to 1, when an agent is available for duty, and to 0 when otherwise (e.g. for technical reasons).

Dimensions corresponding to knowledge and resource nodes indicate whether or not agents have access to a particular knowledge or resource: $\sigma_i = 1$ means 'Yes'; $\sigma_i = 0$ denotes 'No'. Dimensions corresponding to organization nodes indicate the membership of agents in a particular organization: $\sigma_i = 1$ means 'Yes'; $\sigma_i = 0$ stands for 'No'.

Depending on the number of knowledge areas, resources and organizations included in the network model, dimension N could be relatively high. In the SARnet model N = 90. Note that dimension N could be increased to include more node attributes, if needed.

The unit values found in the first 15 positions of string σ indicate the agent class, the environmental domain (realm) in which this agent operates, and the operational domain to which the asset belongs, i.e.:

$$\exists i, j, k \begin{cases} \sigma_i = 1 & (i = 1, \dots, 5) \\ \sigma_j = 1 & (j = 1, \dots, 6) \\ \sigma_k = 1 & (k = 1, \dots, 4) \end{cases}$$
(2)

We say that agents belong to the same *heterotype*, if they are identical in the first 15 positions of string σ , i.e. $[\sigma_1, \ldots, \sigma_{15}]$. The number of heterotypes within the network can be used to measure the network heterogeneity.

If K is the number of heterotypes and X_k is a fraction of agents of heterotype k (k = 1, ..., K), then network entropy can be defined as follows:

$$S = -\frac{1}{\ln K} \sum_{k=1}^{K} X_k \ln X_k.$$
 (3)

In Eq. 3, network entropy S is normalized by its maximum value $S_{max} = \ln K$.

As follows from Eq. 3, network entropy $S \in [0, 1]$. The minimum value S = 0 corresponds to a network composed of one heterotype. The maximum value S = 1 corresponds to a network composed of agents evenly distributed between all K heterotypes (i.e. $X_k = 1/K$ for $k = 1, \ldots, K$). As network entropy approaches 1, the agent distribution between heterotypes becomes uniform.

According to our analysis of the operations of the Canadian Arctic SAR system [7], the number of heterotypes increases and distribution of SAR assets between heterotypes becomes less balanced in the course of search and rescue operations.

Heterotypes can also be used as auxiliary constructs in simulations, as we will see in Subs. 2.4.

2.2 Assortative Networks

The concept of assortative mixing comes from social network science [5] (e.g. linking of entities, according to age or income). In this paper, assortative networks are invoked to visualize multidimensional agents of heterogeneous systems.

If each agent is represented by string σ with scalar variables σ_i (i = 1, ..., N), then various networks can be developed that link pairs of agents with similar values of one of these variables (e.g. through calculating a correlation coefficient for scalar σ_i). Taking a step further, one can also perform assortative mixing according to a substring $\tau \subset \sigma$. Assortative networks, obtained in this manner can be used to visualize and then visually examine multidimensional agents of heterogeneous systems. They can also be used to store assorted information in simulations.

Figure 1 shows 6 assortative networks that display SAR assets, according to environmental domains (realms), as

represented by substring ($\tau = [\sigma_6, \ldots, \sigma_{11}]$). In the figure, link color matches the color of the realm. Controllers are shown by red, actors by blue, sensors by grey, routers by green, and databases by yellow dots. Because of space constraints, some agent ID labels are omitted and nodes with similar attributes are collated into meta-bodes. Thus the CRPG patrol meta-node represents ground patrols of the 1^{st} Canadian Ranger Patrol Group, operating in the North. The AOC meta-node represents the Air Operations Centres of 5 RCAF SAR squadrons, included in SARnet. The SAR Tech meta-node is the only actor agent that can operate in more than one realm. The *ORA software [8] was used to develop and visualize assortative networks.



Figure 1: SAR assets grouped according to realms.

Assortative networks, shown in Figure 1, can be used to quickly identify SAR assets, operating in a specific realm. SARnet also includes assortative networks that group SAR assets according to agent classes, SAR domains, access to knowledge and resources, and membership in organizations (see [7] for details).

2.3 The Standby Network

The Standby network is a network that represents the standby posture of the system and its everyday functions, such as maintenance and training. It links agents based on working relationships, e.g. 'who works/communicates with whom', 'who reports to whom' (for actor agents), 'who tasks whom' (for controller agents), and 'who has info about whom' (for databases).

Figure 2 shows the Standby network of the SAR system. In the figure, agents are colored according to agent classes: controllers are shown by red, actors by blue, sensors by grey, routers by green, and databases by yellow dots. Link color matches the color of the source node. Some agent ID labels are omitted, because of space constraints. The network is fragmented into several segments that consist of controller agents linked to SAR assets at their disposal (e.g. actors and databases). Thus, the controller agent AOC 424 (i.e. the Air Operations Centre of the 424 SAR squadron) is the hub of the network segment, which comprises SAR aircraft and teams of SAR techs. Isolates represent agents that are not normally part of the Standby network, but are created or join the operational network in the course of a response. The *ORA software [8] was used to visualize the Standby network.



Figure 2: The Standby network of the SAR system.

The Standby network, shown in Figure 2, recreates a 'bird's eye view' of the SAR system, which is easy to grasp and analyze. With the network representation, various aspects of the system can be visualized and then visually examined.

In addition to visualization capabilities, described above, the network representation offers an efficient framework for analyzing system dynamics, as we will see in the next section.

2.4 Simulation of System Response

The Operational network is dynamically created on the nodes of the Standby network, when a response is initiated by the system [6], [7]. The dynamics of the Operational net differs to that of the Standby net, as the architecture of the former evolves at a much shorter time scale (minutes or hours instead of months or even years). The agent heterogeneity is the main driver for the formation of the Operational net. Such agent attributes as the agent class, realm, and operational domain play an important role in selecting most appropriate responders. The agents' skill sets and access to resources are taken into account, too. Network measures of the Operational network serve as performance indicators for the entire system [6], [7].



Figure 3: The Search Master's Sphere of Influence (from [7]).

In the SAR context, the Operational network is a step-bystep dynamical representation of SAR response. In addition to the factors mentioned above, the location information of a vessel/aircraft in distress plays a role in what assets will be deployed.

The Search Master (which was an isolate in the Standby network, as can be seen from Figure 2) becomes the hub and the most influential entity of the Operational network. Figure 3 shows the Search Master's Sphere of Influence within the Operational network of a real SAR incident in the Arctic, considered in [7]. (A Sphere of Influence of a node is a sub-network of radius 1 that includes all nodes to whom that node has direct connections plus connections between those nodes.) According to [7], the Search Master's Sphere of Influence encompassed 67% of the entire Operational network. For comparison, the Sphere of Influence of the next most influential node of that network contained about 14% of agents. High centralization of the SAR Operational net is identified in [7] as a contributing factor to the operational efficiency. The *ORA Sphere of Influence report [8] was used to derive spheres of influence for key entities of the Operational network.

Conceptually, the process of the generation of the Operational network can be subdivided into the following steps:

- 1) Search for the most appropriate agent to perform the task at hand;
- 2) Issue an order to deploy; and
- 3) Task and deploy the identified agent.

In general, the above steps will need to be repeated, according to the prescribed sequence of tasks.

Computationally, the above steps translate into the following procedure. Let \mathcal{H} be a set of agent heterotypes, included in the model, and \mathfrak{h} be a subset of \mathcal{H} , associated with a particular task $\mathfrak{t}_i \in \mathcal{T}$, where $\mathcal{T} = {\mathfrak{t}_1 \rightarrow \mathfrak{t}_2 \ldots \rightarrow \mathfrak{t}_k}$ is a sequence of tasks of the response action. Then the generation of the Operational network corresponds to the following steps:

- Search for available instances of h ⊆ H associated with t_i;
- 2) Change the node status to active;
- 3) Generate new links in the Operational net;
- 4) Proceed to task t_{i+1}

The above steps will need to be repeated until sequence \mathcal{T} is exhausted.

In what follows, we describe a simulation algorithm that we have developed for automated generation of the Operational network. Our simulation algorithm has its theoretical basis on Generative Network Automata (GNA) [9], [10]. GNA is a dynamical network modeling framework that can represent both state-transition of nodes/edges and topological transformation of the network in a uniform manner as a sequence of graph rewriting events. The dynamics of a GNA is formulated in the form of two mechanisms: extraction and replacement [9], [10]. In each iteration, the extraction mechanism selects a subgraph out of the entire network, and then the selected subgraph is replaced by the replacement mechanism (Figure 4).



Figure 4: An example of GNA updating by extraction and replacement (from [10]).

In this paper, we adopted and revised this GNA framework so that it can represent discrete heterogeneous systems and their dynamics efficiently. The simulation code was developed in Python, and NetworkX was used for network representation and analysis [11].

Our model represents the dynamical network of heterogeneous agents – a directed network made of discrete nodes, each of which has its own attributes and internal variables that are direct translations of the agent representation described in Subs. 2.1. The dynamics of this network are described as a set of possible rewriting events (which is written in an external spreadsheet file for convenience of frequent editing). A rewriting event is defined as an establishment of a new link between two agents, possibly involving changes of their states. Each possible event is specified by the following eight properties:

- 1) Conditions: (Optional) Logical expression(s) that indicate when this event can be executable.
- 2) Source: Agent from which the new link departs.
- 3) Destination: Agent to which the new link points.
- 4) Link type: Type of the event (i.e., interaction between the two agents). The following three types are allowed:
 - "Request": The source agent requests the destination agent for specific information.
 - "Flow": The source agent sends specific information to the destination agent.
 - "Task": The source agent commands the destination agent to do a particular task.
- 5) Knowledge required: (Optional) List of internal variables the source agent needs to have in order for the event to occur.
- 6) Knowledge transferred: (Optional) List of internal variables whose values are requested or shared between the two agents during the "Request" or "Flow" event.
- 7) Duration: Amount of time the event takes.
- 8) Duration variation: Amount of stochastic variation for the duration.

Our simulation software reads the set of possible rewriting events given in the above format. The algorithm of simulation of this network proceeds in the following steps:

- 1) Select all the events that are currently executable (i.e., all conditions are met and the source agent has all the knowledge required).
- Make the selected events active and set a duration time (with stochastic variation added according to the duration variation property of the event) to their respective internal time counters.
- 3) Decrease the time counters of all of the active events by a unit time.
- 4) If the counter of any of those active events hits zero, establish a new directed link from the source agent to the destination agent in the network. Also, depending

on the type of the event, update the internal variables of both agents. Then deactivate the event.

5) Repeat the process above until no more executable events exist.

The operational network will emerge as the simulation progresses and more agents are connected by information exchange and task allocation. Our simulation software implements interactive GUI by which the user can operate and inspect the simulation status. Figure 5 presents the current design of the GUI.

3. Conclusions

In this paper, we have presented a modeling framework that combines visualization and simulation capabilities for representation and analysis of discrete spatially-distributed systems. Agents or active players in such systems are highlyspecialized entities that differ by their roles/specializations, environmental domains/realms, and also by other casespecific attributes.

This study is concerned with the systems that perform two main functions: detection of a significant event and execution of an appropriate response action. Under our approach, the system is viewed and modeled as a dynamical network with embedded heterogeneous agents. The emphasis is given to the relationships (links) and to the dynamics of how these relationships change in the course of a response action.

Each agent is encoded as a multidimensional string σ (Eq. 1), where each position σ_i represents an agent attribute or property (Subs. 2.1). The first five positions of σ represent agent classes: sensor, router, actor, database, and controller. The next six positions denote realms in which agents operate: maritime, land, air, space, cyber, and cognitive. Substring $[\sigma_{12}, \ldots, \sigma_{15}]$ is reserved for operational domains: Air, Ground, Maritime, and Joint. The number of dimensions is not limited and can include as many case-specific attributes as needed. Agents that belong to the same agent class, realm and operational domain (i.e. they are identical in $[\sigma_1, \ldots, \sigma_{15}]$) are grouped in heterotypes. The number of heterotypes within the network is used to measure the network heterogeneity (Eq. 3). Heterotypes are also used as auxiliary constructs in simulations.

We have shown how the visualization capabilities of the modeling framework can be used to visualize the architectural make-up of the system and to explore its multidimensional agents. Assortative networks have been invoked to group and display agents, according to their attributes, as represented by string σ (see Subs. 2.2). The concepts of the Standby network has been used to graphically represent the system architecture (see Subs. 2.3).

In addition to the visualization capabilities, discussed in Subs. 2.2 – 2.3, the modeling framework allows for a straightforward application of the Generative Network Automata (GNA) – a network generation technique for automated generation of new networks depending on the



Figure 5: Graphical User Interface of the GNA-based simulation software.

dynamical state of the nodes [9], [10]. We have adopted and revised the GNA approach to develop a simulation software for generation of operational networks on standby networks of the SARnet dynamical network model of the Canadian Arctic Search and Rescue System [7]. The simulation algorithm is described in Subs. 2.4.

The presented modeling framework provides a virtual laboratory in which the performance of discrete distributed heterogeneous systems can be efficiently evaluated. Such performance aspects as resilience to technical failure or attack can be efficiently assessed through simulating the generation of operational networks on standby networks in 'what-if' scenarios.

Other potential applications include tools for automated generation of coordinated rapid response, such as SAR response, which requires a quick identification of appropriate responders among multiple diverse assets.

Acknowledgments

The development of the GNA-based simulation software was supported by the Canadian Government contract W7714-125419. The GNA theoretical modeling by Sayama and Wong was supported in part by the U.S. National Science Foundation (Award # 1027752).

References

- M. Newman, A.-L. Barabási and D.J. Watts, *The Structure and Dynamics of Networks*, Princeton University Press, Princeton, 2006.
- [2] M. Aldana, S. Coppersmith and L. P. Kadanoff, "Boolean dynamics with random couplings," arXiv:nlin.AO/0204062 (electronic publication), vol. 2, pp. 1–71, 2002.
- [3] T. Gross and H. Sayama, "Adaptive Networks," In: Adaptive Networks: Theory, Models and Applications, (T. Gross and H. Sayama, Eds.), pp. 1–11, 2009.
- [4] K. M. Carley, "Dynamic Network Analysis," In: *The NRC Workshop on Social Network Modeling and Analysis*, (R. Breiger and K.M.Carley, Eds.), National Research Council, pp. 1–13, 2008.
- [5] M. E. J. Newman, "The structure and function of complex networks," SIAM Review, Vol. 45, No. 2, pp. 167–256, 2003.
- [6] I. Pestov and M. Pilat, "Modelling Search and Rescue Systems with Dynamical Networks," Proc. IEEE Symp. Series on Computational Intelligence, W/shop on Computational Intelligence in Security and Defence Applications (CISDA 2011), 8 pp., 2011.
- [7] I. Pestov and M. Pilat, "A Net-Enabled Approach to Modelling Joint Interagency Systems: The Arctic Search and Rescue Network," DRDC CORA TM 2012-067, Ottawa, Canada, 86 pp., March 2012.
- [8] K.M. Carley, D. Columbus, M. DeReno, J. Reminga and I.C. Moon ORA User's Guide 2008., CMU-ISR-08-125, Carnegie Mellon University, Pittsburgh, PA, 2008.
- [9] H. Sayama, "Generative network automata: a generalized framework for modeling complex dynamical systems with autonomously varying topologies," Proc. IEEE Symposium on Artificial Life (CI-ALife 2007), pp. 214–221, 2007.
- [10] H. Sayama and C. Laramee, "Generative Network Automata: A Generalized Framework for Modeling Adaptive Network Dynamics Using Graph Rewritings," In: *Adaptive Networks: Theory, Models and Applications*, (T. Gross and H. Sayama, Eds.), pp. 311–332, 2009.
- [11] NetworkX High productivity software for complex networks, Los Alamos National Laboratory (http://networkx.lanl.gov).

Modeling of DNA Replication

Xiaoli Yang¹, Rong Ge¹, Yifan Cai¹ and Charles Tseng²

¹Department of Electrical and Computer Engineering ²Department of Biological Sciences Purdue University Calumet Hammond, IN, USA

Abstract - DNA replication is a necessary step prior to cell division, so that the genetic material can be duplicated for equal distribution in the daughter cells. Although in the course of evolution, cells have developed specific mechanisms to ensure the fidelity of the process, faulty enzymes and mutagens may cause changes in DNA sequences, leading to a variety of diseases including cancers. As important as the DNA replication process is, however, teaching and learning of the subject have been difficult. The present paper describes an innovative computer program that stresses inquiry based learning through visualization, cognitive feedback and handson interactions. It is one of a series of interactive computer modules for learning genetics at both high school and college levels.

Keywords: DNA replication, modeling, computer program

1 Introduction

DNA replication is a fundamental property of all living organisms. Prior to cell division, DNA must be replicated, so that after cell division, each of the resulting daughter cells ends up with the same amount of genetic material as the original cell. This process ensures the constancy and continuity of genomic DNA during cell reproduction. Like most biological processes, the detailed mechanisms of DNA replication have not been completely worked out, although a great deal of effort has gone into their elucidation. Our current knowledge is based primarily on the study of bacteria such as E. coli (1). However, since similar proteins involved in DNA replication have also been identified in eukaryotic cells (e.g., yeast and other eukaryotic cell cultures), it seems safe to say that the major DNA replication processes in prokaryotic and eukaryotic cells are similar except for minor details (2).

Although DNA replication is a subject that is taught in both high school and college biology courses, students of all levels still find this subject difficult. Part of the difficulty lies in the intricate and abstract nature of the molecular processes (3-7). Textbooks these days have detailed illustrations that are quite helpful for learning (1), but in the end, the learning is not active. Recent multimedia tools such as DVDs and computerized animations represent a new way of teaching and learning (8-10). However, these multimedia based learning methodologies do not emphasize the interaction of eyes, mind, and hands in the learning process. The present paper describes an innovative computer program that stresses inquiry based learning through visualization, cognitive feedback and hands-on interactions. The DNA replication module is one of several modules developed (11-15) for learning genetics using an interactive computer program. Our specific aim is to provide a useful learning tool for a number of high school and college level courses in the areas of general biology, genetics, cell biology, and molecular biology.

2 Model Development

2.1 Overview

The DNA model has three levels in its structural hierarchy and is composed of many independent ball-shaped elements. Each element has a position, a color, and a radius. Linked together, the elements can interact with one another and move uniformly. A smoothing algorithm, which adds a square outline to the linked elements, is used to fill the gap between two elements. A string of elements forms a rod, representing a DNA strand (Fig. 1).



Fig. 1. DNA modeling : from balls to rod

2.2 Basic model element: node

The node class represents the basic element of the DNA molecule – the ball. The "ball" is nothing but structural data in a linked list. Every ball along the linked list is regarded as a node. A node is characterized by the Cartesian coordinates X and Y (location), a radius (size), and a color (identity) (Fig. 2). Initially, a node is created according to preset parameters. After creation, the color and the radius remain static, while the coordinates may change from time to time during the simulation process.



Fig. 2. Node class overview

It should be pointed out that the nodes cannot overlap with one another in the coordinate system. Each node

occupies its own location so that there is no ambiguity. To form a bidirectional linked list, every node must have two pointers – one points to the last node, while the other points to the next node. Only the first node (head node) has a void pointer. Every node interacts with others based on these relationships (Fig.3). The interaction between small internal nodes makes the rod move flexibly.



2.3 Smoothing algorithm

As mentioned above, to avoid the appearance of discontinuous segments, a smoothing algorithm, which adds a square outline to the linked nodes, is used to fill the gap between every two nodes. This makes the list more like a rod than a series of balls. To create the square outline, four points need to be determined. Fig.4 shows how the 4 points are calculated using a pair of homothetic triangles.

The known variants are the nodes' coordinates and radius. From the property of homothetic triangles, we know that

$$\frac{2R}{dy} = \frac{R}{a} \qquad \frac{2R}{dx} = \frac{R}{b} \tag{1}$$

Thus,



Fig. 4. Smoothing square and the smoothed rod

Assume that (N1x, N1y) and (N2x, N2y) are the coordinates of the nodes. Point 1 is (N1x-a, N1y-b), Point 2 is (N2x-a, N2y-b), Point 3 is (N2x+a, N2y+b), and Point 4 is (N1x+a, N1y+b).

2.4 Node movement

The node itself is not capable of doing complex movement. In fact, only two types of movement are allowed: 1) teleporting the node to a specified location and 2) connecting a node to another nearby node (what we call stepping-up movements). These two movements are one-time movements; there are no intermediate states during the movements. The algorithm for stepping-up is also based on two homothetic triangles (Fig. 5). In order to step up, SX and SY are calculated as follows:

$$Sx = (Dxy - 2R) * \frac{dx}{Dxy}$$
$$Sy = (Dxy - 2R) * \frac{dy}{Dxy}$$
(3)

where *DXY* is the distance between two nodes and dx & dy the distance along x & y axis. 2R is the sum of the radii of two adjacent nodes.



Fig. 5. The stepping-up algorithm for the node

3 Program Contents

Content design is based on three fundamental concepts: 1) Unlike RNA polymerase, DNA polymerase is unable to initiate synthesis of a new strand de novo, that is, it requires a preexisting primer. The major role of DNA polymerase is, therefore, primer extension. In the cell, the primer is synthesized by the enzyme RNA primase. 2) DNA replication is a protein-controlled process. Numerous proteins are involved in changing the topology of the molecule and separating the two strands of the double helix. The proteins are aggregated in a complex "factory" through which the DNA duplex passes (individual proteins do not "travel" to the duplex) and are recognized/bound by individual proteins in the factory for specific reactions. 3) Due to the antiparallel nature of the DNA duplex, semiconservative replication must proceed in the opposite directions on the two template strands. For the two core enzymes of DNA polymerase to stay together, the lagging strand template moves differently than the leading strand template so that the two core enzymes can perform both strand synthesis without falling apart (see details below).

This module is designed to emphasize inquiry based learning (16); learning is achieved through questioning and hands-on interactions. In each of the learning steps, dynamic models of DNA molecules undergoing changes mediated by various proteins are presented for visualization, cognition, and operation. Completion of the program requires comprehension of the entire concept and thus ensures the success of learning experiences.

3.1 Antiparallel organization and semiconservative replication of DNA

Each of the two intertwined strands of the DNA double helix is made of many basic units called nucleotides, which are composed of a 5-carbon sugar (deoxyribose), a phosphate group attached to the 5'C of the sugar, and a nitrogenous base (A, C, G, or T) attached to the 1'C of the sugar. At the opposite end of the phosphate group is an OH group attached to the 3'C of the sugar. Therefore, each strand of DNA has two ends: The 5'P end and the 3'OH end (Fig. 6).



Fig. 6. DNA double helix with antiparallel organization

After separation of the complementary DNA strands, each strand serves as a template for DNA synthesis. Fig. 7 shows 2 new strands being synthesized in opposite directions. The resulting two DNA duplexes each consists of an old strand and a new strand; this is known as semiconservative replication.



Fig. 7. New DNA strands synthesized in opposite directions.

3.2 Individual steps: protein facilitated DNA replication

DNA replication involves the following steps in sequence: a) Denaturation of double stranded DNA by helicase (Fig. 8a), b) Binding of single strand binding proteins (SSBs) to prevent renaturation of newly separated DNA strands (Fig. 8b), c) Binding of RNA primase to initiate the synthesis of a short RNA primer in the 5' to 3' direction; the strand that serves as a template for continuous DNA synthesis is called the leading strand template, while the strand that serves as a template for discontinuous synthesis is called the lagging strand template (Fig. 8c), d) Extension of the RNA primer by DNA polymerase III (core enzyme) known as DNA synthesis or elongation (Fig. 9); the discontinuous synthesis of the lagging strand is now evident and each fragment of the newly formed DNA is known as an Okazaki fragment (Fig. 9 and 10a), e) Removal of the RNA primers by the enzyme RNase H (Fig. 10b and c), which degrades the RNA nucleotides in the 5' to 3' direction one by one until the last RNA nucleotide which is then removed by an exonuclease (Fig. 10a). f) Filling of the gap (after the removal of the RNA primer) by DNA polymerase I through synthesis of a short piece of DNA (Fig. 10c); the new DNA segment is not connected to the neighboring Okazaki fragment, resulting in a nick that is then sealed by DNA ligase (Fig. 10d), which catalyzes the formation of a phosphoester bond.







Fig 9. Primer extension by DNA polymerase III: leading strand (dark blue, left), lagging strand with two Okazaki fragments (light blue, right)



Fig. 10. Replacement of RNA primer with new DNA. a) 3
Okazaki fragments (right) with an exonuclease (grayish blue structure with spikes) for removing the last RNA primer nucleotide of the oldest Okazaki fragment, b) RNase H (circular gold structure) for removing the RNA primer, c)
DNA polymerase I (orange red oblong structure) for synthesis of a new, short strand of DNA to replace RNA primer, d) DNA ligase (triangular gold) for sealing the nick (phosphoester bond formation).

3.3 Replication fork: the factory of DNA replication

The replication fork is the junction between the double stranded DNA and the newly separated single stranded DNA. Typically, there are two replication forks, one on each side of ori, the origin of replication (the point where DNA starts unwinding for replication). The two replication forks move away from each other until replication is completed. Functionally, the replication fork serves as a factory containing numerous proteins to facilitate the DNA replication process. According to the individual replication steps described above, however, two core enzymes of DNA polymerases III (one for synthesizing the leading strand and the other the lagging strand) move in opposite directions. How can the two core enzymes stay together in the same protein factory? The lagging strand template must first move forward and pass through the replication protein complex (factory) in one direction (from right to left) and then retract backward in the opposite direction (from left to right) during the lagging strand synthesis so that both the leading and the lagging strand syntheses appear in the same direction (Figs. 11a-d).



Fig. 11. a) Two core enzymes of DNA polymerase III are joined by a β -clamp loader with a β -clamp (left); also shown is a DNA duplex with a leading strand template (upper, dark red) and a lagging strand template (lower, light red) which is being encircled by the helicase for separating the duplex, b) While the leading strand is being synthesized as the leading strand template travels through the core DNA polymerase III with a β -clamp, the lagging strand template is curved and a RNA primer is made by the RNA primase, c) The lagging strand template has entered and passed through the core enzyme of DNA polymerase III and is about to be locked by the β -clamp, d) Lagging strand synthesis occurs as its template moves backwards (from left to right).

3.4 Prokaryotic and eukaryotic DNA replications

The typical prokaryotic DNA is circular, and replication starts at *Ori*. The general mode of DNA replication is shown in (Fig. 12).

The typical eukaryotic DNA is linear. During replication, the entire chromosomal DNA molecule may be divided into many segments, each with an *Ori*. In the early S phase of the cell cycle, DNA replication starts at each *Ori* and extends laterally until the replicated DNA duplexes meet and join (Figs. 13).



Fig. 12. Prokaryotic DNA replication: a) a circular DNA duplex, b) DNA replication starting at the top with two

replication forks moving away from each other, c) completion of replication with two daughter DNA duplexes.



Fig. 13. Eukaryotic DNA replication: a) chromosomal DNA synthesis at three origins of replication (*Ori*), b) completion of DNA synthesis (only one daughter DNA duplex is shown).

3.5 Telomere and telomerase

Telomeres are specialized structures located at both ends of eukaryotic chromosomes. Telomeres are important structures that protect and stabilize the chromosomes. Since a telomere is at the end of a chromosome, it contains the 5'end of one strand and 3'end of the other. During DNA replication, lagging strand synthesis requires periodic syntheses of primers ahead of DNA elongation. Once the last RNA primer, which is at the 5' end of the lagging strand, is removed, there is no way that it can be replaced with a new DNA segment. Consequently, eukaryotic DNA gets shorter and shorter after each replication until, eventually, the essential DNA coding sequence near the telomere is affected. In other words, eukaryotic chromosomes become shorter after every cell division until the cell dies.

The enzyme telomerase can elongate eukaryotic DNA at the 3' end so that it can serve as a template for synthesizing a new strand, replacing the lost segment due to the removal of the RNA primer. In normal cells the telomerase activity is relatively low. However, in cancer cells the telomerase activity is high, so that cancer cells may divide and live indefinitely; and chromosomes in cancer cells are not shortened after each cell division. How can telomerase accomplish this task? It turns out that telomeric DNA contains tandem repetitive units at the chromosome ends. In humans, for example, there are tandem repeats of 5'TTAGGG3' totaling 10 to 15 kb long. Since telomerase is a RNA-protein complex, it contains the sequence 3'AAUCCC5', which acts as a template for the 5'TTAGGG' repeat. When telomerase binds to the terminal 3' end of the telomeric DNA, only part of the telomeric RNA is paired with the telomeric DNA; the part near the 3'end of the RNA remains as a free single stranded end, serving as a template for telomeric DNA synthesis. The enzyme then moves to expose the RNA sequence at 3'end as a free template for another round of telomeric DNA elongation. As the process continues, the 3' end of the telomeric DNA is lengthened and serves as a template during next round of DNA replication, recovering the previously shortened DNA (Fig. 14).



Fig. 14. Sliding movement of telomerase (a, b, c) to create a free 5' end of the telomeric RNA template for elongating the 3' end of telomeric DNA to replace the lost 5' end telomeric DNA after the next round of DNA replication

4 Conclusions

DNA replication is an extremely complicated process. It requires coordination of many enzymes to assure the fidelity of replication. Missteps in DNA replication lead a variety of diseases including cancers (17). This important subject is generally taught as a unit right after the introduction to DNA structure and again as an integral part with another unit on mutation in most biology textbooks. Students are usually excited about the subject at the beginning, since DNA is not only the fundamental molecule of life but also related to diseases. After encountering the intricacies of the replication process, however, many find the subject difficult to grasp and become disenchanted with biology, an unfortunate situation that needs to be improved.

This computer program is written with the hope that teaching and learning DNA replication becomes an easy and interesting, motivating beginners and also serving as basis for a variety of topics in genetics. This module has recently been tested in a Genetics course at Purdue University Calumet, along with other genetic modules. The initial feedbacks were positive. The ultimate goal of the project is to complete a whole series of interactive computer modules for learning genetics.

5 References

Leslie Griffiths A. J. F., S. R. Wessler, S. B. Carroll, J. Doebley (2012) Introduction to Genetics Analysis, 10th ed, W. H. Freeman and Co. New York, NY.

[2] Pursell, Z. F., I. Isoz², E-B. Lundström, E. Johansson and T. A. Kunkel (2007)"Yeast DNA Polymerase ϵ Participates in Leading-Strand DNA Replication". Science 317 (5834):127–130.

[3] Tibell, L. A. E. and C. J. Rundgren (2010) "Educational challenges of molecular life science: characteristics and implications for education and research" CBE - Life Sci. Educ. 9: 25-33.

[4] Huang, P. C. (2000) "The integrative nature of biochemistry: challenges of biochemical education in the USA" Biol. Educ. 28:14-17.

[5] Bahar, M., A. H. Johnstone, and M. H. Hansell (1999) "Revisiting learning difficulties in biology" J. Biol. Educ. 33: 84-86.

[6] Brig, J. (1996) "Enhancing teaching through constructive alignment" Higher Education, 32:347-364.

[7] Sheley, S. M. and T. R. Mertens (1990) "A Survey of Introductory College Genetics Courses" J. Heredity 81: 153-156

[8] Essential Biochemistry - DNA Replication

[9] www.wiley.com/college/pratt/.../animations/dna_replica tion/index.ht

- [10] DNA Replication Process-YouTube
- [11] www.youtube.com/watch?v=teV62zrm2P0
- [12] DNA makes DNA Cell Biology Animation
- [13] www.johnkyrk.com/DNAreplication.html

[14] Yang X., G. Rong, C. Tseng (2011) "Modeling of DNA Replication" The 2011 International Conference on Modeling, Simulation and Visualization Methods, p.146-149, Las Vegas.

[15] Wu W., X. Yang, C. Tseng (2011) "Effective Algorithms for Altering Human Chromosome Shapes" The 2011 International Conference on Modeling, Simulation and Visualization Methods, p. 257-261, Las Vegas.

[16] Yang X., R. Ge, Y.Yang, H. Shen, Y. L and C. Tseng (2009) "Interactive Computer Program for Learning the Genetic Principles of Segregation and Independent Assortment through Meiosis" The 31st Annual International Conference of the IEEE Engineering in Medicine and Biology Society (EMBC 2009), p. 5842-5845, Minneapolis.

[17] Wu W., X. Yang, B. Chen, Z. Zhao, J. Lacny and C. Tseng (2009) "Computer Based Simulation of Chromosome Abnormality" The 2009 World Congress in Computer Science Computer Engineering and Applied Computing (WORLDCOMP 2009) p. 359-363, Las Vegas.

[18] Yang, X., D. Wen, Y. Cui, X. Cao, J. Lacny and C. Tseng (2009) "Computer Based Karyotyping" The 3rd International Conference on Digital Society (ICDS 2009), 310-315, Cancun, Mexico.

[19] Inquiry Based Learning: www.thirteen.org/edonline/concept2class/inquiry/

[20] Helleday Thomas, T., E. Petermann, C. Lundin, B. Hodgson and R. A. Sharma (2008) "DNA repair pathways as targets for cancer therapy" Nature Reviews Cancer 8:193-204.

[21] Ree Source Person. "Title of Research Paper"; name of journal (name of publisher of the journal), Vol. No., Issue No., Page numbers (eg.728—736), Month, and Year of publication (eg. Oct 2006).

Exploring Information Stewardship with the Cloud Ecosystem Model

A. Baldwin¹, Y. Beres¹, L. Carrotte², T. Koulouris¹, B. Monahan¹, D. Pym³, S. Shiu¹, and C.Y. Yam¹

¹HP Labs, Bristol, England, UK ²Nomos Media Ltd, Bristol, England, UK ³University of Aberdeen, Scotland, UK

Abstract. The emergence of cloud computing has transformed the way in which enterprise IT is delivered and creates new challenges around risk management, security strategy, and control over policies and information. For a particular example, the economies of scale that can be achieved by large cloud service providers are encouraging ecosystems of service providers in which the marketplace (rather than the consuming enterprises or individuals) determines security standards. The problem goes beyond core security concerns, since all cloud stakeholders will be relying on others to steward their information, and so will be concerned with the overall sustainability and resilience of the ecosystem, from both security and business perspectives. To help cloud consumers and other stakeholders explore the impact of the cloud ecosystem on their business decisions, we have developed a system that combines systems modelling, simulation, and 3D visualization techniques. At the heart of this system is a model of the cloud ecosystem, built by combining economic and system modelling approach. The model uses utility theory as the unifying vocabulary for stakeholders to express their decision-making. Simulations of the model representing several years of operations are then performed, with various shocks — such as economic downturns and security attacks — introduced at certain points in time. The visualization component has been built specifically for this model and consists of interactive 3D graphics that can be used in any compatible web browser, so allowing stakeholders to interact with and explore the model and simulation results easily.

Keywords: Cloud Computing, Information Stewardship, Ecosystem Visualization

1 Introduction

The typical risk management lifecycle involves risk assessment, policy setting, investment programs, design and deployment of controls, procuring and managing infrastructure, and monitoring and audit to ensure controls and policies are effectively mitigating risk. A challenge with cloud computing is that the activities of this lifecycle become disaggregated and are performed by different, third parties, each with different incentives. As control over security policies seeps out, and the organization becomes dependent on multiple stakeholders in the ecosystem, security concerns develop into information stewardship issues [25,11,23].

The challenge for enterprises is how to judge the risks involved in consuming cloud services, and to understand the options, and their consequences, that are available. This challenge is more complex than outsourcing, where consumers have been able to dictate terms and conditions: the large scale and cost structures of cloud providers will tend to allow vendors, and the marketplace, to dictate standard, onesize-fits-all, security service levels. Moreover, in order to leverage scale and associated cost reductions further, service providers will tend to bundle cloud services, standardized offering terms and conditions.

It is not just cloud consumers that are challenged. Each stakeholder in the ecosystem is vulnerable not only to changes that may occur within that ecosystem, but also to any to changes in the external environment that may impact on parts of the ecosystem's operations. For one example, the activities criminals that target a particularly successful enterprise, causing high impact security incidents, may affect many different supply chains. For another, skill shortages and liquidity shocks will affect multiple groups in different ways, with potentially large impact on the whole ecosystem. Such incidents can affect the reputation and trust of many participants, as well as the whole system.

So, all cloud stakeholders must understand the options they have to improve stewardship outcomes. How should regulators and policy makers impose rules and regulations? How much influence does a single consumer have? How does this change if they act as a group? How much transparency into operations should be demanded by consumers and offered by providers? How should all the stakeholders act to deal with factors exogenous to the market, such as the state of the global economy, business trends, technology changes, or shifts of human skills?

The lack of effective ways for stakeholders to explore their assumptions about value, risk, and operational uncertainty will lead to inaccurate perceptions and, potentially, poor decisions. Clearly, the challenge is how to provide guidance and support for stakeholders to comprehend the risks associated with cloud, and so and form good strategic and operational responses.

In this paper, we present an interactive modelling and simulation tool that provides a step towards addressing this

challenge. At the heart of this is a mathematical system model that explores numerous aspects of the emerging cloud ecosystem. The system model consists of (hundreds of) firms consuming IT, (hundreds of) firms offering services, and several cloud platform providers offering IT resource capacity. In addition, the model explores the implications of exogenous and endogenous factors on the ecosystem and on information stewardship. The system model simulates the following: consuming firms switching from internal IT to the cloud, or changing service providers; new service providers entering the market with different cost and security properties; and new platforms offering different conditions for the service providers.

So that various stakeholders, be they cloud consuming organizations or cloud services providers, can easily explore the model and the results obtained from its simulations, we have also developed the associated visualization tool that supports simulation play-back features, and allows interaction with the model using 3D graphics that can be used in any compatible web browser.

Section 2 describes our modelling methodology, combining system models with some elementary ideas from utility theory. Section 3 presents the model developed in order to capture the complex cloud ecosystem. Section 4 covers the visualizations of the model and simulation results. Section 5 reviews some related and Section 6 gives our conclusions.

2 Modelling Methodology

We have developed a methodology for combining economic and system models to help organizations with risk assessment, security analysis, and decision-making [1, 9, 27, 28, 29]. Economic models, represented within system models using simple utility functions, are used to help stakeholders think about and share their preferences and priorities for different business outcomes. We then use structural models to help stakeholders think about and share their assumptions for how different investment and policy choices will affect the outcome. Finally, we use a discrete process simulation tool [2, 3, 27] that allows stakeholders to explore and predict consequences of different assumptions. Figure 1 provides a schematic of this methodology.



Figure 1: A framework for using economic and system models to support organizational decision-making.

We have also conducted a series of customer case studies [4, 5, 9] to develop and refine this process. An early example was to help a large enterprise decide between a range of policies and investments to manage risks from software vulnerabilities [6]. The structural models help ensure shared understanding between stakeholders, so they can discuss, say, whether scheduling is significantly delaying patch testing, or when and how often the assessment team should accelerate patch processes. However, with such a complex system of inter-dependent concurrent processes and actions, it can be very difficult to see or reason about the cause and effects. To address this, we use a simulation-modelling tool, Gnosis [2, 3 26, 7, 27], to create an executable mathematical model of the system. Gnosis builds on an underlying mathematical analysis of systems that is based on structural models of location, resource, and process [2, 7, 26, 27] and stochastic representation of environment [2, 27].

Using Gnosis we can run Monte Carlo-style simulations to explore the interactions and their effect on time to mitigate risks. By varying parameters stakeholders can see the (model) predicted effect of different investment choices. Results are typically shared as histograms showing, say, the difference in time taken to mitigate risk for different investment choices. Further experiments can explore the effect based on different assumptions about the threat environment, or to differentiate on different types of mitigation.

Most security decisions involve multiple trade-offs between mitigating different kinds of risks, maintaining services, and minimizing costs [9, 28, 29]. To frame this issue, we encourage stakeholders to define utility functions that express their preferences between the multiple outcomes. To make this approach deployable in practice, we built on the approaches to decisions with multiple objectives developed by Keeney and Raiffa [8]. We developed some simple tools for preference elicitation and then use the system models to explore the effect of different security choices on these other outcomes [9].

Our experience is that focusing on utility (of outcomes), in the context of system models, provides a constructive way to engage multiple stakeholders (with different knowledge and incentives) in the complex process of risk assessment and choosing security investment and policy. Providing evidence for this is difficult as organizations, people, and problems vary so much. We have done some preliminary studies that suggest our methodology affects the justifications security professionals might use, and which fits with why it might be useful for multi-stakeholder decision-making [10]. We are currently looking at further cognitive studies to generate more precise hypotheses about how and why economic and system modelling affect security decision-making.

3 The Cloud Ecosystem Model

In previous modelling work, using the methodology described above, we have explored security decisions by looking at one or two interacting processes. In the considering decision-making in cloud ecosystems, we must consider a much more complex situation, in which there are many interacting entities that must be modelled. In [24], we have considered using real options modelling techniques from financial economics to examine an individual company's decision as to whether to outsource its IT to cloud. This work suggests, unsurprisingly, that the decision depends on the company's expectation of the value cloud will bring to its operations and the uncertainty about whether the chosen service will deliver that value.

The cloud ecosystem will consist in large numbers of customers and service providers with just a few platform providers. The actions of one entity, and exogenous events, may affect the way the overall system functions through both direct influences and feedback loops. In a previous paper [11], we have suggested a conceptual framework for information stewardship in the cloud and how to model cloud as an ecosystem, drawing on the analogy of ecological ecosystems [12]. In [23], we set this approach into the context of enterprise risk management. Here, we briefly describe elements of a cloud ecosystem model being visualized.



Figure 2: Dynamics of the cloud ecosystem.

Within our cloud model (see Figure 2), we have a number of companies who consume IT in order to run their business processes (for example, accounts payable and supply chain management). Each company has a choice: it can run these using its own internal IT experts and data centers (or outsource them) or it can use a cloud service, in which case it must choose a particular service offering, with associated terms and conditions. There is a group of Software-as-a-Service (SaaS) providers who would offer these services. Perhaps in the past they would have produced shrink-wrapped enterprise software. Here we assume that they have little infrastructure themselves, and instead rely on public cloud providers, such as Amazon, HP, or Microsoft, to provide raw computational power and storage. The service providers need to decide on the terms and conditions they offer to their customers based on their costs and perceived needs of their customers and the infrastructure properties (such as security and resilience) they gain from the cloud platforms. Platform (or Infrastructure-as-a-Service, IaaS) providers must make decisions about the basic technology offerings, as well as when to provision new data centers to create new capacity.

The cloud ecosystem sits within a wider world, and here we look at modelling the overall effects of the economy and the threat environment as stochastic variables. All of these internal and external factors relate in a number of feedback loops that determine how the overall ecosystem functions. This allows us to run the ecosystem in different environments to see how it will fare and to shock the system to see how resilient it is to serious external or internal disruptions. For example, we could look at a credit crunch where investment capital becomes limited — how does that effect the way decisions are made, there outcomes, and hence the overall functioning of the cloud ecosystem.

Our approach to modelling is to describe just enough of the structure of the systems that we are modelling to capture the important aspects of their behaviour for the questions of interest to us. We build models using the Gnosis modelling language [2, 3, 27] — based on concepts of process, resource and location — that allows us to model the structure of distributed systems. Hence when we talk of the feedback loops within the ecosystem our model does not explicitly code them but they emerge due to the processes that each of the entities within the model run.

Figure 3 depicts the various entities, states, and processes involved. A detailed description of the model will be the subject of another paper; here we aim to describe enough to give the reader an idea of how the model being visualized works.

Consider, for example, a company consuming IT services. For each service, there will be a review process that periodically looks at the value the company is getting from its IT provision and whether this could be improved. The idea of value here is wrapped up in a utility function for the company that includes productivity improvements for the business, as well as potential costs arising from risks such as the exposure of information, loss of integrity of the process, loss of availability of the data, or failure to meet regulations. When reviewing a decision perhaps comparing internal IT with different cloud services, the customer will look at how well different options meet its utility along with the costs of the different options. Here, a customer may take a slightly lower utility where the cost is much cheaper.

In making the decision and looking at how its utility is met, the customer will consider what information is available within the system. For example, in assessing risks of moving to the cloud the customer may look at the overall reputation of the cloud along with any views they have on the particular cloud providers they are considering. Reputation figures will be derived from others' experiences and events such as service downtime and security incidents. These are reported through other processes happening within ecosystem, so creating the feedback loops. Other examples of information used in making decisions may be staff costs (dependent on labour availability and rewards for cloud start-ups) or the availability of software to run internally, or the ability to raise capital for investing in new IT systems. When comparing different cloud offerings, a company may look at both its own

utility (how secure it needs the system to be) and the threat environment (including published incidents).

the effects of vendor lock-in in the ecosystem for different mixes of open- versus closed-source cloud service providers.

PLATFORM

State held

Profit Status

State used

Process • Review existence

Profile

 Sector, capital

Economy
Demand for cloud

Resources used

Capital
 Event generated

Event consumed

SERVICE

Process

State held

status

profit

State used

fund

Service type

reputation
Pricing & margin
quality delivered

demand forits service
 Resources used

transaction capacity
 Event generated

Event consumed

regulation

response to service failure response to data leak response to regulation change

response to market competition
 exit market

response to technological change

Shocks–economy, cyber attack,

Service: exit mark et

review profitability
 review utility & cost

review service provision & pricing
 review platform
 Rate platform

total customers (i.e. Providers)
 technology provided (to provider)

cost = operation + incident weighted
 service provision level

Own reputation (for pricing)
 customers' service's status
 customers' service's transaction volume

transaction volume & capacity available
 fund

Exit market

- Stewardship requirement - Growth profile

- Agility





Figure 3: Components of the cloud ecosystem model.

4 Visual Representation

Note that the model described here does not contain a model of the market for cloud services (i.e., prices, supplydemand interaction). That level of economic sophistication will be necessary in future developments, but is not necessary to support the visualization thread of our work, which is the focus of this paper. Here, pricing information is represented using simple stochastic variables.

Currently, we run simulations based on the developed model over a seven-year horizon. In the first instance, we explore scenarios that examine the evolution of the ecosystem and its services based on the different amounts of economic growth (or lack thereof). We plan to follow up with scenarios in which the threat environment worsens and also to explore

Our aim in developing the ecosystem model has been to enable various stakeholders to explore the evolution of the ecosystem and the outcomes of their own decisions under various exogenous and endogenous factors. However, with the model being so complex, our previously developed structural representations did not illustrate the effects of complex firm inter-dependencies and richness of autonomous and group behavior in the desired level of detail.

We have decided to create a totally separate modelvisualization component, specifically aimed at capturing the cloud ecosystem as an element that evolves over time. The requirements are to be able show individual entities in the ecosystem, relationships between entities, activity between entities, progress over simulation time, top-level statistics, trends over time, and external influences on the system.



Figure 4: Visualization of the cloud ecosystem model.

We have devoted considerable effort to selecting the good (helpful, reliable) visual metaphors and operations to make the ecosystem concepts accessible to audiences of various backgrounds, while maintaining richness of the visual information. The developed tool supports multiple displays and concurrent views, including three-dimensional, global views of the cloud supply chains, 'drill-down' views of particulars firms, and graphical representations of statistics of interest at various levels of detail. Visualized simulation scenarios can be controlled in real-time using a dedicated control interface, including jumping to particular points in the simulated timeline, pausing, forwarding or reversing to assist in analyzing model effects. Finally, this visual front-end is supported 'behind the scenes' by a scalable cloud-based simulation engine which handles the massive computational workloads required to execute models of this scale and to render the visuals. Figure 4 shows a screenshot of the 'global view' displayed by the model visualization tool.

The ecosystem is represented as a three-dimensional sphere with consumers as blue dots, service providers in purple, using infrastructure provided by platforms, in pink. The sphere has two atractors, the top representing cloud and the bottom in-house IT service provision. Using this metaphor, platforms cluster around the topmost, with cloud service providers organized in a sphere around them. In turn, consumers are positioned on the surface of the larger sphere according to how much cloud or in-house IT they consume in aggregate. During simulation, consumer firms float upwards as more cloud is consumed or, conversely, gravitate towards the bottom if in-house IT becomes more attractive. This can be demonstrated statically when Figures 4 and 5 are compared. Any changes in the ecosystem, such as the introduction of a shock, the arrival or departure of firms, or the effect of service procurement, have a corresponding effect on the placement of firms in the ecosystem sphere and the reorganization of nodes. A wide range of ecosystem effects, ranging from subtle to large-scale events, can be visualized in this manner.



Figure 5: The cloud ecosystem during a period of high cloud adoption.

Fluid visualization of the ecosystem's behaviour is complemented with the inclusion of a 'headline metrics' panel that provides a quick overview of key 'ecosystem health' indicators.¹ The panel's contents can be rearranged and expanded to assist in monitoring particular effects. Key metrics include the following: cloud usage (aggregate cloud usage as a proportion of global IT transactions); internal IT usage (aggregate in-house IT usage as proportion of global IT transactions); security incidents (the number of overall attempted/successful in cloud versus attempted/successful in internal IT security incidents); IT transactions (overall number of transactions); cloud satisfaction (aggregate satisfaction level, calculated as the difference between desired utility and received utility for each service); switching cost/vendor lock-in (aggregate of switching costs calculated for each consumer); credit availability (externally available credit); economy (value created inside or because of the ecosystem); cloud utilization (cloud usage compared to overall capacity); and cloud reputation (overall reputation of the market ecosystem as viewed by consumers).



Figure 6: Detailed view of an individual consumer firm.

¹ The charts in the presented screenshots show some preliminary results from simulations, and as such are presented for the purpose of example and not for the use of interpretation of the model and results from simulations.

Finally, the tool allows for zooming-in on an individual consumer, and for exploring the supply-chain dependencies developed between it and its service providers and platforms, as well as more specific metrics such as particular utility and satisfaction, as shown in Figure 6.

5 Related Work and Future Directions

There has been considerable amount of work on cloud architectures, and various properties within them, but not as much work exploring the digital cloud ecosystem and the implications on its various stakeholders. A number of cloud service models (e.g., Business-process-as-a-Service, Software-as-a-Service, Platform-as-a-Service, Infrastructure-as-a-Service), as well as cloud deployment and management models (private, public, community and hybrid), have been explored [15, 16, 17, 19]. However, organizations need mechanisms to understand better the implications of these cloud service models and their impact on application, process, and data interoperability. Towards that aim is the work by Briscoe and Marinos [14], in which they present a sociotechnical conceptualization for sustainable cloud computing and explore briefly the tensions between open source and proprietary software use within. Related to this is also the work by Gill and Bunker [18] on the context-aware cloud adaptation framework that is aimed at enabling organizations to better self-assess, select, and adopt an appropriate cloud computing model.

On the other hand, there is the work exploring economic implications and the various economic models in the cloud [20, 21, 22]. However, we are not aware of any work in this area that enables the various stakeholders to explore interactively the cloud ecosystem using models and simulations. Our next aim with the work presented here is to be able to use the model and visualization tool in a scenarioplanning workshops, involving various stakeholders such as IT managers responsible for procuring services from cloud or security practitioners that care about how security and stewardship requirements can be met in the cloud ecosystem.

6 Conclusions

As organizations employ cloud services, they rely on others not only to provide those services, but also to protect their information and appropriately control its interactions and evolution. In the cloud, IT operations will be purchased from highly inter-connected ecosystems of services, consumers, and platform providers. Changes in one part of the ecosystem can affect many other parts, in complex ways that will, typically, be difficult to conceptualize. We contend that modelling and simulations, of the kind we have described in this paper, together with interactive visual tools can help decision-makers to understand these complex relationships and dependencies. From the perspective of managing the security lifecycle, organizations can use this information to understand how different events in different components of the ecosystem may affect the systems for which they are responsible. From the wider perspective of the stewardship of the ecosystem itself, this approach can help explore how the ecosystem can be managed to be sustainable and resilient.

Acknowledgements

This work has been partially supported by the 'Cloud Stewardship Economics' project [13], funded by the Technology Strategy Board of the UK Government. We are grateful to Marco Casassa Mont, Christos Ioannidis, Martin Sadler, and Julian Williams for valuable discussions of our ideas.

7 **References**

[1] Pym, D., Shiu, S., Coles, R., van Moorsel, A., Sasse M. A., Johnson, H.: Trust Economics: A systematic approach to information security decision-making. Final report for the UK Technology Strategy Board-funded 'Trust Economics' project. Hewlett-Packard, 2011. http://www.hpl.hp.com/news /2011/oct- dec/Final_Report_collated.pdf

[2] Collinson, M., Monahan, B., Pym, D.: Semantics for Structured Systems Modelling and Simulation. *Proc. Simultools 2010*. ACM Digital Library and EU Digital Library. ISBN: 978-963-9799-87-5.

[3] Core Gnosis. System available for download (subject to HP licence agreement) at: http://www.hpl.hp.com/research/systems_security/gnosis.html.

[4] Baldwin, A., Casassa Mont, M., Shiu, S.: Using Modelling and Simulation for Policy Decision Support in Identity Management. *Proc. IEEE Policy 2009 Symposium, London*: 17-24.

[5] Squicciarini, A.C., Rajasekaran, S. D., Casassa Mont, M.: Using Modelling and Simulation to Evaluate Enterprises' Risk Exposure to Social Networks. IEEE Computer 44(1), 66-73, January 2011.

[6] Beres, Y., Griffin, J., Shiu, S., Heitman, M., Markle, D., Ventura, P.: Analyzing the performance of security solutions to reduce vulnerability exposure windows. *Proc. 2008 Annual Computer Security Applications Conference* (ACSAC), IEEE Computer Society Press, 2008, 33-42.

[7] Collinson, M., Monahan, B., Pym, D.: A logical and computational theory of located resource. *Journal of Logic and Computation* 19(b): 1207-1244, 2009.

[8] Keeney, R. L. and Raiffa, H.: Decisions with Multiple Objectives: Preferences and Value Tradeoffs. Wiley, New York, 1976. Reprinted, Cambridge University Press, New York, 1993.

[9] Beres, Y., Pym, D., Shiu, S.: Decision Support for Systems Security Investment. In *Proc. Business-driven IT* *Management* (BDIM) 2010, Network Operations and Management Symposium Workshops. IEEE Xplore, 2010.

[10] Shiu, S., Baldwin A., Beres, Y., Casassa Mont, M, Duggan, G.,Johnson, H., Middup, C.: Economic methods and decision making by security professionals. In Bruce Schneier (editor), *Proc. 10th Workshop on the Economics of Information Security (WEIS)*, 2011. Springer, 2012. In press. Preprint available at http://weis2011.econinfosec.org/papers/ Economic%20methods%20and%20decision%20making%20b y%20security%20profession.pdf.

[11] Baldwin, A., Pym, D., Sadler M., and Shiu, S.: Information stewardship in cloud ecosystems: towards models, economics and delivery. *Proc. 3rd IEEE International conference on Cloud Computing*, Athens, 2011. IEEE Conference Publications, 784-791, 2011. doi: 10.1109/CloudCom.2011.121.

[12] Chapin III, F. S., Kofinas, G. P., Folke C., (editors): *Principles of Ecosystem Stewardship: Resilience-Based Natural Resource Management in a Changing World*. Springer-Verlag, 2009

[13] The 'Cloud Stewardship Economics' project, IISP, https://www.instisp.org/SSLPage.aspx?pid=463.

[14] Briscoe, G., Marinos, A.: Digital Ecosystems in the Clouds: Towards Community Cloud Computing. *Proc 3rd IEEE International Conference on Digital Ecosystems and Technologies*, 2009.

[15] Mell, P., Grance, T.: The NIST Definition of Cloud Computing, 2009. NIST Special Publication 800-145. http://csrc.nist.gov/publications/nistpubs/800-145/SP800-145.pdf

[16] Shroff, G.: *Enterprise Cloud Computing: Technology, Application and Architecture.* Cambridge University Press, 2010.

[17] Whyld, D.C.: Moving to the Cloud: An Introduction to Cloud Computing in Government, 2010. http://www.business ofgovernment.org/report/moving-cloud-introduction-cloud-computing-government.

[18] Gill, A.Q., Bunker, D.: Conceptualization of a Context Aware Cloud Adaptation Framework. *Proc. Ninth IEEE International Conference on Dependable, Autonomic and Secure Computing*, 2011.

[19] Lenk, A., Klems, M., Nimis, J., Tai, S., Sandholm, T.: What's Inside the Cloud? An Architectural Map of the Cloud Landscape. ICSE Workshop on Software Engineering Challenges of Cloud Computing, CLOUD '09, 2009.

[20] Lindner, M.A, Vaquero, L. A, Rodero-Merino, L., Caceres, J.: Cloud economics: dynamic business models for

business on demand. International Journal of Business Information Systems, 2010.

[21] Yam, C-Y., Baldwin, A., Ioannidis, C., Shiu, S.: Migration to Cloud as a Real Option: Investment decision under uncertainty. *Proc. IEEE TrustCom 2011 Symposium & Workshops*.

[22] Hongyi Wang, Qingfeng Jing, Rishan Chen, Bingsheng He, Zhengping Qian, Lidong Zhou: Distributed systems meet economics: pricing in the cloud. *Proc. 2nd USENIX Conference on Hot Topics in Cloud Computing* (HotCloud'10), 2010.

[23] Baldwin, A. Pym, D. and Shiu, S.: Enterprise information risk management: Dealing with cloud computing. To appear: Privacy and Security for Cloud Computing: Selected Topics, Siani Pearson and George Yee (editors), Springer, Computer Communications and Networks series, 2012.

[24] Yam, C-Y. Baldwin, A., Ioannidis, C., Shiu, S.: Migration to Cloud as a Real Option: Investment decision under uncertainty. In Proc. IEEE TrustCom 2011 Symposiums & Workshops.

[25] Pym, D., Sadler, M.: Information Stewardship in Cloud Computing. *International Journal of Services Science, Management, Engineering, and Technology* 1(1), 50-67, 2010.

[26] Collinson, M., Pym, D.: Algebra and Logic for Resource-based Systems Modelling. Mathematical Structures in Computer Science 19:959-1027, 2009.

[27] Collinson, M., Monahan, B., Pym, D.: A Discipline of Mathematical Systems Modelling. College Publications, 2012.

[28] Ioannidis, C., Pym, D., Williams, J.: Investments and Trade-offs in the Economics of Information Security. In *Proc. Financial Cryptography and Data Security* 2009, LNCS 5628: 148-162, Springer, 2009.

[29] Ioannidis, C., Pym, D., Williams J.: Information Security Trade-offs and Optimal Patching Policies. *European Journal of Operational Research*, 216(2):434-444, 2012. doi:10.1016/j.ejor.2011.05.050.

Modeling of the Movement of the Endoscopy Capsule inside G.I. Tract based on the Captured Endoscopic Images

Guanqun Bao, Yunxing Ye, Umair Khan, Xin Zheng and Kaveh Pahlavan

Center for Wireless Information Network Studies Worcester Polytechnic Institute Worcester, MA, 01609, USA (gbao, yunxingye, uikhan, xzheng, kaveh)@wpi.edu

Abstract - Wireless capsule endoscopy (WCE) is a noninvasive technology that provides excellent images of the intestinal lumen as the capsule moving along in the gastrointestinal (G.I.) tract. However, the biggest drawback of this technology is its incapability of localizing the capsule when an abnormality is found by the video source. Existing localization methods based on radio frequency (RF) and magnetic field suffer a great error due to the non-homogeneity of the human body and uncertainly of the movement of the endoscopic capsule. To complement the existing localization techniques, in this paper, we developed a series of image processing and visualization based algorithms to model the movement of the endoscopic capsule. First, a 3D map of the G.I. tract is generated to navigate the transition of the capsule. Then, by comparing the local similarity and feature matching between the consecutive frames, the speed and rotation angels of the capsule can be roughly estimated. Finally, by mapping the pattern of the capsule's movement onto the 3D G.I. tract map, we are able to simulate the entire transition of the endoscopic capsule in the 3D space. Empirical results show that the proposed method has a good estimation of the capsule's movement.

Keywords - capsule endoscopy, localization, movement modeling, rotation estimation

1 Introduction

Wireless capsule endoscopy (WCE) [1] has been in the clinical arena for 12 years. It provides a noninvasive imaging technology of the entire G.I. tract. Current devices, for example devices developed by Given Imaging and Olympus American, are able to provide excellent images of the lumen of the intestines as they moving along in the intestinal tract. However, none of these devices could provide accurate location information of the capsule when an abnormality is found by the video source. An accurate measurement of the capsule's position is of great benefit to the physicians in terms of reducing diagnosis time and taking immediate clinical management to obscure gastrointestinal bleeding [2]. During the past few years, many efforts had been done to develop reliable localization technique inside human body. The Given Image Pillcam capsule [3] was originally developed with the potential capability of localizing the capsule on a 2-D plane at a twice-per-second rate. Eight external antennae were fixed to the anterior abdominal wall to detect the UHF-band signal that is emitted by the capsule. The position of the capsule on the 2-D plane of the abdomen is estimated depending on the signal strength received by each antenna with an accuracy of 6 inches. However, the clinical use of this software found that the crude localization result generated by the software was not helpful and this approach was soon abandoned. Another commonly used approach for capsule localization is to assume the capsule travels at a constant speed and the approximate position of the capsule is calculated according to the time of travel away from some pre-defined land marks such like pylorus and ileocecal valve. Apparently, when using this approach, the further the capsule moves away from the land marks, the greater the likelihood of error is. Especially after the video capsule has entered a few centimeters of the small intestine, the localization error will increase dramatically. This is mainly due to the high complicity level of the shape of the small intestine. The distribution of small intestine is like a curled snake with its length varies from 4.6m to 9.8m [4] (the average value for human being is 7m) and the tendency of loops is highly undistinguishable. Besides, the intestinal motility is not consistent. Peristalsis may make the wireless capsule sometimes move quickly, sometimes stop or sometimes even reverse and then progress with any combination of the movement above. Furthermore, the transition of the capsule itself is not axial. It may rotate with different angles or get flipped by 180°. The unpredictable angulation of the wireless capsule creates difficulties in RSS based localization in terms of changing the antenna gain. Thus, knowing how the capsule moves inside human body will help us analyzing the radio channel and thereby enhance the accuracy of the RF localization. Besides. considering the 3-dimensional distribution of the small intestine, a 3D localization technique instead of 2D is also needed to provide sufficiently accurate spatial location information of the capsule.

The rest of the paper is organized as follows: In section 2, we generated a 3D intestinal tract map by extracting the central line of the existing 3D G.I. models. This map provides us a clear view of how the capsule transits inside human body. Then, by comparing the local similarity and matching the feature points between the consecutive endoscopic frames, information such like speed and rotation angles of the endoscopy capsule can be roughly estimated to model the capsule's movement inside the G.I. tract. In section 3, both the experimental results and analysis of the proposed method are given explicitly. Conclusion and future work are addressed in section 4.

This research was funded by the National Institute of Standards and Technology (NIST), USA, under contract FON 2009-NIST-ARRA-MSE-Research-01, entitled "RF Propagation Measurement and Modeling for Body Area Networking".

2 Movement of the Endoscopy Capsule

2.1 3D Map Generation

In every localization technique, map always plays an important role in terms of refining the localization results. Existing literature [5] reported that a clear street map is able to reduce the GPS localization error from tens of meters to several meters in the urban area. In case of the localization inside human body, "map" is even more important since everything goes through the G.I tract follows the same route. Knowing a clear pattern of the intestinal tract will greatly enhance the localization accuracy. Therefore, tracing the path of intestinal tract is essential to the accurate capsule localization. Given a 3D CAD model of the G.I. tract as shown in Fig. 1, we want to trace the center of the intestinal volume so we can model the movement of the capsule inside the tract. In this sub section, some 3D image processing techniques are applied to accomplish this goal. For the large intestine (as shown in the middle of Fig. 1), since it already has a very clear pattern which looks like a big hook, we applied 3D skeletonization technique [6] to extract the path of it. As for the small intestine, since the shape of the small intestine is much more complicated (the trend of the small intestine can be hardly recognized by human eyes), we developed an element sliding technique to trace the path. The basic idea behind this technique is to define an element shape (ES) with its radius automatically adjustable to the radius of the small intestine. This ES is propelled forward by a factor associated proportional to the average distance between the vertices within certain range and the physical center of the ES. As the ES goes along the small intestine, the position of its physical center is recorded and this will give us a clear path of the small intestine. The preliminary result of the path extracted from the 3D model is shown on the right of Fig. 1.



Fig.1. 3D path generation from a 3D G.I. tract model

2.2 Speed Estimation of the Endoscopy Capsule

After the endoscopy capsule was swallowed by the patient, it travels through the G.I. tract propelled by peristalsis [4]. From the dataset we observe that whenever the intestinal lumen contracts, the difference between consecutive frames is high, which to some level reflects a higher speed of the capsule's movement. Based on this observation, we developed an automatic peristalsis detection method based on the color histogram similarity between the two consecutive frames [7]. The reason why we choose color histogram as the similarity metric is that it's more robust in this particular application and insensitive to the texture noise caused by the bubbles and little pieces of food. However, the overall color histogram is not a good descriptor of the color feature since frames with different contents may have the similar histogram distribution. One example is given in Fig.2, although the contents of the two frames are totally different, the overall histograms of the two still look similar. Thus, instead of using the overall color histogram, we developed a local color histogram comparison algorithm. The algorithm divides the captured image into 16 non-overlapping blocks and calculates the local similarity for each pair of the corresponding blocks. Besides, the original frames captured by the endoscopy capsule are encoded by RGB color space, which is difficult to describe the nature of a color by the amounts of each channel. Thus, we convert the histograms into HSV color space, which is more similar to human's visual perception [8, 9]. Then, the similarity between two frames can be calculated by:

$$Sim(Img_1, Img_2) = \frac{1}{M \times N} \sum_{i=1}^{M} \sum_{j=1}^{N} \left(1 - \frac{|H_{i1}(j) - H_{i2}(j)|}{Max(H_{i1}(j), H_{i2}(j))} \right)$$
(1)

where H_1 and H_2 represent the color histograms of the corresponding blocks. M is the number of the non-overlapping blocks which equals to 16. N is the sample number of the histogram which equals to 768. Since the WCE uses a pinhole camera, when calculate the similarity, only the area covered by the red circle (as shown in Fig. 3) is under consideration. A sample partitioned endoscopic frame is shown in Fig.3.





Fig.3. Local similarity comparison based on partitioned endoscopic image

2.3 Rotation Estimation of the Endoscopy Capsule

If the similarity between two frames is high, it can be inferred that the endoscopy capsule stayed still in the same position or just rotated slightly during measured time slot. In the case when the same pattern appears in both frames, it's possible to calculate the rotation angel between the two frames by feature matching. First, a set of feature points needs to be identified in the reference (first) frame. The simplest statistical measurement is to calculate the variance σ^2 of grey levels in a square neighborhood ($P \times P$, P=16) centroid on a pixel. By sliding the square window back and forth on the frame, points with maximum local grey level variance (mostly on the edges) are selected to be the initial feature points.

$$\sigma^{2} = \frac{1}{P^{2}} \sum_{a=-P/2}^{P/2} \sum_{b=-P/2}^{P/2} (f(i+a,j+b)-\mu)^{2}$$
(2)

$$\mu = \frac{1}{P^2} \sum_{a=-P/2}^{P/2} \sum_{b=-P/2}^{P/2} f(i+a,j+b)$$
(3)

where σ^2 is the variance of the local grey level and μ is the average grey level within the sliding window.

Then, to find the corresponding feature points in the second frame, a cross-correlation matching technique [10] is applied. In signal processing, cross-correlation is a classical method of estimating the degree to which two series of signals are correlated. In 2D pattern recognition, cross-correlation can be used for identifying the target pattern in the image. Consider the image below in black and the mask shown in red. The mask is centered at every pixel in the image and the cross correlation is calculated, this forms a 2D array of correlation coefficients. The un-normalized correlation coefficient at position (i, j) on the image is given by:

$$r(i,j) = \sum_{a=-Q/2}^{Q/2} \sum_{b=-Q/2}^{Q/2} \left(mask \left(a + \frac{Q}{2}, b + \frac{Q}{2} \right) - \overline{mask} \right) (f(i+a,j+b) - \overline{mage})$$

$$(4)$$

where Q is the size of the mask, f(i, j) represents the intensity value at (i, j), \overline{mask} is the mean value of the mask pixels and \overline{mage} is the mean value of the image pixels covered by the mask. The mask itself is a cropped image which needs to have the same appearance as the pattern to be found. If the mask and the pattern being sought are similar, the cross correlation between the two will be high. The peak r(i, j) is the position of the best match in the searching image.



Fig.4. Cross correlation sliding window

Knowing the positions of the corresponding points in two consecutive frames, to find how much the capsule rotated, we need to calculate the rotation matrix between the two frames. Rotation matrix, generally indicated by \mathbf{R} , is a 3×3 matrix shown as below:

R =

$$\begin{bmatrix} \cos\theta\cos\phi & \sin\psi\sin\theta\cos\phi - \cos\psi\sin\phi & \cos\psi\sin\theta\cos\phi + \sin\psi\sin\phi\\ \cos\theta\sin\phi & \sin\psi\sin\theta\sin\phi + \cos\psi\cos\phi & \cos\psi\sin\theta\sin\phi - \sin\psi\cos\phi\\ -\sin\theta & \sin\psi\cos\theta & \cos\psi\cos\theta \end{bmatrix}$$
(5)

where ϕ , θ , and ψ represent the rotation angel around x-axis, y-axis and z-axis respectively.

There are many ways to recover the rotation matrix. In this paper, we used Singular Value Decomposition (SVD) [11] method due to its easy implementation. The basic principle of SVD is to decompose a matrix (defined as H, in this particular application H is a 3×3 square matrix) into 3 separate matrices:

$$[\boldsymbol{U}, \boldsymbol{S}, \boldsymbol{V}] = SVD(\boldsymbol{H}) \tag{6}$$

$$H = U S V^T$$
(7)

where the columns of U are the left singular vectors, S has singular values and is diagonal and V^T has rows that are the right singular vectors. The SVD represents an expansion of the original data in a coordinate system where the covariance matrix is diagonal.

The next step involves accumulating a matrix, called H. One thing needs to be pointed out during this step is the recenter of both dataset so that both centroids can be placed at the origin, like shown below:

$$\boldsymbol{H} = \sum_{i \in N} (P_A(i) - \bar{A}) (P_B(i) - \bar{B})^T$$
(8)

$$\bar{A} = \frac{1}{N} \sum_{i=1}^{N} P_A(i) \tag{9}$$

$$\bar{B} = \frac{1}{N} \sum_{i=1}^{N} P_B(i)$$
 (10)

where P_A and P_B are corresponding point sets in first frame and second frame respectively in $[x, y, z]^T$ style. This step removes the translation component, leaving only the rotation component to deal with. After *H* is factorized by SVD, the rotation matrix can be calculated by multiplication of *V* and U^T :

$$\boldsymbol{R} = \boldsymbol{V} \, \boldsymbol{U}^T \tag{11}$$

where **R** can be expressed in the following form:

$$\boldsymbol{R} = \begin{bmatrix} R_{11} & R_{12} & R_{13} \\ R_{21} & R_{22} & R_{23} \\ R_{31} & R_{32} & R_{33} \end{bmatrix}$$
(12)

From the matrix **R**, the rotation angel ψ around the optical axis, as illustrated in Fig. 5, can be calculated by:



Fig.5. Rotation of the capsule around main optical axis

3 Experimental Results

According to the clinical data provided by Umass Memorial Hospital, the average transit time of the capsule from the duodenum to the cecum was 240 ± 40 min. Our experimental results show that average 2.4 peristalsis were detected per min from the endoscopic video and each peristalsis takes up around 6 ± 2 seconds. In terms of time, around 20% of the time the similarity between consecutive images drops below 65%, which means the capsule proceeding very fast propelled by the peristalsis. Around 30% of the time the similarity between consecutive images stays beyond 75%, which means the capsule either stays still or rotates very slowly. Fig.6 shows a sample video clip of 50 seconds. During this transit time, two peristalsis were detected.



Fig.6. Speed estimation according to the similarity between consecutive images

Fig. 7 shows two consecutive endoscope images and the corresponding feature points found by the algorithm introduced in section 2. Overall 52 feature points were detected in the example shown in Fig.6 where circles indicate the original positions of the feature points and lines pointed to the rotated positions in the next frame. The calculated results for the rotation matrix and rotation angle are shown in table 1.



Fig.7. Rotation Estimation by feature matching between two consecutive images

Table 1. Calculated parameters for the example shown in Fig. 7

| Feature Points Detected | Rotation Matrix | Translation Matrix | Rotated Angle |
|-------------------------------|--|---|------------------|
| 52 | $\begin{bmatrix} 0.9635 & 0.2677 & 0 \\ -0.2677 & 0.9635 & 0 \\ 0 & 0 & 1 \end{bmatrix}$ | $\boldsymbol{T} = \begin{bmatrix} -14.8879\\ 19.3272\\ 0 \end{bmatrix}$ | 15.5354 |



Fig.8. (a) original frame (b) rotated by $10^{\rm o}$ (c) rotated by $20^{\rm o}$ (d) rotated by $30^{\rm o}$

| Table 2. Calculated | parameters for the | example shown | in Fig. 8 |
|---------------------|--------------------|---------------|-----------|
| | 1 | 1 | 0 |

| Real Rotation Angel | Feature Points Detected | Rotation Matrix | Calculated Rotation Angle | Error |
|---------------------------|-------------------------------|---|---------------------------------|----------|
| 10° | 79 | $\begin{bmatrix} 0.9852 & -0.1713 & 0 \\ 0.1713 & 0.9852 & 0 \\ 0 & 0 & 1 \end{bmatrix}$ | 9.8697° | 0.1303 |
| 20° | 18 | $\begin{bmatrix} 0.9244 & -0.3814 & 0 \\ 0.3814 & 0.9244 & 0 \\ 0 & 0 & 1 \end{bmatrix}$ | 22.4324° | 2.4324 |
| 30° | 10 | $\begin{bmatrix} -0.9065 & -0.4222 & 0 \\ 0.4222 & 0.9065 & 0 \\ 0 & 0 & 1 \end{bmatrix}$ | 180.0913° | collapse |

To verify the accuracy of the results, we artificially generated a set of rotated frames with different angles from the same endoscopic image. The rotated angles are 10° , 20° and 30° respectively. From Table 2 it can be seen that when the rotation angle is small ($<30^{\circ}$), the calculated results are very accurate with an average error below 1.2° . However, as the rotation angle goes up, the accuracy drops down, finally the algorithm collapses at the rotation angle $>30^{\circ}$.

4 Conclusion

In this paper, we presented a novel image processing based approach to analyze the movement of the endoscopy capsule inside human body. The major contribution of this work includes introducing the concept of "3D map" into the localization inside human body and modeling the movement of endoscopy capsule. The proposed technique is very easy to implement, low cost, and with high accuracy. No extra device is needed for this technique other than the video camera itself. The experimental results show that the proposed speed and rotation estimation methods have a good performance especially when the capsule moves slowly. In the future, we will focus on refining this algorithm according to the clinical data and combining this technique with the existing RF localization approaches to provide a hybrid solution to the localization inside human body.

Acknowledgement

The authors would like to thank Dr. David Cave for his precious suggestions and the colleagues at the CWINS laboratory for their directly or indirectly help in preparation of the results presented in this paper.

Reference

- [1] D. Faigel and D. Cave, "Capsule Endoscopy", November 29, 2007
- [2] L. Liu, C. Hu, W. Cai and M.Q.-H. Meng. "Capsule endoscope localization based on computer vision technique," Engineering in Medicine and Biology Society, Annual International Conference of the IEEE. pp. 3711-3714, Nov. 2009
- [3] F. De Lorio, C. Malagelada, F. Azpiroz, M. Maluenda, C. Violanti, L. Igual, J. Vitri &, J.-r. Malagelada. "Intestinal motor activity, Endoluminal Motion and Transit", Neurogastroenterology & Motilit, Vol. 21, Issue 12, pp. 1264-1271, Dec. 2009
- [4] R. Fu, Y. Ye, K. Pahlavan and N. Yang, "Doppler Spread Analysis of Human Motions for Body Area Network Applications", 22nd Annual IEEE international symposium on personal, indoor and mobile radio communications PIMRC 2011, Toronto, Canada, Sep. 2011
- [5] K. Pahlavan and P. Krishnamurthy, "Principles of Wireless Networks", Prentice Hall, 2002
- [6] S. Andrei, L. Thomas, S. Ariel and K. Leif. "On-the-fly Curve-skeleton Computation for 3D Shapes", Computer Graphics Forum, Vol. 26, Issue 3, pp. 323-328, Sep. 2007
- [7] C. Poh, Z. Zhang, Z. Liang, L. Li and J. Liu. "Feature selection and classification for Wireless Capsule Endoscopic frames," International Conference on Biomedical and Pharmaceutical Engineering, 2009. ICBPE '09, pp. 1-6, Dec. 2009
- [8] B. Li, M.Q.-H. Meng, "2010 IEEE International Conference on Capsule Endoscopy Images Classification by Color Texture and Support Vector Machine," Automation and Logistics (ICAL), pp. 126-131, Aug. 2010
- [9] C. Poh, T. Htwe, L. Li, W. Shen, J. Liu, J. Lim, K. Chan and P. Tan. "Multi-level Local Feature Classification for Bleeding Detection in

Wireless Capsule Endoscopy Images", Cybernetics and Intelligent Systems (CIS), 2010 IEEE Conference on, pp. 76-81, June 2010

- [10] J. N. Sarvaiya, S. Patnaik and S. Bombaywala, "Image Registration by Template Matching Using Normalized Cross-Correlation", International Conference on Advances in Computing, Control, & Telecommunication Technologies, 2009. ACT '09, pp. 819-822, 28-29 Dec. 2009
- [11] M. Jun, K. Parhi and E. Deprettere, "A Unified Algebraic Transformation Approach for Parallel Recursive and Adaptive Filtering and SVD Algorithms", Signal Processing, IEEE Transactions on, Vol.49, Issue.2, pp. 424-437, Feb. 2001

Modeling Steady and Unsteady High Viscosity Piston Skirts EHL in Initial Engine Start Up

Syed Adnan Qasim¹, M. Shoaib Ansari², and M. Afzaal Malik³

¹College of Electrical and Mechanical Engineering, National University of Sciences and Technology

(NUST), Rawalpindi, Punjab, Pakistan

²School of Mechanical and Manufacturing Engineering, National University of Sciences and Technology (NUST), H-12, Islamabad, Pakistan

(NUST), H-12, Islamabad, Pakista

³Department of Mechanical and Aerospace Engineering, Air University, E-9, Islamabad, Pakistan

Abstract - The initial engine start up wear occurs at low load and speed when an IC engine experiences transient transverse piston eccentricities in the absence of an elastohydrodynamic lubricating (EHL) film. The unsteady squeeze and steady-state wedging alter the piston load-carrying capacity. The 2-D steady-state computational piston EHL models are developed. The unsteady squeeze is introduced to assess its effects on pressures and film profiles. Reynolds equation in steady-state and unsteady-state conditions is solved numerically by generating a finite difference computational mesh to produce hydrodynamic pressures. The graphical simulations are generated for better process visualization. The computational analysis show the effects of wedging and squeeze on piston eccentricities and contact geometry. A comparative simulations analysis shows that squeeze alters hydrodynamic pressures significantly as compared to steady-state conditions. The load-carrying capacity of lubricant is adversely affected, increasing the chances of a physical contact and wear of piston skirts and liner.

Keywords: Modeling, Simulation, Piston skirts, EHL,

Initial engine start up

1 Introduction

The mathematical modeling of the lubrication of piston and rings assembly is a challenging area in the initial internal combustion (IC) engine start up conditions. In the initial engine start up the absence of a fully established EHL film reduces the effectiveness of piston lubrication, invites wear and reduces the engine life considerably. To find ways to prevent the engine wear the hydrodynamic processes of transient squeeze and the steady-state wedging must be coupled with the secondary piston dynamics and then modeled numerically. The unsteady squeeze process involves the timemarching for the given initial and boundary conditions. An engine lubricant flows between the interacting surfaces in relative motion normal to the direction of flow. The unsteady squeeze and the steady-state wedging improve the lubrication and performance of many machines, components and engine bearings [1]. An effective piston skirts EHL requires an appropriate viscosity-grade for easy engine start up for which a high-viscosity engine lubricant is considered in the computational at a low speed. To check the changes in the

squeeze EHL the couple stresses influencing the squeeze performance of a journal bearing may be applied on the piston skirts [2]. To model these processes by solving the set of mathematical equations simultaneously a computer code is needed to be developed, necessitating the use of a professional engineering software. In view of this the 2-D numerical models are developed in matlab in the isothermal conditions. In the steady and the un-steady state conditions separately the 2-D Reynolds equation is discretized and solved numerically to generate the hydrodynamic pressures. In the computational EHL models the pressure-viscosity relationship is introduced and the elastic surface displacements are determined.

1.1 Assumptions

A comparison is made between the simulation results of the steady and un-steady lubrication models for the analysis. For the computational models the assumptions are:

1.1.1 Newtonian oil flooding and thermal effects neglected.

1.1.2 Surface roughness and waviness are neglected.

1.1.3 Zero pressure is assumed at the inlet of the contact zone

1.1.4 The effects of turbulence are neglected in the flow.

1.2 Nomenclature

- C = Radial clearance between piston and bore =10 microns $C_g = \text{Distance of center of mass and piston pin = 0.0002 m}$ $C_p = \text{Gap between gudgeon pin and piston-axis = 0.001 m}$ $E_1, E_2 = \text{Young's moduli of skirts and cylinder = 200 \text{ GPa}}$ F = Force acting as Normal on the skirts $F_f = \text{Force due to friction that acts on the surface of skirts}$ $F_{fh} = \text{Force of Friction due to the hydrodynamic film of oil}$ $F_G = \text{Gas force (function of crank angle)}$ $F_h = \text{Normal force due to hydrodynamic pressure in the film}$ $F_{IC} = \text{Transverse Inertia force due to piston mass}$ $F_{IP} = \text{Transverse Inertia force due to piston-pin mass}$ $F_{IP} = \text{Reciprocating Inertia force due to piston-pin mass}$ $I_{Pis} = \text{Piston inertia about its centre of mass}$
- L = Length of the skirts = 0.0338 m
- M = Moment generated that acts on the skirts
- M_f = Friction moment acting on skirt surface

- R = Half of piston diameter = 0.0415 m
- a = Height from top of skirt and gudgeon-pin = 0.0125 m
- b = Height from skirt top to piston centroid = 0.0015 m.
- e_t , e_b = Piston eccentricities at skirts top and bottom surface
- h = Film Thickness, l = Connecting rod length = 0.133 m
- m_{pis} =Piston mass = 0.295 kg, m_{pin} =Piston-pin mass =0.09 kg
- p = Pressure (hydrodynamic), r = Crank radius = 0.0418 m
- $\boldsymbol{\omega}$ = Constant crankshaft speed, corresponds to engine speed τ = Shear stress, U_1 , U_2 = Skirts and liner Poisson's ratio = 0.3
- η = Dynamic viscosity = 0.1891 Pa.S.
- Φ = Connecting rod angle, ψ = Crank angle
- $\theta = \theta 1 + \theta 2$ = Total piston skirts angle = 75 degree

2 Mathematical Model

The IC engine considered in our model is the same as used by Zhu et al [3]. The primary sliding motion of the piston and its secondary dynamics involving the small eccentricities, velocities and acceleration are incorporated in the model. The piston displacement, velocity and acceleration during its primary motion depend on the crank angle ' ψ '. The primary velocity and acceleration at a given crank speed ω , are [3]:

$$U = \dot{Y} = r\omega \sin\Psi + r\omega B \cos\Psi (l^2 - B^2)^{-0.5}$$
(1)
where $B = C_p + r \sin\Psi$ (2)
$$\dot{Y} = r\omega^2 \cos\Psi + (r\omega B \cos\Psi)^2 (l^2 - B^2)^{-1.5} + ((r\omega\cos\Psi)^2 - r\omega^2 B \sin\Psi) (l^2 - B^2)^{-0.5}$$
(3)

The secondary acceleration terms of the top and bottom skirts are calculated as [4]:

$$\begin{bmatrix} m_{pis} \left(1 - \frac{a}{L}\right) + m_{pis} \left(1 - \frac{b}{L}\right) & m_{pis} \frac{a}{L} + m_{pis} \frac{b}{L} \\ \frac{I_{pis}}{L} + m_{pis} (a - b) \left(1 - \frac{b}{L}\right) & m_{pis} (a - b) \frac{b}{L} - \frac{I_{pis}}{L} \end{bmatrix} \begin{bmatrix} \ddot{e}_t \\ \ddot{e}_b \end{bmatrix} = \begin{bmatrix} F + F_s + F_f \tan \emptyset \\ M + M_s + M_f \end{bmatrix}$$

$$(4)$$

where,
$$F_s = tan \Phi (F_G + F_{IP} + F_{IC})$$
 (5)

$$M_s = F_G C_p - F_{IC} C_g \tag{6}$$

The hydrodynamic film thickness of the lubricant is [4]:

$$h = C + e_t(t) \cos x + [e_b(t) - e_t(t)] \frac{y}{L} \cos x$$
(7)

where, $\theta = x / R$

The 2-D steady-state Reynolds equation is used to calculate hydrodynamic pressure[4,5]:

$$\frac{\partial}{\partial x}\left(h^{2}\frac{\partial p}{\partial x}\right) + \frac{\partial}{\partial y}\left(h^{2}\frac{\partial p}{\partial y}\right) = 6\eta U \tag{8}$$

The 2-D average Reynolds equation involving the squeeze action and the boundary conditions are [3]:

$$\frac{\partial}{\partial x} \left(h^2 \frac{\partial p}{\partial x} \right) + \frac{\partial}{\partial y} \left(h^2 \frac{\partial p}{\partial y} \right) = 6\eta U \frac{\partial h}{\partial y} + \qquad (9)$$

$$\frac{\partial p}{\partial x_{\theta=0}} = \frac{\partial p}{\partial x_{\theta=\pi}} = 0$$

$$p = 0 \quad \text{when } \theta_1 \le \theta \le \theta_2$$

$$p \left(\theta, 0 \right) = p \left(\theta, L \right) = 0 \qquad (10)$$

The dimensionless form of the important variables and the Reynolds equation are [4,5]:

$$h^{*} = h/c; \ x^{*} = x/R; \ y^{*} = y/L; \ p^{*} = (pc^{2})/6U\eta R$$
(11)
$$\frac{\partial}{\partial x^{*}} \left(h^{*3} \frac{\partial p^{*}}{\partial x^{*}}\right) + \left(\frac{R}{L}\right)^{2} \frac{\partial}{\partial y^{*}} \left(h^{*3} \frac{\partial p^{*}}{\partial y^{*}}\right) = \frac{\partial h^{*}}{\partial y^{*}} + \frac{\partial y^{*}}{\partial t^{*}}$$
(12)

The normal hydrodynamic force and its moment about the piston pin are [4]:

$$F_{\rm h} = R \iint_{A} p(x, y) \cos x \, dx dy \tag{13}$$

$$M_h = R \iint_A p(x, y)(a - y) \cos x \, dx \, dy \quad (14)$$

The hydrodynamic friction force and its moment about the piston pin are calculated as [4]:

$$F_{fh} = R \iint_{A} \tau(x, y) dx dy$$
(15)
$$M_{fh} = R \iint_{A} \tau(x, y) (R \cos x - C_p) dx dy$$
(16)

The non-dimensional shear stress is calculated as [5]:

$$\tau = \left(\frac{\upsilon\eta}{c}\right) \left(\frac{1}{h^*} + 3h^* \frac{\partial p^*}{\partial y^*}\right)$$
(17)

In the EHL regime the bulk elastic displacement is taken. The equation representing the thickness of the film is [4]:

$$h_{ehl} = h + f(\theta, y) + v \tag{18}$$

where $f(\theta, y)$ shows the profile of the surface of piston. It represents any deviations from the desired manufacturing tolerances and neglected here. The deformation of the surface in the differential form is [4, 6]

$$dv = \frac{1}{\pi E'} \frac{p(x,y)dydy}{\left(\frac{f}{r}\right)^2 + \left(\frac{f}{r}\right)^2} \tag{19}$$

where
$$\hat{\mathbf{r}} = \sqrt{(x - x_0)^2 + (y - y_0)^2}$$
 (20)

$$\frac{1}{E_{\ell}} = \frac{1}{2} \left[\frac{(1 - v_1^{-1})}{E_1} + \frac{(1 - v_2^{-1})}{E_2} \right]$$
(21)

the elastic displacement at a particular point (x_o, y_o) is [6, 8]:

$$v(x_0, y_0) = \frac{1}{\pi E_i} \iint_A \frac{p(x, y) dx dy}{\acute{\mathbf{r}}}$$
(22)

3 Numerical Solution

Equation (4) constitutes an initial value problem to find the secondary eccentricities $e_t(t)$ and $e_b(t)$ of the skirts [4]. The simulation of the differential equations show the hydrodynamic lubrication of skirts at the respective time steps or crank angles [4]. The average Reynolds equation and the film thickness equation are solved simultaneously. A 21x21 size finite difference mesh is generated and an explicit forward time central space differencing scheme is employed to determine the hydrodynamic pressures, using the iterative computational Gauss Seidel approach [7, 8]. The equations for Reynolds pressures, and the thickness of oil film are solved simultaneously to calculate all forces and moments in equation (4). We compute accelerations \dot{e}_p \dot{e}_b and then satisfy them from the solution of velocities \dot{e}_p \dot{e}_b at the old and new time steps. If we fail to satisfy them then we adjust the present solution of the velocities. For this we use Runge-Kutta scheme. The piston position at the end of the current time step after getting the satisfactory values of the secondary velocities \dot{e}_p \dot{e}_p , \dot{e}_b , is [4]:

$$\mathbf{e}_{t}(t_{i} + \Delta t) = \mathbf{e}_{t}(t_{i}) + \Delta t \dot{\mathbf{e}}_{t}(t_{i}); \ (t_{i} + \Delta t) = \mathbf{e}_{b}(t_{i}) + \Delta t \dot{\mathbf{e}}_{b}(t_{i})$$

The solution converges in the sixth crank cycle. The simulation results of the computational skirts lubrication models correspond to the 720 degree cycle and show the film profiles and pressures. For the numerical solution of the EHL model, the Reynolds equation is integrated for the geometry of the surface interacting with the liner [4, 6]. The inverse solution technique is used to solve the Reynolds equation [6]. The elastic surface displacements in the piezo-viscous regime are calculated and then incorporated in the film thickness equation to produce the modified film profiles. The thickness of the EHL film is calculated by using the three basic numerical processes. The first is the integration of Reynolds equation where as the second is the inverse approach to find a solution of the Reynolds equation [6]. It gives the profile of the film and the elastic deformation calculations. In the solution of the Reynolds equation we locate the inflexion points at two locations on the curve of the pressure due to hydrodynamic action. One point is found on the sweep of pressure at the inlet and the second is close to the outlet. The point of minimum film thickness is found near the outlet end of the pressure curve. The condition of zero second order pressure changes in the sliding direction is applied at a particular point to get the film thickness. At the other points the film thickness is determined. The cubic equation is solved such that [6]:

$$K\overline{h^{3}} + \overline{h} - 1 = 0$$
(23)
where $K = \frac{h_{m}^{2}}{12\eta u} \frac{\partial p}{\partial y}; \quad \overline{h} = \frac{h}{h_{pmax}}$

4 **Results and Discussion**

An understanding of the prevailing operating conditions in the start up of an engine at a low-speed engine must be developed. In view of this, the mathematical models of the lubrication of the piston skirts cater for the steady-state and the transient engine start up conditions. In the steady-state lubrication models the secondary piston displacements are coupled with the steady-state wedging and the resultant hydrodynamic pressure profiles are studied in the stated lubrication regimes. In the un-steady-state lubrication models the time-dependence of the film thickness is coupled with the secondary piston displacement rates. The simulation results are generated that cater for the steady-state hydrodynamic wedging and the time-dependent squeeze action, separately. In each case the graphical results are plotted against the 4-stroke cycle.

4.1 Steady State Piston Lubrication

In the steady-state piston lubrication models the behavior of a high-viscosity grade Newtonian engine lubricant is studied in terms of the profiles of the lubricant film and pressures generated in the hydrodynamic and EHL regimes. The magnitudes of the secondary piston displacements or eccentricities introduce the probability of a contact of the opposing surfaces of the skirts and the cylinder during the 4stroke cycle at 600 rpm speed.

4.1.1 Eccentric Displacements and Velocities

The dimensionless eccentric displacements or eccentricities of the piston skirts in the hydrodynamic and EHL regimes are shown in (1). The eccentricity profiles of the surfaces at the top and bottom are represented as E_t and E_{b_1} , respectively. These profiles are plotted against the 4-stroke crank cycle. The sub-figures 1 (a) and (b) are similar generally and show the three horizontal lines. The upper horizontal line is at 1.0 and represents the non-thrust side of the cylinder liner. The lower line is at -1.0, which shows the thrust side of the liner. The middle line is at 0 and shows the concentric position of the piston with the liner axis. If the E_t or the E_b curve touches either the upper or the lower line, adhesive wear will occur. The induction, compression, expansion and exhaust strokes of the piston are from 0-180, 181-360, 361-540 and 541-720 degrees, respectively. In the hydrodynamic regime, the piston travels concentrically with the liner axis up to the mid-induction stroke and attains the maximum cyclic velocity. Then, it decelerates and gets displaced eccentrically towards the non-thrust side of the liner.At the end of the induction stroke the piston with the negligible cyclic speed changes its direction of translation and commences its journey in the compression stroke. These developments displace the piston further towards the nonthrust side eccentrically. In the compression stroke the piston compresses the air-fuel mixture. At the end of the stroke the highly compressed mixture is ready for combustion, which occurs at 372 degree crank angle. A tremendous amount of energy is transferred to the piston crown, which displaces it eccentrically from the non-thrust to the thrust side of the liner. The thrust generated due to combustion subsides as the piston moves past the mid-expansion stroke. The exhaust valves start opening and become fully open in the exhaust stroke. No significant energy transfer during the exhaust stroke precludes the possibility of any further eccentricities. In the expansion stroke the E_t curve touches the thrust side and remains in contact until the completion of the exhaust stroke. In EHL

regime, the pressures rise very high. The lubricant viscosity rises exponentially with the pressures. The piezoviscous effects and surface deformations reduce the chances of adhesive wear, as shown in 1(b). The eccentricity rate or velocity profiles are generated to study the energy transfer between the skirts and the liner. Figure (2) shows the dimensionless velocity profiles of the top and bottom of the skirts. $E_t dot$, and $E_b dot$ represent the velocity profiles of the top and bottom skirts, respectively. The velocity profiles in the positive and the negative quadrants imply the energy transfer from the skirts to the liner and vice versa. The extraordinarily high elastohydrodynamic pressures enhance the capacity of the oil film to carry the loads effectively.





4.1.2 Film Thickness

Figure (3) shows the film profiles in the hydrodynamic and EHL regimes. In the rigid lubrication regime the film is very thick prior to the side leakage and squeeze out effects and its profile is represented as *Max. Hyd. Film.* After the leakage and squeeze out the left over film carries the actual loads and is the *Min. Hyd. Film.* With no eccentricities the films are quite thin up to mid-point of the induction stroke. The eccentricities in the induction and compression strokes thicken the films. Combustion occurs as the expansion stroke commences. The tremendous combustion thrust does not allow any other effects to influence the primary piston translations in the third (expansion) and the fourth (exhaust) strokes. A noticeable effect of combustion is the sharp decline in the hydrodynamic film thickness. The un-steady effects increase the thickness of the film appreciably. The squeeze effects reduce the thickness of the maximum film significantly. The resultant minimum film thickness is around one-third of the maximum film thickness. However, after the reduction in the thickness the minimum film carries the hydrodynamic loads. The hydrodynamic pressures rise substantially high in the EHL regime and decrease the thickness of the film considerably. The high rising pressures displace the interacting surfaces elastically. The thickness of the hydrodynamic film subjected to extreme hydrodynamic pressures before the surface displacements is plotted and shown as the *h pmax* profile. The displacements of the surfaces create an additional space which is readily occupied by the lubricant due to flooding. It enhances the thickness of the film. The EHL film profile in (3) shows the film after the elastic displacements. The different magnitudes of the film profiles shows the extent of the elastic surface deformation. The film thickness profiles in the hydrodynamic regime are of particular interest during the compression and exhaust strokes. The eccentricities of the skirts enhance the wedging action to cause a corresponding increase in the film thickness. The very low hydrodynamic pressures in the induction stroke increase rapidly higher values in the compression and expansion strokes. The amplitude of the combustion gas force is very high in the beginning of the expansion stroke. It changes the direction of the secondary piston motion which decreases the thickness of the film sharply. The very high pressures decrease the hydrodynamic film from a few to a fraction of a micron and affect the viscosity. The viscosity of the lubricant becomes pressure-dependent and rises exponentially with EHD pressures. The instantaneous spikes are due to the sudden pressure drop and rise in the beginning of the second (compression) and the second part of power stroke.



Fig 4. Steady Hydrodynamic Pressures at (a) 180 deg (b) 270 deg (c) 360 deg (d) 450 deg (e) 540 deg (f) 630 deg

4.1.3 Hydrodynamic and EHD Pressures

In the steady-state conditions the hydrodynamic pressures are generated primarily by the wedging action. In the induction stroke the pressures are of low intensity and the amplitudes are reasonably low. In the compression stroke the pressures start rising as the air-fuel mixture is compressed and the intensity of the gas pressure force increases. The hydrodynamic pressures build up over the surface of the skirts during the 4-stroke cycle. The intensity and magnitude of pressures vary at each angle of the 720-degree crank rotation cycle. The fairly gentle slopes of the low-intensity pressures in the induction stroke transform into steep ones in the compression and the expansion strokes. It happens due to the transformation to high-intensity pressures in the second-half of the compression stroke. Combustion occurs in the beginning of the expansion stroke. It generates a very high intensity thrust, which intensifies the pressures to very high proportions. The rapidly rising very high intensity pressures imply the extended load-carrying capacity of the lubricant in the face of the tremendous thrust produced by combustion. Figure 4 shows the 3-D pressure fields at some of the critical piston positions in the hydrodynamic regime. The sub-figures 4(a), (b), (c), (d), (e) and (f) show the pressure fields at the end of induction stroke, mid-compression stroke, end of compression stroke, mid-expansion stroke, end of expansion stroke and mid-exhaust stroke, respectively. The simulation results show the buildup of the positive pressures over the surface of the skirts. The low-intensity pressures develop in the induction stroke, fairly high pressures with steep slopes in the compression stroke and the excessively high pressures with steeper slopes in the first half of the expansion stroke. The positive pressures remain biased towards the top surface of the skirts till the piston starts decelerating in the second half of the expansion stroke. Since the top surface of the piston skirts has the significant eccentric displacements as compared to the bottom surface, hence in case of starvation the chances of a physical contact cannot be ruled out completely. The bias of the positive pressures shift to the bottom surface in the second half of the expansion stroke. By that time a physical contact gets established between the top surface and the liner. In the



Fig 5. EHL pressure rise under steady-state conditions

EHL regime the pressures rise very high and deform the surfaces. The rising pressures that cause the surfaces to deform elastically are plotted over the surface of the skirts and shown in (5). The maximum rise is noticed near the mid-skirts surface. The pressures rise more than 16% of the peak

hydrodynamic pressures in the rigid hydrodynamic regime. This is the extent of the rise, which is experienced by the interacting surfaces prior to their deformation elastically.

4.2 Unsteady Piston Lubrication

The time-dependent squeeze action occurs when the piston skirts and the liner surfaces come close to each other. The speed at which the piston skirts come closer to the liner surface is the squeeze velocity. The squeeze motion allows the excess lubricant to squeeze out from the sides of the interacting surfaces. The amount of lubricant trapped in between does not permit the interacting surfaces to establish a physical contact with each other. The entire phenomenon is transient in nature and is expected to create a more pronounced effect during the initial engine start up at low speeds and under the low-load conditions. Due to the squeeze effect the thickness of the lubricant film changes with time. Hence, it is required to be recorded at every instant of time to evaluate the effects created by the squeeze action. The squeeze out of the lubricant film.



4.2.1 Unsteady Eccentricities and Velocities

The primary sliding motion of the piston generates significant secondary oscillations in a few initial engine start up cycles. The secondary piston displacements in the transverse direction determine the possibility of a physical contact between the piston skirts and the cylinder liner surfaces. The squeeze effect is anticipated to be more pronounced during the secondary piston motion. When the squeeze action is considered, the secondary piston displacements are simulated and the results plotted in figure 6. In this figure the dimensionless eccentric displacements of the top and the bottom surfaces of the piston skirts are plotted in the hydrodynamic and EHL regimes in the way similar to that adopted in the case of the steady-state conditions. The eccentric displacement profiles in the hydrodynamic regime

show that the squeeze action generates a beneficial effect in preventing a physical contact between the interacting surfaces. The action is more pronounced when the piston is at or near the respective dead centers. In the hydrodynamic regime merely the wedging action does not stop the piston to come closer to the non-thrust side in the compression stroke. More so, the tremendous thrust produced by combustion facilitates a physical contact, if the squeeze action is absent during the expansion stroke. In the EHL regime very high hydrodynamic pressures are generated. The wedging action produce such exceptionally high pressures under the steady-state conditions that do not allow the interacting surfaces to establish a physical contact with each other. However, under the unsteady state conditions the thin lubricant film gets squeezed out when the interacting skirts and the liner surfaces are experiencing already very high loading. Under these conditions the film thickness is expected to get reduced further to still smaller values. A situation might occur when after the squeeze out action a very thin left out film fails to prevent a physical contact of the surfaces effectively. A similar case is the initial engine start up at the relatively low load and speed conditions. In the EHL regime a comparison of the eccentricity profiles in the steady and un-steady conditions shows the visible effects of squeeze in the expansion and exhaust strokes. Squeeze brings the interacting surfaces closer enough to almost touch each other and wear may occur.

4.2.2 Un-Steady Film Thickness

In the un-steady time-dependent conditions the lubricant film should be thick enough to carry the hydrodynamic loads conveniently and effectively under the variable engine operating conditions. Figure 8 shows the unsteady film thickness profiles in the hydrodynamic and the EHL regimes separately. The respective film thickness profiles are plotted against the 720 degree crank rotation cycle. In the hydrodynamic regime the maximum and the minimum film thickness profiles are generally similar to those plotted in the under the steady-state conditions. However, there are differences in terms of the magnitudes and the instantaneous variations in the shape of the respective films. These differences are due to the effects of the transient squeeze action. The maximum and the minimum film profiles in the hydrodynamic regime show that the squeeze action has a visible effect on the thickness in the expansion stroke. During the expansion stroke the relatively steep gradients and the instantaneous spikes at the end of the stroke as compared to that in the steady-state case are due to the squeeze action. Moreover, the reduction in the thickness is extended at the end of the expansion stroke due to the thrust of combustion and the squeeze effect. In the EHL regime the film gets thinner at the maximum pressure and after the elastic surface displacements. There is an instantaneous spike of low intensity during the expansion stroke. Its comparison with the spike in the steady-state case shows a considerable reduction in the magnitude and duration under un-steady squeeze action.







Fig 10. EHL pressure rise under un-steady-state conditions

4.2.3 Un-steady Hydrodynamic and EHD Pressures

The pattern, magnitude, shape and the bias of the hydrodynamic pressures vary considerably when the timedependent squeeze effects are introduced in the hydrodynamic lubrication model. In the un-steady model the 3-D pressure fields are plotted at the same piston positions as is the case in the steady-state model. A comparison of the models show that during the induction and the compression strokes the squeeze effects are not clearly visible. It is despite the fact that there are minor variations in the magnitudes and the intensities of pressures. The squeeze effects are clearly visible on the hydrodynamic pressures after combustion occurs in the beginning of the expansion stroke. Unlike in the steady-state model the hydrodynamic pressures rise very sharply with very high intensities and steeper gradients when the piston reaches at the mid-expansion stroke. Moreover, the bias of the peak pressures also shift closer to the right edge of the top surface of the skirts. The piston starts decelerating after crossing the mid-stroke location of the liner due to which the effects of wedging action start subsiding. Resultantly there is a gradual reduction in the magnitudes of the generated pressures as the piston slides towards the bottom dead-center. In the case of the steady-state model the effects are visible in the form of the relatively gradual build up of pressures with gentle slopes. A very noticeable change is the shift in the bias of pressures towards the bottom surface of the skirts before the piston

completes the expansion stroke. When the squeeze action is introduced in the model the pressure intensities do not subside and instead, the gradients become steeper. The bias of the pressures does not shift towards bottom surface of the skirts. The net result is the sharp build up of the fairly intense positive pressures. With the piston travel in the second half of the expansion stroke and in the exhaust stroke the peak pressures continue to intensify under the squeeze action. The hydrodynamic action is primarily under the influence of the time-based squeeze film formation during the exhaust stroke. It is the squeeze effect, which does not allow a complete contact between the interacting surfaces during the expansion and the exhaust strokes.

5 Conclusions

In this work the effects of time-dependent squeeze action are analyzed in the un-steady EHL of piston skirts in the initial engine start up. The simulation results with and without squeeze effects are compared and analyzed. In the hydrodynamic regime the results show that squeeze produces some beneficial effects. The physical contact between the interacting skirts and the liner surfaces is managed to be avoided in the expansion stroke as compared to the steadystate conditions when wear takes place due to such a contact. In the steady-state conditions the bias of hydrodynamic pressures shift towards the bottom surface in the expansion stroke and become less intense. Under the un-steady conditions the squeeze effects aid in maintaining the highintensity hydrodynamic pressures in the expansion and exhaust strokes. The high-intensity pressures remain biased towards the top surface of the skirts. The squeeze action affects the hydrodynamic film thickness after combustion in the expansion stroke. In the EHL regime the squeeze action is not beneficial as it causes the surfaces to almost contact each other and wear takes place. The EHL pressures do not rise as high under the un-steady conditions as is the case in the steady-state model. The squeeze action alters the EHL film thickness profile due to the variations in the amplitudes of EHD pressures. Before considering any definite conclusions further studies are recommended by considering the different viscosity-grade lubricants and initial engine start up speeds under the stated operating conditions.

6 References

[1] E. A. M. Almas and F. A. P. De Silva, "Finite Difference Automatically generated Non-uniform Grids in the Numerical Solution of the Reynolds Compressible 1-D Squeeze-Film Equation ", Proc. Instn Mech. Engrs Part J: Journal of Eng. Tribology., Vol. 217, pp-244-248, 2004.

[2] Jaw-Ren Lin, "Squeeze Film Characteristics Of Finite Journal Bearings: Couple Stress Fluid Model", Tribology International, Vol. 31(4), pp. 201–07, 1998.

[3] Zhu D., Cheng H. S., Arai T. and K. Hamai, "A Numerical Analysis for Piston Skirt in Mixed lubrication. ASME J. Tribol., Vol. 114(3), pp. 553-62, 1992.

[4] S. Adnan Qasim, M. A. Malik, M. A. Khan, and R. A. Mufti, "Low Viscosity Shear Heating in Piston Skirts EHL in the Low Initial Engine Start Up Speeds", Tribology International, Vol. 44(10), pp.1134-43, Sep 2011.

[5] Gwidon Stachowiak and A. W., Batchelor "Engineering Tribology", 3rd ed., Elsevier, 2005, pp. 112-290.

[6] D. Dowson and G. R., Higginson "Elasto-hydrodynamic lubrication: the fundamentals of gear and roller lubrication", 1st ed. Pergamon Press; pp. 55-106, 1966.

[7] Hoffmann K. A. and Chiang S. T., "Computational Fluid Dynamics", Vol. 1. 4th ed. EES, 2000, pp. 29-96.

[8] S. Adnan Qasim, M. A. Malik and U. F. Chaudhri, "Analyzing Viscoelastic Effects in Piston Skirts EHL at Small Radial Clearances in Initial Engine Start Up", Tribology International, Vol. 45(1), pp.16-29, Jan 2012.

[9] Coy R. C., "Practical Application of Lubrication Models in Engines", Tribology Transactions, Vol. 31(10), pp. 563 -71, 1998.

[10] M. A. Malik, S. Adnan Qasim, B. Rashid and Shahab K., "Modeling and Simulation of EHL of Piston skirts considering Elastic Deformations in Initial engine start up "Proc. STLE/ASME IJTC, Trib2004-64101, pp. 859-67 2004.

[11] S. J. Hupp, "Defining the Role of Elastic Lubricants and Micro Textured Surfaces in Lubricated, Sliding Friction", PhD Thesis, Massachusetts Institute of Technology (MIT), USA, Feb 2008.

[12] I. Kudish, G. Ruben and M. Coritch, "Modeling of Lubricant Degradation and EHL", In Snidlle R. W., Evens H.
P., editors. IUTAM Symposium on EHD and Micro-EHD, 1st ed. The Netherlands, Springer, pp. 149-174, 2006.
Nonlinear Modeling and Numerical Analysis for Fabricating High-performed Flexible Yarns over a Moving Solid Structure

Jie Feng¹, Bingang Xu¹ and Xiaoming Tao¹

¹Institute of Textiles & Clothing, The Hong Kong Polytechnic University, Hong Kong

Abstract: *Flexile yarn structure made from natural fibres,* nano-fibres, carbon nanotubes or other types of fibrous materials is all formed by twisting an assembly of short or long fibres. During the fabrication process, flexible yarn is subject to different processes and contacts with various machine parts, which causes variable tension forces and frictional resistances to torsional stresses of the yarn, leading to different yarn tension and twist distributions. In the previous studies, many attentions have been attracted to analyze yarn twist and tension distributions on a stationery cylindrical solid. However, few of studies were carried out on a moving cylindrical solid structure. Therefore this paper presents a theoretical modelling and numerical analyses for dynamic twist and tension distributions of high-performed flexible yarns over a moving cylindrical solid under steady condition. Nonlinear equations of equilibrium are established for conducting numerical simulations and the influences of various system parameters on yarn tension and twist distributions are discussed in detail.

Keywords: Yarn dynamics, theoretical modelling, twist distribution, tension distribution

1 Introduction

Flexible yarns are the fundamental materials for making a wide range of functional structures and composites such as wearable electronics, smart and conductive structures, geotextiles, 3D auxetics and protectors, and membranes. For the flexible yarn, twist is an important element for providing cohesive forces into fibres to form a high-performed yarn structure with sufficient strength for various end-uses. During fabrication, the yarn is subject to various tension and frictional resistances caused by contacting with machine parts, as simply depicted in Figure 1. After passing through a pair of delivery rollers, twist is generated by a twister and it propagates upwards into the yarn. When the varn contacts with surface of the solid structure, frictional resistance is generated to impede the twist propagation, leading to different twist distributions $(T_2>T_1)$.



Figure.1 Schematic view of twist blocking

In the fabrication and handling process of flexible yarns, most of friction solids structures are stationery with cylindrical cross section, so a number of previous studies have concentrated on distributions of yarn tension and twist on such a friction solid [1-2]. However, few of studies were carried out on the analysis of yarn twist and tension distributions on a moving cylindrical solid. Therefore, this paper aims to provide a theoretical method to quantify yarn dynamic behaviour in terms of yarn tension and twist distributions on a moving cylindrical solid structure. We will firstly establish equations of dynamic equilibrium on the moving cylindrical solid. Then we will numerically simulate yarn tension and twist distributions with steady conditions and discuss the influences of various system parameters on yarn tension and twist distributions.

2 Theoretical Modelling

Figure 2 shows a finite element of flexible yarn running over a moving cylindrical solid. In Figure 2, a cylindrical reference frame *O-XYZ* was selected, in which **R(s, t)** is defined as the yarn element position vector. An arbitrary point *P* on the yarn element is projected onto the Plane *XYO*, as presented by *P'*. Then, the angle formed by the *OP'* with the *X* axis is defined as the wrapping angle (ψ). According to the base vectors $\mathbf{e}_{\psi}\mathbf{e}_{z}$, which represents the yarn tangible direction of increasing ψ and the *z* axis direction, respectively, and $\mathbf{e}_{r} = \mathbf{e}_{\psi} \times \mathbf{e}_{z}$, the yarn element position vector can be expressed by

$$\mathbf{R} = r_0 \mathbf{e}_{\mathbf{r}} + z \mathbf{e}_{\mathbf{z}} \tag{1}$$

In the following analysis, three kinds of external forces exerted on the yarn element are considered, namely internal tension force \mathbf{T} , unit friction force \mathbf{F} and normal reaction force \mathbf{N} .



Figure.2 A yarn element on a moving cylindrical solid

In the steady state, considering all acting forces mentioned above, the full vector form of the time-independent equation of motion [3-5], for a yarn element, is:

$$m[V^{2} \frac{\partial^{2} \mathbf{R}}{\partial s^{2}} + 2\mathbf{w} \times V \frac{\partial \mathbf{R}}{\partial s} + \mathbf{w} \times (\mathbf{w} \times \mathbf{R})]$$
$$= \frac{\partial |\mathbf{P}|}{\partial s} \frac{\partial \mathbf{R}}{\partial s} + |\mathbf{P}| \frac{\partial^{2} \mathbf{R}}{\partial s^{2}} + \mathbf{F} + \mathbf{N} + A \qquad (2)$$

As shown in Figure 3, in order to simply the solution of the differential equation above and denote the direction of all acting forces, a unit vector \mathbf{e}_{β} and a moving coordinate system $Ph\tau v$ with its base vectors $\mathbf{e}_{\mathbf{h}}\mathbf{e}_{\mathbf{r}}\mathbf{e}_{\mathbf{v}}$ were introduced. In Figure 3, π_{1} and π^{2} are the tangent plane to the surface of the moving cylindrical solid and the osculating plane of the yarn curve at *P*, respectively, and \mathbf{e}_{β} is the principle normal to the yarn axis. The unit vectors \mathbf{e}_{τ} and $\mathbf{e}_{\mathbf{h}}$ are the yarn tangent vector in the direction of yarn motion and the unit vector normal to the tangent plane π_{1} , and the vector $\mathbf{e}_{\mathbf{v}}$ can be expressed by $\mathbf{e}_{\mathbf{v}}=\mathbf{e}_{\mathbf{h}} \times \mathbf{e}_{\tau}$.



Figure.3 Moving coordinate and forces acting on a yarn element

In the moving coordinate system, unit normal reaction force N can be written as

$$\mathbf{N} = |\mathbf{N}|\mathbf{e}_{\mathbf{h}} \tag{3}$$

In Figure 3, the angle α is defined as the angle between the direction of friction force and the axis perpendicular to the yarn axis, then the unit friction force **F** can be expressed as

$$\mathbf{F} = -|\mathbf{F}|(\cos\alpha\mathbf{e}_{\mathbf{v}} + \sin\alpha\mathbf{e}_{\mathbf{h}}) \tag{4}$$

In the cylindrical coordinate system, the unit tangent vector \mathbf{e}_{τ} can be written as

$$\mathbf{e}_{\tau} = \frac{d\mathbf{R}}{ds} = r_0 \frac{d\psi}{ds} \mathbf{e}_{\psi} + \frac{dz}{ds} \mathbf{e}_{z}$$
(5)

As shown in Figure 3, the torsional moment for a yarn element can be expressed by the following equation

$$2\pi IDnds = \frac{\partial \mathbf{M}}{\partial s} ds + m_f ds \tag{6}$$

where m_f is the external moment, which can be written by $m_f = |\mathbf{F}| \cos \alpha R$.

In this analysis, a linear relationship between the twist and torque of the yarn (M=KT), together with twist flow equilibrium [6], is assumed. Therefore, in the steady state, Equation (6) can be re-written by

$$(2\pi I V^2 - K) \frac{\partial T_W}{\partial s} = m_f \tag{7}$$

Figure 4 shows the coordinate system and geometrical relationship on a moving cylindrical solid.



Figure.4 Coordinate system and geometrical relationship on a moving cylindrical solid

Since it is assumed that there are no relative rotation is involved between the yarn and the cylindrical solid, we can obtain a general force equilibrium equation

$$\frac{\partial |\mathbf{P}|}{\partial s} \frac{\partial \mathbf{R}}{\partial s} + |\mathbf{P}| \frac{\partial^2 \mathbf{R}}{\partial s^2} + \mathbf{F} + \mathbf{N} = \mathbf{0}$$
(8)

As shown in Figure 4, the vector of the cylindrical solid is given by

$$\mathbf{R} = \{r_0 \cos\psi, r_0 \sin\psi, z\}$$
(9)

where the angle of ψ ranges from 0 to π .

Following the relationship between two base vectors of the reference frame and the moving coordinate system, Equation 8 can be further written as

$$\begin{aligned} \left| \partial |\mathbf{P}| / \partial s - |\mathbf{F}| \sin \alpha &= 0 \\ \left| \mathbf{P} | K_n + |\mathbf{N}| &= 0 \\ \left| \mathbf{P} | K_g + |\mathbf{F}| \cos \alpha &= 0 \end{aligned}$$
(10)

where K_n and K_g are the normal curvature and geodesic curvature, which can be expressed by

$$\begin{cases} K_n = K \cos \gamma = -\frac{1}{r_0} \cos^2 \lambda \\ K_g = K \sin \gamma = -\frac{d\lambda}{ds} \end{cases}$$
(11)

where λ is the angle between the tangent direction of yarn and the tangent vector of ψ .

When yarn moves over the cylindrical solid, yarn motions are composed of three components: circumferential rotation around yarn axis $2\pi Rn\mathbf{e}_v$, yarn forward velocity $V\mathbf{e}_{\tau}$ and the surface speed of the belt $V_B\mathbf{e}_z$, as shown in Figure 5.



Figure.5 Yarn dynamic analysis

Therefore, the angle between the direction of frictional force and the axis perpendicular to the yarn axis α can be obtained by

$$\tan \alpha = \frac{V_B \sin \lambda + V}{V_B \cos \lambda - 2\pi nR}$$
(12)

where *n* is the yarn rotational speed, expressed by $n=VT_{W^0}+n_0$ $VT_{W^0}+n_0$ through the twist flow equation $\partial n / \partial s = V \partial T_W / \partial s$ [6]; and n_0 and T_{W^0} are initial rotational speed and twist level of the yarn.

Since the friction moment acting on the yarn element is $m_j = |\mathbf{F}| \cos \alpha R$, together with $|\mathbf{F}| = \mu |\mathbf{N}|$, the equation of twist distribution on the moving cylindrical solid becomes

$$(2\pi I V^2 - K) \frac{\partial T_W}{\partial s} = \mu |\mathbf{N}| \cos \alpha R \qquad (13)$$

where μ is friction coefficient between the yarn and moving cylindrical solid.

Equation 13 describes the twist distribution in the moving frictional contact zone. Solving equations 10, 11, 12 and 13 simultaneously yields yarn tension \mathbf{P} and twist level \mathbf{T}_{w} .

3 Numerical Simulations

In this section, three cases are used for numerical simulations of yarn twist and tension distributions and the data to be used in the following calculations are given in Table 1. In the case study, friction efficient between the yarn and moving cylindrical solid μ , initial yarn tension P₀ and yarn torsional rigidity K are considered as variables and simulation results were shown in Figures 6, 7 and 8.

Figure 6 shows the distributions of yarn twist and tension at various friction coefficients between the yarn and moving friction solid. In Figure 6(a), it was clear that for all three curves, yarn tension decreases to a minimum value at around 5 degrees, and after this point, gradually increases. Higher friction coefficients can result in larger increase of yarn tension. Also as shown in Figure 6(b), all three curves of yarn twist distributions present an approximately linear relationship with warping angle on

| No. | T _{w0} (turns /cm) | $\begin{array}{c} n_0 \\ (turns \ /s) \end{array}$ | P ₀ (cN) | Ψ (°) | μ | R (mm) | r ₀ (mm) | K (10 ⁻⁶ mNm ²) | V (m/s) | V ₀ (m/s) |
|--------|-----------------------------------|--|------------------------|----------|-----|-----------|------------------------|--|------------|-------------------------|
| | 15 | 186 | 10 | 40 | 0.7 | 0.1 | 3 | 2.0 | 0.12 | 65.7 |
| Case 1 | 15 | 186 | 10 | 40 | 0.8 | 0.1 | 3 | 2.0 | 0.12 | 65.7 |
| | 15 | 186 | 10 | 40 | 0.9 | 0.1 | 3 | 2.0 | 0.12 | 65.7 |
| | 15 | 186 | 8 | 40 | 0.8 | 0.1 | 3 | 2.0 | 0.12 | 65.7 |
| Case 2 | 15 | 186 | 10 | 40 | 0.8 | 0.1 | 3 | 2.0 | 0.12 | 65.7 |
| | 15 | 186 | 12 | 40 | 0.8 | 0.1 | 3 | 2.0 | 0.12 | 65.7 |
| | 15 | 186 | 10 | 40 | 0.8 | 0.1 | 3 | 1.0 | 0.12 | 65.7 |
| Case 3 | 15 | 186 | 10 | 40 | 0.8 | 0.1 | 3 | 2.0 | 0.12 | 65.7 |
| | 15 | 186 | 10 | 40 | 0.8 | 0.1 | 3 | 3.0 | 0.12 | 65.7 |

Table.1 System and machine parameters for numerical simulations

the moving friction solid. Higher friction coefficients can lead to a larger reduction of yarn twist.



Figure.6 Case one: Distributions of yarn twist and tension on a moving friction solid at various friction coefficients



Figure.7 Case two: Distributions of yarn twist and tension on a moving friction solid at various initial tensions

Figure 7 shows the distributions of yarn twist and tension on a moving friction solid at various initial tensions. In Figure 7(a), it was clear that for all three curves, there are no obvious changes before 10 degrees,

after this point, yarn tension increases. Higher initial tensions can result in slightly larger increase of yarn tension. Also as shown in Figure 7(b), all three curves of yarn twist distributions present an approximately linear relationship with warping angle on the moving friction solid. Higher initial tensions can lead to a larger reduction of yarn twist.

Figure 8 shows the distributions of yarn twist and tension on a moving friction solid at various yarn torsional rigidities. In Figure 8(a), it was clear that there are no obvious changes of yarn tension under various yarn torsional rigidities. Also as shown in Figure 8(b), all three curves of yarn twist distributions present an approximately linear relationship with warping angle on the moving friction solid. Lower yarn torsional rigidities can lead to larger reduction of yarn twist.



Figure.8 Case three: Distributions of yarn twist and tension on a moving friction solid at various yarn rigidities

4 Conclusion

vist

A nonlinear theoretical model is developed to analyze the twist and tension distributions in fabrication of highperformed flexible yarns over a moving friction solid under steady state. A series of governing equations have been formulated to describe yarn tension and twist distributions on a moving frictional contact zone. These equations can be solved simultaneously when the calculation parameters are known.

Numerical analysis of simulation results showed that there is a non-linear relationship between yarn tension and warping angle, while yarn twist presents an approximately linear relationship with the changes of warping angle. Higher friction coefficients and initial tension can result in larger reduction of yarn twist with larger increase of yarn tension. Lower yarn torsional rigidities can lead to larger decrease of yarn twist, while yarn torsional rigidity has no obvious influence on yarn tension.

5 Acknowledgment

This research was funded in part through a research grant and postgraduate scholarship by the Hong Kong Polytechnic University.

6 References

[1] Zhang, W. G. "Theories on yarn formation", Textile Industry Press, Beijing, 1983.

[2] Guo, B. P. "Mechanism of yarn twist blockage caused by frictional contact", Ph.D. Thesis, The Hong Kong Polytechnic University, Hong Kong, 2006. [3] "Theoretical mechanics", Tongji University Press, Shanghai, 2005.

[4] Xu, B. G. and Tao, X. M. "Integrated approach to dynamic analysis of yarn twist distribution in rotor spinning: part I: Steady state", Textile Research Journal, Vol.73 (1), 79-89, 2003.

[5] Guo, B. P. Tao X. M. and Lo, T. Y. "A mechanical model of yarn twist blockage in rotor spinning", Textile Research Journal, Vol. 70 (1), 11-17, 2000.

[6] Miao, M. H. and Chen, R. Z. "Yarn twisting dynamics", Textile Research Journal, Vol. 63 (3), 150-158, 1993.

Overall Structural Design of Jet Engine Based on Master Model

Fan Jiang, Zhang Rui, Shen Xiu-li, Wang Rong-qiao, Hu Dian-yin

MDO Group, Beihang University, Beijing, PR China

Abstract - The approach of jet engine's overall structural design based on the master model technique is mainly discussed. The master model of jet engine is divided into control structure and detailed structure according to different design phases. The control structure determines the structural frame of engine and the detailed structure is further designed under the constraints of the control structure that can be driven by the control structure. By means of the control structure method, it solves the problem of combining conceptual phase with detailed phase and facilitates transmitting modifications of design from top to down during overall structural design of jet engine. Finally, the design model is automatically updated to the best result in self-compiled program after optimization process. A parallel and efficient design process of jet engine's overall structure is achieved.

Keywords: CAD/CAE; structural design; master model; control structure; optimization

1 Introduction

The overall design of jet engine is a bottom-up serial design process traditionally ^[1]. The structural design of aero engine mainly depends on designer's experiences and lessons from previous design programs, so designers need to repeatedly check the result with every relevant department to meet general requirements, which will cost lots of manpower, resulting in the prolonged design cycle and increased cost. Thanks to the development of modern CAD/CAE technology, an integrated and collaborative platform has been made available for product design. Especially, the master model technique enabled seamless connection between

design process and data exchanging.

Master model is virtually a central database that contains a lot of geometric and non-geometric information in product's life cycle ^[2]. It ensures the same cal source throughout design process, and coordinates various design phases. Consequently, the master model designed in top-down manner keeps all design levels consistent and relevant. As bottom data come from top-level data, the designed master model can update from top to down when top data changes. In addition, master model can provide for strength, vibration, weight and other disciplines with analysis models and data, which ensures the consistency and relevance of the data analysis between analytical and design departments. Therefore, the process of overall structural design based on master model as a data source can be transformed from serial manner into a parallel process.

2 Control structure

Control structure ^[3] is actually an assembly structure. Control structure method separates conceptual from detailed design of master model, which will form two relevant yet independent structures. This facilitates parallel work for general department and divisional teams. The control structure of master model is designed by the general department, which generally decides the layout of engine parts and defines design criteria of various subsystems. Detailed structure is carried out by divisional teams, which is completed under control structure's constraints. Due to two structures' correlation, when general department changes design, detailed structure will update driven by control structure. Both control structure and detailed structure should follow the principle of top-down design manner, which will relate all levels of master model through control objects.



Figure 1: Control structure and detailed structure of master model

2.1 The system level

The output of flow channel module, cross-section parameters of various components is taken as the input of overall structure module. Therefore, geometric elements and parameters of flow channel design are treated as top-level control objects of overall structural design for the coordination between overall structural and flow channel, which also provides the interface of further flow channel design.

The top level of control structure includes main benchmark and geometric elements, such as the positioning reference of various components, the axial length and import and export parameters of flow channel's cross-section, which controls the layout of overall structure as well as coordinates relative location between components.



Figure 2: The sketch and parameters of top level **2.2 The component level**

For rotor system designing, taking fan for example, the control objects of fan assembly are established under the constraint of flow channel lines and two benchmarks of the fan cross-section. The control objects of fan are mainly to determine rotor series, leave's envelope spaces ^[4] and rank gaps at all ranks. The fan department carries on the detailed design work of fan blades, disks and other parts under the constraint of control objects of components as mentioned above.



Figure 3: The control objects of fan

The control objects of various components are established on the second level according to the different design manner of different components. Control structure reflects the layout and design rules of engine and provides a design benchmark, key parameters and outline for detailed design.

3 The detailed structure

It is component designing departments that design the detailed structure of master model in accordance with control structure following the top-down design method. Detailed structure design is divided into three aspects to introduce briefly. They are bearing-support design, load frame design and rotor system. In this paper, the content focuses on the support structures and rotor system design.

3.1 The bearing-support design

To meet parametric design needs, bearing model should follow these requirements: optional bearing type, numerical bearing location, convenient and reusable bearing model. In this paper, the bearing model is design in light of the UDF^[5] method.

After building bearing template by UDF wizard, the bearing template is stored in user-customized database. Bearings are available from the store to meet designer's requirements. For example, low pressure rotor's supporting style is 0-1-1 and high pressure rotors' is 1-1-0, as shown in figure 4.



Figure 4: The bearing-support style of rotors

3.2 Load frame design

Under the constraint of low-pressure compressor outlet lines and high-pressure compressor inlet lines in the control structure, sketches of outer casing, flow ring and bypass ring are drafted as shown in figure 5. By this correlative design method, the under level objects will contact with the upper objects. So the master model will be capable of updating top-down when the upper objects are modified.



Figure 5: The sketch of outer casing, flow ring and bypass ring

The web of frame contacts the flow ring and bearings in the aspect of its structure, so the flow-ring's base wall is taken as a positional reference of the web's top-wall and the bearing-seat of the web is posited by position parameters of bearing. As shown in Figure 6, bearing-seat inner diameter R1 and R2 are equal to bearing outside diameter D and bearing-seat posited parameters X1 and X2 are equal to bearing posited parameters by the "expression between parts" ^[6] in associated method. What's more, the web structure can

be summarized and divided into two forms, open web and closed web, according to the internal space requirements of current engine types.



Figure 6: The open web sketch **3.3 Rotor system design**

As rotor system is very complicate and involves many aspects, here the main parts discussed are disks and shafts, as well as the connection between parts of the rotor, like the connection of shaft and bearing, the connection of shaft and disk.

3.3.1 Parametric design of disk

Disk features consist of major feature and complementary feature ^[5]. The major feature includes main shape of disk and the complementary feature is designed to complement the major feature by being attached to the main feature, such as installing side, cooling hole, chamfer and so on. According to reference data ^[7], it is straight line and arc method used to describe the shape of web. Leaf envelope space objects in the control structure are used to link the disk sketch and as the design datum of rotor system.



Figure 7: The major feature of disk



Figure 8: The install edge of disk

3.3.2 Relevant design of shaft

Disk installing edge has a surface contacting relationship with the corresponding shaft. The method of parameter referencing is used to build the connection between disk and shaft. As shown in figure 9, HS_X and HS_R are the axial and radial parameters of shaft installing edge, which respectively reference the axial parameter N_2 _X and radial parameter N_2 _R of disk installing edge so as to achieve the relevant design of disk and shaft. In order to facilitate parameter referencing, parameters of disk installing edges in every stage are stored in a parameter table, including N_1 _X, N_1 _R, N_2 _X, N_2 _R and so on. Through referencing different parameters in table, shaft installing edge can contact with any disk in all stages, as shown in figure 10 and figure 11.



Figure 9: Connection between shaft and disk



Figure 10: Shaft contact with the first disk



Figure 11: Shaft contact with the second disk

On the other hand, shaft and the corresponding bearing and have a concentric constraint. As shown in figure 12, shaft outside diameter R is equal to bearing inner diameter D_IN by the "expression between parts" function. The front part of shaft is under the control of bearing. So the shaft will also be updated with the bearing.



Figure 12: Relevant design of shaft and bearing

4 The optimization

Turbine disk is the most demanding and difficult part to design because of the harshest working conditions. Therefore, in this chapter the optimization process based on the master model is described by taking turbine disk for example. As the structural strength calculation needs simplified finite element models, the master model can provide a required context models for the analysis of turbine disk. The process of calculation and updating to the best results is carried on in the optimization platform automatically.

4.1 The context model

Context model ^[8] is an analysis model attached to the master model in a discipline. As the data is from the master model, when the master model released by design departments is modified, the context model analyzed by analytical departments will be updated synchronously, as shown in figure 13 and figure 14. This facilitates transmitting modification through different departments and ensures the possibility of parallel design. Using "create a new hyperlink part" function ^[6] can rapidly create an analysis model of a discipline, without changing original features of the master model. This saves large computing resources and computing time under the premise of accuracy, due to abandoning unnecessary information of the analysis model.



Figure 13: Master model and two context models



Figure 14: Updating after modified **4.2 The optimization process**

1) The program of extracting parameters

Optimization platform needs a parameter file as the input file. The parameter extracting program getpara.dat reads expressions from the master model and outputs them to the text file exp.ini, which is as the parameter input file.

2) The program of updating

The updating program ug.dat needs three files as input files, respectively parameter file, master model file and context model file. Ug.dat reads parameter values in exp.ini and updates the master model disk.prt, then loads and updates context model disk_cm.prt in order. Finally, it gives an output of analysis model disk_cm.x_t from disk_cm.prt for the calculation program. 3) The program of calculating

The program calculate.dat reads calculation command file input.txt and analysis model file disk_cm.x_t automatically, then outputs result file cal_result.txt at backstage.



Figure 15: The process of optimization

4.3 The result of optimization

The entire optimization runs 30 minutes, 85 iterative steps totally. The best result is the 69th step, as shown in Figure 16.



Figure 16: The target curve

As shown in Table 1, the weight of optimized disk reduces by 17.72% compared with the initial figure,

which makes a great improvement. Although the circumferential stress figure increases by 3.2% and the radial stress increases by 10.45%, the result is less than the strength criterion and meets design requirements. During the optimization process, the turbine disk in the master model updates to the best solution automatically by the updating program, as shown in figure 17.



Figure 17: The disk before and after optimized

| Parameter | Remark | Initial Value (mm, mpa) | Optimized Value (mm, mpa) | Contrast |
|-----------|------------------------|----------------------------|------------------------------|-----------------|
| HRIM1 | Rim height 1 | 8.000 | 7.000 | |
| HRIM2 | HRIM2 Rim height 2 | | 9.973 | |
| WWEB1 | Web top width | 22.000 | 15.000 | |
| WWEB2 | Web base width | 28.000 | 25.000 | |
| HHUB | Hub height | 20.000 | 12.000 | |
| WHUB | Hub width | 84.000 | 70.896 | |
| SX | radial stress | 711.509 | 785.867 | Increase 10.45% |
| SZ | circumferential stress | 961.499 | 992.302 | Increase 3.2% |
| V | volume | 6324251.4000 | 5203344.590 | Reduce 17.72% |

Table 1: The design parameters

5 Conclusions

The control structure method more conforms to the characteristics of uncertainty of structural design. It facilitates modification during the whole design process. What's more, the master model can provide a parametric design model and required analysis models, which makes overall structural optimization possible. It provides an efficient and parallel design method that will reduce the design cycle and cost in result.

6 References

 Aero Engine Design Manual editor committee. "Aero engine design manual". Aviation Industry Press, Vol. 5, 2001.

[2] Yuan, Q.K. &Zhang, D.L. "The research of the product master model". China Mechanical Engineering, Vol. 10, Issue 1, 1999.

[3] Niu, G. "The use of UG/WAVE in aero engine design". The 14th structural strength and vibration symposium of China Aviation Society, 206-211, 2008.

[4] Zhou, Q.Y. "Study on master modeling and application in MDO of turbine disk and blade". Machinery Design and Manufacture, Issue 2, 2011.

[5] Zhu, K., & Wang, R.Q. "Analysis and realization of UDF in design of the blade". Machinery Design and Manufacture, Issue 1, 30-32, Feb 2003.

[6] Unigraphics Solutions Inc. "UG/WAVE student guide". Qinghua University Press, 2002.

[7] Xiao, R., & Ma, M., & Lin, X.R. "Aero engine structure optimization". Beihang University Press, 1991.

[8] Li, H.J., & Deng, J.T. "Use of context model in collaborative design". China Mechanical Engineering, Vol. 14, Issue 22, 1947-1950, Nov 2003.

SESSION DESIGN, ANALYSIS, EVALUATION METHODS

Chair(s)

TBA

Computational Architectural Neuroscience: Towards the Computational Theory of the Human Brain Interactions with Architectural Design

H. Mirkia¹, A. Sangari², M. Nelson¹, M Sarmadi¹ and A. Assadi³

¹Design Studies Department, University of Wisconsin, Madison, WI, USA ²Electrical and Computer Engineering Department, University of Wisconsin, Madison, WI, USA ³Mathematics Department, University of Wisconsin, Madison, WI, USA

Abstract - This article illustrates the opportunities and the challenges in an emerging multidisciplinary research area in study of human brain interactions with architectural design: Computational Architectural Neuroscience. We propose to apply main-stream methods from computational neuroscience to transform the commonly static formulation of the interior design evaluation problem into a dynamical systems formulation, which is the hallmark of modeling brain response and human behavior. While navigating indoors, observers experience perceptual features that often represent a synthesis of their visual perception of natural scenes, which are shaped and guided by the design and creativity of the architects. We propose to analyze the dynamic patterns of brain-generated signals in relation to perception of harmony of a design and decision making by human observers. We argue the feasibility of the methodology and efficiency of proposed algorithms to provide a quantitative theory of interior design evaluation in the cases considered below.

Keywords: Interior design, Computational neuroscience, Optical flow, Visual perception

1 Introduction

A challenging problem of theoretical and practical significant in design studies is whether it is possible to develop quantitative measures for evaluating alternative design solutions on a single project. Is it possible to develop quantitative measures for objective classification of separate design projects that might be considered comparable? Aesthetics and human response to aesthetics in architecture are sometimes described in relation to mathematics and geometry. However, the strength of this relationship is difficult to test empirically. If evidence could be shown that this relationship holds true, then computational methods using mathematics to analyze spatial geometry could be used to predict the human response to aesthetics as they exist in a wide range of design solutions. Furthermore architecture, especially interior space, is perceived dynamically during one's movement through space. At the same time study has shown that visual perception is shaped by eye movement that creates a dynamic series of spatial vignettes. While this dynamic series might seem to be difficult to express as a two dimensional static construct, artists have successfully tackled similar problems. This paper proposes a computational model that would predict human response and aesthetics through empirical analysis of the geometry of dynamic two dimensional images of a design that are linked with neuroimaging of the brain occurring while looking at those images.

The content of this article are as follows. We outline the basic background from computational neuroscience to investigate how the brain is likely to respond to harmony of design in interior spaces. The first observation, which we introduce as a way to study a design is the role of "perceptional dynamics" in the human-centric evaluation and appreciation of design elements. This is accomplished by elucidating the manner that eye movement provides the sensory input in to the observer's visual cortex. Eventually, such information is transmitted to higher cognitive areas, and interacts with the affective centers of the brains such as the amygdala. This can be measured by using eye tracking equipment, thus, the availability of a rich source of empirical data from highly relevant brain activities. The second observation concerns a theory for visual perception of perspective in interior spaces. The brain area used to interpret perspective can be measured by using neurological imaging. Accordingly, the area of the brain (MSTd) that is known to contain the neurons and local circuits for perception of optical flow appears to be the site of local circuits that make the most significant contributions to detect perspectivity in interior spaces. By bringing these two observations together in a computational model, we simulate a simplified version of the brain processes for perception of interior space as a dynamic event. Thus, we demonstrate the mathematical feasibility for a quantitative theory of design in which some of the fundamental design principles found in high quality human centric design could be subject to systematic measurement in the context of whole designs. Further, we discuss an approach for practical implementation of algorithms to quantify certain aspects of harmony and facility in communication of design elements.

The contemporary architect implicitly considers the psychological factors in design of interior spaces to invoke the sense of harmony and positive affect in the observers. We outline a unified framework for quantifying the perception of geometry and motion in interior spaces, situated in the psychological foundation of integrative affect and cognition, as a way to evaluate anthropocentric interior design. A systematic mathematical-behavioral framework to model design of interior spaces is Perceptual Geometry [2], [4]. This interdisciplinary research area focuses on study of geometry from the perspective of visual perception, and in turn it applies such geometric findings to the ecological study of The aesthetic aspects of interior design beyond vision. geometry require a deeper study into the affective dimensions of psychology than what is touched upon below. Perceptual geometry is used to investigate quantitative and cognitive aspects of a design with the goal of attempting to answer fundamental questions about the synthesis of perception of form and representation of space. This is paired with biological theories of visual perception and geometric theories of the physical world. The contribution of this paper is to start with the view that perception of form, space and motion are facets of the same computational steps in the observer's brain, and then to argue a hypothesis that mathematical modeling and machine learning can replicate many of the same computational steps as the brain. Our approach incorporates mathematical translation of the following indicators of harmony in interior design: (1) The observer's attention is guided through a configuration of lines; (2) the design elements are communicated easily to the observer to facilitate remembering the information. We propose adding a third indicator, which should be quantified in our model using brain imaging and psychophysics; namely, (3) a significant part of the affective response of the observer arises from the observer's attention to salient features of the design.

3 **Methods and Discussion**

In this section, we sketch the steps for development of a mathematical model for "quantifying design," that is, to assign numerical measures that evaluate the above-mentioned three principles indicators of harmony in design. This hybrid of empirical mathematics, computational neuroscience and perceptual psychology proceeds to quantify design in several steps, as outlined in Figure 2 below. To begin with, we need to validate the assertion that the line-drawing abstraction of a design project faithfully conveys the simplest visualcognitive representation of that project. This step uses experiments with tracking of observer's eye movements and analysis of such data. The (dynamic) data sets are collected in two parallel cases of visual perception: First, perception of the line-drawing as a candidate for simplified medium for communication of design in the project without loss of perceptually critical information; second, eye movement

The Central Theme: Quantifying Design tracking in visual perception of the original design project. The output of this step is selection of a collection of "optimally simplified line drawings" of the project that we refer to as the "design perceptual minimal models." The implementation of this step requires a sufficiently large database of design projects, which are processed by algorithms for extraction of line drawings. These linedrawings must be incrementally "simplified" without loss of critical information about the design, through repetitive psychophysics experiments that uses evaluation of eye movement records as the fundamental tool for comparing two line drawings, and deciding which is "simpler without loss of information," or, having reached the choice between a minimal model and the line-drawing that is over-simplified by removal of a single line from the minimal model. The mathematical result of these operations is an algorithm for "model reduction," which provides the basic elements for the Empirical-Mathematical Theory for Quantifying Design. Due to space limitation, we omit the discussion for an efficient strategy of using functional brain imaging (or an equivalent method) to test and validate the outcomes of the mathematical model/theory.

> The workflow of our model (Figure 1 below) yields outputs that are observer-dependent. Moreover, there could be several "simplest models at the psychophysical threshold," thus several minimal models for the same design project. The design expert then selects one or two minimal models as the "best representatives of design elements."



Figure 1: The workflow of the model for quantifying design.

We measure the brain activation level of the observer by fMRI (or other noninvasive brain imaging). We associate a cost function from the fMRI experiment to the line-drawing that is "perceptually minimal." The lowest value of the cost functions is the quantitative cognitive/perceptual measure for the design. There is a minimal exposure time (threshold time) that the observer's brain needs to process the sensory stimuli and for the perceptual subtasks to take place, in order to reliably distinguish the design element in the project; and a lesser exposure time could confound the observer with a large statistical error. We hypothesize that the observer's brain must extract a minimum level of information in order to achieve this. The perceptual state of the brain, we hypothesize to have constructed an entity, which we refer to as design perceptual minimal model. The dynamic processes for

2

extracting the information for design perceptual minimal model are commensurable to the dynamic processes that the brain generates to reconstruct the Gestalt of design elements (Just Noticeable Difference in psychophysics).

We also propose the existence of neuroanatomical loci that play a pivotal role in perception of perspective, and demonstrate a mechanism for visual perception of interior spaces based on eye movement. We conclude that learning to estimate motion from optical flow is similar to learning to estimate a single vanishing point. [1]

The study of 3D perception from 2D images is the main issue of vision science in both biological and computer systems. Depending on the projection method used for mapping 3D space to a 2D plane, different images may be resulted. Among these projection methods, perspective projection is known as an accurate mapping for reconstructing geometrical features of real objects. One important concept in perspective projection is the "vanishing point". The vanishing point of a straight line under perspective projection is that point in the image beyond which the projection of the straight line cannot extend. Therefore, the vanishing points correspond to the "points at infinity" in the receding direction of the parallel straight lines in space.

Computationally, the vanishing point in interior spaces with a single vanishing point can be found by finding the point where most of the salient straight lines in the image intersect with each other. The salient straight lines can be found using processes such as edge detection and line extraction. The mechanism of edge detection in the brain is related to the lateral inhibition in the receptive field of the visual sensory system [8]. A sample of how edge detection occurs in an interior space is shown in Figure 3, center.

The Gestalt principle of continuity shows how the brain can fill in the gaps between the dots in the image and extract the major lines. The Hough transform can be used to extract the slope and the constant term of lines in the image, based on black and white result of edge detection algorithms [3]. The line drawings made from the interior spaces of Figure 2, center carry pieces of information that the brain needs to estimate the orientation of surfaces and compute the perception of the interior space in perspective. By extending these lines and finding the intersection points, we can computationally find the vanishing point (Figure 2, right).

In real word experience, additional information such as texture, color and cues provide much more useful information for richer perceptual experience of interior spaces [2], [4]. However, this method focuses on simpler geometric objects as illustrated in the line drawings below.



Figure 2: Two steps in detection of vanishing point: edge

detection of local circuits in area V1 and the thalamic feedback system provide the information about the lines as above. We have simulated the extraction of lines using Hough transform. (a) Left: an image from the Wisconsin Institute for Discovery (WID). (b) Center: binary output of edge detection.

(c) Right: overlay of the Hough transform output onto the original image.

4 Algorithms and Computations

The conventional static view in decision making regarding interior design alternatives focuses on the interactions among design elements and the established principles for harmony and esthetics. To transform this static formulation into a dynamic framework, we consider the human brain response to navigation indoors while interacting visually with design elements and eventually arriving at decision making among alternatives. A key step in computational neuroscience modeling of observer's decision making is to elucidate the neuronal mechanisms that incorporate the observer's prior experience of perception of navigation indoors in order to form top-down flow of information (in the Bayesian sense). Optical flow is a primary computational theory for observer's mechanism perceptual governing navigation. Our computational model for an observer demonstrates that prior experience from optical flow considerably improves the observer's performance in identifying perspective geometry, compared to the performance of the "naïve observer model" that must learn the geometry of lines ab intio. We would expect that the measurement of the observer's brain activity when presented an interior design project, and the record of the eye movements (see the work flow) would be an unbiased representation of the correlations among perception of perspectivity (the primary factor in interior space design) and navigation indoors. Therefore, comprehensive data collection as discussed above is very likely to lead the anticipated results for quantification of dynamic patterns of brain activity in response to interactions of the observer and design elements. Finally, we invoke the hypothesis that the facility in reaching the Gestalt of a design project (shortest time to form stable perception, while performing better or within the designated error rate) indicates the observer's decision regarding which design alternative could be the most harmonious. Below, we outline the computational model for identifying perspective geometry in the context of optical flow. which simulates experiments in quantifying performance of observer's perception of interior space.

We have developed a unified framework for perception of geometry and motion in indoor environments. We designed some 3 dimension interior spaces and import them in a virtual reality environment. We use the principal component analysis for reducing dimension of successive frames of this simulation and take some of the most informative components in training a neural network (Optiflonet) for estimation of direction of movement in the image plane. Then, we show this neural network with slightly different configuration can extract the vanishing point in still images.

In our model, we have simulated the sensory information by architectural rendering of the some simple designs (such as Fig. 3). In order to test motions with more complex kinetics, we used virtual reality techniques in rendering these simple interior spaces. We compare extraction of lines from simulated scenes versus extraction of lines in natural scenes.



Figure 3: This is an example for training validating and testing a neural network that learns to extract the vanishing point as a significant piece of geometric information for extraction of orientation of surfaces from directions of lines (edges) in the sense.

We now discuss learning perception of self-movement characteristics from optical flow. When an observer moves in an interior space, optical flow provides a robust mechanism to estimate the motion and heading direction. The optical flow vectors indicate the distribution of stimulated local circuits that detect motion heading direction. Optical flow vectors, (the arrows in the figure 4) can provide the input for training our neural network *Optiflonet*. Validation and testing of neural network are performed using 600*300 arrays with sparse sampling of the optical flow vectors. As it is illustrated in figure 4, for simple interior spaces, optical flow vectors can be approximated by a spars matrix. Therefore, we use principal component analysis for extracting the most informative components of optical flow vectors in all frames.



Figure 4: Optical flow vectors in simulated interior space movement. Right image shows the area inside rectangle with magnification in order to emphasize the small optical flow vectors.

We use the first 60 principal components of image frames in order to reduce the required number neurons in our neural network and in the meanwhile avoid the reconstruction errors. Figure 5 shows the whole structure of *Optiflonet* system for estimation of motion direction in the image plane, also the architecture of neural network that we used is presented in figure 6.



Figure 5: Flowchart of data processing in *Optiflonet* from image frames to estimation of motion direction in image plane.



Figure 6: This figure shows the architecture of *Optiflonet* neural network. It has two layers of neurons and in the hidden layer there exits 100 neurons.

To test our theory we designed a neural network that learns to extract vanishing points from perspectives called *Perspectinet*. This network incorporates a simplified model of eye movement, which passes sampling of the scene as a flow of consecutive snapshot of the interior spaces from an eye with movement. Consequently the subsequent layer of *Perspectinet* receives dynamic information from the scene. Design of *Optiflonet* is based on the figure 7. The optical flow from movement on the retina provides the input to the next layers of visual sensation and visual perception of motion.

We propose a model that an observer brain estimates perspective from prior experience of navigation in indoor space in the context of optical flow. We train an Optiflonet using examples of navigation in an interior space. The model uses training sets that are temporal sequences of image samples, where each image frame is a snapshot of the interior space in perspectives projection. The observer's movement indoors provides different projections on the retina; consequently, different perspectives are registered with possible shifts in the vanishing point. The perspective could be tested in the following way: The dynamic scene (perspective images) turn into a movie that represents the effect of eye movement. Subsequent blocks have layers of the network and connection-weights that are inherited from training the Optiflonet. Our problem is to introduce and measure the performance of the Perspectinet when test images are from interior space without any training by feedback from this category of interior space forms. In other words, the local circuits of Perspectinet in the visual motion processing layers are used as local neural circuits that are

employed for computations in image sequences of eye movement snapshots from interior spaces forms.



Figure 7: Flowchart of data processing in *Perspectinet* from the still image to extraction of vanishing point in the image.

5 Results

In order to train *Optiflonet* we used back propagation method on our training data set. Training data set contained the 150 frames from a simulation of walking through a corridor like that in figure 3. We used another set of 150 frames for validating and testing of our neural network.

Also, figure 8 shows the error in estimation of motion heading in horizontal and vertical axis of each frame separately. Since in training and test data we use a simulated walking, the kinetic of movement (velocity and orientation of camera) through interior space was stochastic variable.



Figure 8: Histogram of output error during the testing *Optiflonet* neural network. Left: error in estimation of movement heading (in degree) along vertical axis of frame. Right: error in estimation of movement heading (in degree) along horizontal axis of frame.

In figure 9, we presented the time series of camera velocity and orientation during recording frame data for training. The orientation of the camera is shown by 3 component of the unit vector parallel to camera orientation.



Figure 9: Time series of camera velocity (right plots) and orientation (left plots) during the recording data for neural network training.

Finally, the histogram of error in detection of vanishing point location in different frames is illustrated in the two histograms of figure 10. The left histogram shows the vertical error while the right one shows the error along horizontal axis of frame. The maximum error in horizontal axis is generally less than 1 degree and in the vertical axis is always less than 10 degree.



Figure 10: Histograms of output error during the testing of *Perspectinet* neural network. Left: error (in degree) in estimation of horizontal component of vanishing point along in the image. Right: error (in degree) in estimation of vertical component of vanishing point along in the image.

6 Conclusion

While the analytically based methodology described here has additional nuances and detailed applications beyond the scope of this preliminary paper, we believe that this methodology can be applied to a large sample of images of interior architectural space in order to identify potential standards for quantifying design elements, the central problem of design studies. Furthermore, meaningful computational analysis can be applied as a way to predict the human emotional response to design of interior space.

While the central problem of quantifying design remains to be studied in more detail, we believe this article's computational methodology provides optimism for rapid progress in near future to overcome the theoretical and experimental challenges that Computational Architectural Neuroscience must overcome. In particular, we predict this emerging multi-disciplinary research direction will formulate and develop the computational framework to quantify significance of high quality principled design in human health and well-being.

7 References

[1] Sangari, A. Mirkia, H and Assadi, A. "Perception of Motion and Architectural Form: Computational Relationships between Optical Flow and Perspective", in Proc. Third Doctoral Conference on Computing, Electrical and Industrial Systems (DoCEIS'12) (2012).

[2] Assadi, A.H., Eghbalnia, H., Palmer, S.: A Learning Theoretic Approach to Perceptual Geometry in Natural Scenes, Neurocomputing, vol.38-40; June 2001, p.p. 1077-1085 (2001). [3] Assadi, A.H., Palmer, S, Eghbalnia, H.: Learning Gestalt of Surfaces in Natural Scenes, Neural Networks for Signal Processing IX Proceedings of the 1999 IEEE Signal Processing Society Workshop Cat No98TH8468 (1999).

[4] Assadi, A.H., Palmer, S, Eghbalnia, H., J. Carew: Geometry of the Perceptual Space, In Proceedings of The International Society for Optical Engineering (SPIE) -Vision Geometry VIII, Latecki, Longin J.; Melter, Robert A.; Mount, David M.; Wu, Angela Y. (editors), Vol. 3811, p.p. 130-140 (1999).

[5] Duchamp, M. Quoted in Marcel Duchamp: catalogue of the Duchamp Retrospective at the Palazzo Grassi, Venice, Gruppo Editoriale Fabri, Milan, (1993).

[6] Martinez-Conde S, Macknik SL & Hubel DH: The role of fixational eye movements in visual perception. Nature Reviews Neuroscience; 5, 229-240 (2004).

[7] Newsome, WT.: Deciding about motion: linking perception to action. Journal of Comparative Physiology, 181:5-12 (1997).

[8] Palmer, S.E.: Vision Science: Photons to Phenomenology. MIT Press. (1999).

[9] Rickey, G.W.: The morphology of movement: a study of kinetic art. Art Journal, (1963).

[10] Sangari, A. Mirkia, H and Assadi, A. "Perception of Motion and Architectural Form: Computational Relationships between Optical Flow and Perspective", in Proc. Third Doctoral Conference on Computing, Electrical and Industrial Systems (DoCEIS'12), (2012).

[11] Yarbus, A. L.: Eye Movements and Vision. New York: Plenum Press, (1967).

[12] Zeki, S.: Inner Vision; an Exploration of Art and the Brain. Oxford university press, (1999).

Validating Statics of Long Term Evolution Mobile Communication Systems

Ishtiaq Ahmed Choudhry¹, Nazir Ahmad Zafar² and Mohammed Al-Zahrani¹

¹ Department of Networking and Communications, King Faisal University, Saudi Arabia Emails: {ichoudhry, malzahrani}@kfu.edu.sa
² Department of Computer Science, King Faisal University, Saudi Arabia Email: nazafar@kfu.edu.sa

an. nazarai @Kru.euu.sa

Abstract - Long Term Evolution (LTE) is the latest wireless standard for International Mobile Telecommunication (IMT) system. There are many issues that are yet to be improved due to dynamic complex nature of wireless systems, multimedia software applications and software requirements. Poor service quality, service disconnections due to mobility, seamless handover, handover interruption time and downward compatibility to other Radio Access Networks (RAN) are some of the key issues that are addressed very recently in the scientific literature. Formal method is one of the promising software engineering techniques that assure quality and perfection in software system models. Formal systems are proved before implementation of the models. Formal methods use mathematical language to explicitly specify system specifications and requirements that serve as initial grounds for further development and implementation. It efficiently handles all component connections and resource management parameters using discrete structures. Z Schema language is used to model static aspects of LTE communication system. All the schemas are being verified using Z/Eves toolset. The aim is to provide sound mathematical foundation for system validation and verification that eventually results in a more reliable, scalable and complete software system.

Keywords: LTE Communication Systems, Formal Methods, Emerging Software, Mathematical Language, Z Notation

1 Introduction

Long Term Evolution (LTE) is the latest wireless standard for International Mobile Telecommunication (IMT) system. LTE has been deployed in many countries and many are in pipeline to follow suit. The IMT system is very complex due to different Radio Access Networks (RAN), and their protocols. IMT systems allow users to access information and services at anytime and anywhere. These systems provide mobility environment with different set of constraints, requirements and resources. Poor service quality, service disconnections due to mobility, seamless handover, handover interruption time and downward compatibility to other Radio Access Networks are some of the key issues that are addressed very recently in the scientific literature [1], [2], [3]. There are many issues that are yet to be improved due to dynamic complex nature of wireless systems, multimedia software applications and software requirements.

LTE provides complex communication architecture and offer many real time multimedia applications. Although there is a substantial research progress in LTE systems but due to complexity and flexibility in the system still there is an uncertainty and sufficient improvement are required in it. The telecom industry, with the support and cooperation of academics, is still discovering the system itself, network approaches, and other optimization techniques. Consequently due to complexity of LTE systems planning, non-existence of scientific theory and solutions it is an open problem.

Formal methods are mathematical techniques which are used to model complex systems. By building a mathematics-based model of a system, it is possible to verify its properties using formal tools in a more rigorous and thorough fashion than other traditional testing and simulation techniques. By rigorous descriptions of a system it is promised to improve its reliability and design comprehensively. However, the abstract models described by formal methods have their own limitations but it is observed that use of formal methods is increased in safety, security and complex systems.

Formal methods which are based on discrete mathematical structures enable us to write system specification in such a way that it helps us understanding states of the system and its properties in a rigorous and clear fashion. Formal methods are important abstraction technique for managing the complexity of large, dynamic and distributed systems. The systems requirements, design, specification and limitations are addressed at early stages of systems development. In this way, a thorough analysis, design and development processes leads to a complete, consistent robust model. The most important benefit of using formal methods is that we can gradually move from informal requirements to semi-formal which then further can be transformed to the detailed model maintaining the system properties and conditions.

In the most relevant existing work [4], Milner extended the system from Calculus of Communication System to piecalculus enabling a dynamic communication model which can be used to describe the mobility. In this work, mobility is assumed as movement of links between the connected components. Acharya et. al, [5] have used formal techniques to overcome the consistency issue at various phases of development process in the domain of wireless mobile networks. Formal models are presented in [6], [7] using Z and object oriented Z, however, few discrepancies are identified and the systems are not addressed completely. Duke et. al, [8] have presented several case studies to illustrate application of formal methods. An interesting application is a study regarding mobile phone system which keeps track of several states of mobiles with limited number of operations. The same case study is extended by Battaz in [9] to address the mobility issues explicitly using Z notation. For the same system, Battaz [9] has described handover procedure among two objects in The requirements analysis of LTE more details. communication system is presented in [10]. LTE is developed by the 3rd Generation Partnership Group (3GPP) [11] and its specifications laid out by 3GPP are mainly focused in this paper for model abstraction.

In this paper, formal specification of static structures of the Long Term Evaluation (LTE) system is presented using Z notation. First of all, fundamental definitions used in the LTE model are presented then specification of messages is described. Formal description of user equipment is given based on the above definitions. In the next, formal description of core components, namely, evolved NodeB, Home Subscriber Server and Serving Gateway is presented. Finally, Mobility Management Entity is described after composition of the above structures. Formal specification of the system is analyzed using Z/Eves toolset.

Rest of the paper is organized as follows: Section 2 describes the basic model of the LTE communication system. In section 3, formal model using Z notation is provided. Formal analysis is done in section 4. Finally, conclusion and future are addressed in section 5.

2 Long Term Evaluation System

An introduction to basic components of LTE is provided in this section [12]. The communication system architecture consists of the following functional elements. After analysis of the system and its functional components formal description of the system will be provided.

2.1 Evolved Radio Access Network

The evolved RAN for LTE consists of a single node, i.e., the eNodeB (eNB) that interfaces with the user equipment (UE). The eNB hosts the PHYsical (PHY), Medium Access Control (MAC), Radio Link Control (RLC), and Packet Data Control Protocol (PDCP) layers that include the functionality of userplane header-compression and encryption. It also offers Radio Resource Control (RRC) functionality corresponding to the control plane. It performs many functions including Radio Resource Management (RRM), Admission Control (AC), scheduling, enforcement of negotiated Upload Link (UL), Quality of Service (QoS), Cell Information Broadcast (CIB), ciphering and deciphering of user and control plane data, and compression and decompression of Download Link (DL) user plane packet headers.

2.2 Serving Gateway

The SGW routes and forwards user data packets, while also acting as the mobility anchor for the user plane during intereNB handovers and as the anchor for mobility between LTE and other 3GPP technologies. Its functionality is also terminating S4 interface and relaying the traffic between 2G/3G systems and packet data networks gateway. For idle state of user equipment, the survey gateway terminates the download link data path and triggers paging when the data arrives for the user equipment. It manages and stores UE contexts, e.g. parameters of the IP bearer service, network internal routing information. It also performs replication of the user traffic in case of lawful interception.

2.3 Mobility Management Entity

The MME is the key control-node for the LTE accessnetwork. It is responsible for idle mode UE tracking and paging procedure including retransmissions. It is involved in the bearer activation/deactivation process and is also responsible for choosing the SGW for a UE at the initial attach and at time of intra-LTE handover involving Core Network (CN) node relocation. It is responsible for authenticating the user (by interacting with the HSS). The Non-Access Stratum (NAS) signalling terminates at the MME and it is also responsible for generation and allocation of temporary identities to UEs. It checks the authorization of the UE to camp on the service termination point in the network for ciphering/integrity protection for NAS signaling and handles the security key management. Lawful interception of signaling is also supported by the MME. The MME also provides the control plane function for mobility between LTE and 2G/3G access networks with the S3 interface terminating at the MME from the SGSN. The MME also terminates the S6a interface towards the home HSS for roaming UEs. An architecture of the Long Term Evaluation System, its components and relationship among system and the components is presented in the Figure 1 given below.



Figure 01: An Architecture of LTE System

2.4 Home Subscriber Server

The Home Subscribe Server is the main Internet Multimedia Server IMS database which also acts as database in evolved packet core EPC networks. The HSS is a super HLR that combined legacy HLR and AuC functions together for circuit switched (CS) and packet switched (PS) domains. In the IMS architecture, the HSS connects to application servers as well as the Call Session Control Function (CSCF) using the DIAMETER protocol.

HSS is the master repository for subscriber profiles, device profiles, and state information. As a mandatory control plane function in the LTE Evolved Packet Core (EPC), the HSS manages subscriber identities, service profiles, authentication, authorization, and QoS for LTE and IP Multimedia Subsystem (IMS) networks.

From the above model components it is clear that this system is very much complex in terms of hardware equipment, communication protocols, and multimedia applications any time and anywhere.

3 Formal Specification

In this section, formal specification of static structures of the Long Term Evaluation (LTE) system is presented using Z notation. The schema is a powerful structure in Z notation used for variables definitions, components encapsulation and further used for defining properties in terms of invariants and predicates. In the description, first of all, fundamental definitions used in the LTE model are presented then specification of message is described. In the next formal description of user equipment which is the most important component of the network is given. Then formal descriptions of evolved NodeB (eNB), Home Subscriber Server (HSS) and Serving Gateway (SGW) are presented. Finally, Mobility Management Entity (MME) is described based on the composition of the above structures.

3.1 Fundamental Definitions

Formal definition of message is described below using the schema *Message* in Z notation. The schema consists of six components, namely, message identifier, data stored, sender information, receiving person, sending time and receiving time. The schema consists of two parts in addition to the name of the schema written in the first horizontal line. Definitions of variables used are given in first part of the schema and invariants are defined in the second part. The message identifier is represented by *MessageId*. The data sent or received is defined as a sequence of strings. The source and sending persons are identified by phone identifiers. The sending and receiving times are denoted by *Time*. After giving all of these definitions, it is stated that source cannot be equal to the destination in the message schema.

[MessageId, PhoneId]; [String, Time]

| Message | |
|------------------------------|---|
| mid: MessageId | |
| data: seq String | |
| source, destination: PhoneId | |
| sending: Time | |
| receiving: Time | |
| | _ |
| source \neq destination | |

In LTE, user equipment is a device used directly by an enduser to communicate with other connected users. It can be a telephone, laptop or any other such device connected to the base station. The user equipment is defined by the schema UserEquipment given below. The equipment identifier is denoted by the EquipmentId. The equipment has three states that is active, idle or detached. Further, it has information about its position and contains sim identifier and key. The equipment has other two variables called control and voice frequencies. The value of frequency will be either allocated or null. The inbox and outbox are specified by the sequence of messages in the equipment. The message report has three values that is sent, received or failure. The phone book is described by the function from phone identifier to the person. In the predicate part of the schema, it is described that control frequency has allocated status if and only if the equipment status is either active or idle. The control frequency will be null if and only if the equipment is detached.

| UserEquipment | |
|-------------------------------------|--|
| ueid: EquipmentId | |
| uestate: UEState | |
| position: Position | |
| simid: SimId | |
| simkey: SimKey | |
| cfstatus, vfstatus: FStatus | |
| cfreq, vfreq: Frequency | |
| inbox, outbox: seq Message | |
| msreport: MSReport | |
| pbook: PhoneId \rightarrow Person | |
| | |

 $cfstatus = ALLOCATED \Leftrightarrow uestate \in \{IDLE, ACTIVE\}$ $cfstatus = NULL \Leftrightarrow uestate = DETACHED$

[EquipmentId, SimId, SimKey, Frequency, Person] UEState ::= ACTIVE | IDLE | DETACHED Position ::= EDGE | CENTRE MSReport ::= SENT | RECEIVED | FAILURE FStatus ::= NULL | ALLOCATED

3.2 LTE Components

An evolved Node B is the radio access part of LTE system which contains at least one radio transmitter, receiver and control section. The receivers contain resource management and logic control functions. The formal description of eNB, HSS and SGW are presented here. The evolved NodeB is denoted by the schema *eNB* given below. The schema consists of seven components. The first one node is identifier denoted by *NodeId*. The next three components are frequencies, namely, total set of frequencies, control and voice frequencies. The next component is a set of equipments connected to eNB. Finally, two functions for control frequencies and voice frequencies are described from a set frequency to equipment identifier.

[NodeId]

enb: Nodeld frequencies: \mathbb{F} Frequency cfrequencies: \mathbb{F} Frequency vfrequencies: \mathbb{F} Frequency uequipments: \mathbb{F} UserEquipment calloc: Frequency \rightarrow EquipmentId valloc: Frequency \rightarrow EquipmentId frequencies $\neq \emptyset \land$ cfrequencies $\neq \emptyset$ vfrequencies $\neq \emptyset \land$ uequipments $\neq \emptyset$ cfrequencies \cap vfrequencies $= \emptyset$

 $frequencies = cfrequencies \cup vfrequencies$

 $\begin{array}{l} \forall ue: UserEquipment \mid ue \in uequipments \cdot ue .cfreq \in cfrequencies \\ \forall ue: UserEquipment \mid ue \in uequipments \cdot ue .vfreq \in vfrequencies \\ \forall f1, f2: Frequency; ei: EquipmentId \mid (f1, ei) \in calloc \land (f2, ei) \in calloc \cdot f1 = f2 \\ \forall f1, f2: Frequency; ei: EquipmentId \mid (f1, ei) \in valloc \land (f2, ei) \in valloc \cdot f1 = f2 \\ \forall fr: Frequency \mid fr \in dom calloc \cdot fr \in cfrequencies \\ \forall fr: Frequency \mid fr \in dom valloc \cdot fr \in vfrequencies \\ \forall ue1: EquipmentId \mid ue1 \in ran calloc \\ \cdot \exists ue2: UserEquipment \mid ue2 \in uequipments \cdot ue1 = ue2 . ueid \\ \forall ue1: EquipmentId \mid ue1 \in ran calloc \\ \cdot \exists ue2: UserEquipment \mid ue2 \in uequipments \cdot ue1 = ue2 . ueid \\ dom calloc \cap dom valloc = \emptyset \\ ran calloc \cap ran valloc = \emptyset \end{array}$

Invariants: (i) All the three sets total frequencies, control frequencies and voice are non-empty. (ii) The set of user equipments is also non-empty. (iii) Intersection of control frequencies and voice frequencies is an empty set. (iv) The union of control frequencies and voice frequencies is equal to set of total frequencies. (v) Control frequency of each user equipment, is in the set of control frequencies of evolved node. (vi) Voice frequency of each user equipment, is in the set of voice frequencies of evolved node. (vii) The control and voice frequencies functions are one to one functions. (viii) The domain of control frequency function is subset of control frequencies of evolved node. (ix) The domain of voice frequency function is subset of voice frequencies of evolved node. (x) The range of control frequency function is subset of set of equipments of evolved node. (xi) The range of voice frequency function is subset of set of equipments of evolved

node. (xii) Intersection of domains of control and voice frequencies functions are empty. (xiii) Intersection of ranges of control and voice frequencies functions are empty.

The Home Subscriber Server (HSS) manages subscription related information. The HSS supports the network control layer with subscription and session handling providing capabilities for equipment management, security, user identification handling and, access and service authorization. The HSS consists of various components some of the important information is given below in the schema *HSS*. The first one component is access having value granted or not granted. The next three components are for storing record about sim identifiers, sim keys and equipments identifiers. The relationship between sim identifiers with keys and equipments are also defined. Home location and visitor location registers are described by HLRId and VLRId respectively. Three types of networks, GSM, SGSN and 2G/3G, are considered in the schema.

[HLRId]; [VLRId]

Access ::= GRANTED | NOTGRANTED

RadioAccessNetworks ::= GSM | SGSN | G23

| H | <u>ISS</u> |
|--------------|--|
| acc | ess: Access ; sims: F SimId |
| sim | keys: F SimKey ; equipments: F EquipmentId |
| assi | ign: $SimId \rightarrow SimKey$ |
| srea | cord: $SimId \rightarrow EquipmentId$ |
| hlri | d: HLRId ; vlrid: VLRId |
| netv | vorktype: RadioAccessNetworks |
| | <u></u> |
| $\forall si$ | d: $SimId \mid sid \in sims \cdot sid \in dom \ srecord$ |
| $\forall si$ | $d: SimId \mid sid \in dom \ srecord \cdot \ sid \in sims$ |
| $\forall si$ | $d: SimId \mid sid \in sims \cdot sid \in dom \ assign$ |
| $\forall si$ | $d: SimId \mid sid \in dom \ assign \cdot sid \in sims$ |
| $\forall sk$ | $: SimKey \mid sk \in simkeys \cdot sk \in ran \ assign$ |
| $\forall sk$ | $: SimKey \mid sk \in ran \ assign \cdot sk \in simkeys$ |
| ∀ei | : EquipmentId $ ei \in equipments \cdot ei \in ran \ srecord$ |
| ∀ei | : EquipmentId $ ei \in ran \ srecord \cdot ei \in equipments$ |

Invariants: (i) Each sim identifier in the HSS is in the domain of sim-key record function. (ii) Each sim in the domain of sim-key record function is in the HSS. (iii) Each sim identifier in the HSS is in the domain of sim-equipment record function. (iv) Each sim in the domain of sim-equipment record function is in the HSS. (v) Each sim key in the HSS is in the range of sim-key record function. (vi) Each sim key in the range of sim-key record function is in the HSS. (vii) Each equipment record function. (vii) Each equipment record function. (viii) Each equipment identifier in the range of sim-equipment record function is in the HSS.

The primary functionality of Surving Gateway (SGW) is to manage user mobility and act as a segregation point between the radio access network and the other core networks. The SGW maintains data movement between evolved nodes and the PDN Gateway. In a functional point of view, SGW is a termination point of the packet data network interface. Considering functionality of SGW, formal specification of SGW is given below using the schema *SGW*. The schema consists of three components, namely, evolved nodes, data and networks. The evolved nodes is a power set of *eNB*. The data is defined as a sequence of *Data* where *Data* is a basic set type. The radio access network types are same as defined above in the definition of HSS schema.

[Data]

____SGW_____ enbs: F eNB data: seq Data networks: RadioAccessNetworks

3.3 Formalizing Mobility Management Entity

The Mobility Management Entity (MME) is the key control node for the Long Term Evaluation network. The primary functionality of MME is monitoring tracking and paging procedure for idle mode user equipments. It is involved in activation and deactivation processes and is responsible for choosing the SGW for a user equipment at the initial process. Further, it is responsible for security issues including authentication of the user by interacting with the home subscriber server. It verifies the authorization of the user equipment to site on the service provider and enforces the equipment roaming restrictions. Furthermore, the MME provides the control function for mobility between Long Term Evaluation and other access networks.

| <i>MME</i> | |
|-------------------------------|--|
| enbs: F eNB | |
| hsss: \mathbb{F} HSS | |
| sgws: F SGW | |
| networks: RadioAccessNetworks | |
| | |

 $\forall enb: eNB \mid enb \in enbs$

- $\forall ue: UserEquipment \mid ue \in enb$. uequipments
- $\exists hss: HSS \mid hss \in hsss$
- \exists uei: EquipmentId |uei \in hss . equipments ue . ueid = uei \forall enb: eNB | enb \in enbs
- $\forall enb. env b \mid enb \in enbs$
- $\forall ue: UserEquipment \mid ue \in enb$. uequipments
 - $\exists hss: HSS \mid hss \in hsss$

• $\exists sid: SimId \mid sid \in hss . sims \cdot ue . simid = sid$

 $\forall enb: eNB \mid enb \in enbs$

- $\forall ue: UserEquipment \mid ue \in enb$. uequipments
- $\exists hss: HSS \mid hss \in hsss$
- $\exists sk: SimKey \mid sk \in hss . simkeys ue . simkey = sk$

 $\forall enb: eNB \mid enb \in enbs \cdot \exists sgw: SGW \mid sgw \in sgws \cdot enb \in sgw. enbs \\ \forall sgw: SGW \mid sgw \in sgws \cdot \forall enb: eNB \mid enb \in sgw . enbs \cdot enb \in enbs \\ \end{cases}$

Formal specification of Mobility Management Entity is given above using the schema MME. The schema consists of four

main components, namely, evolved node, home subscriber server, surving gateway and radio access networks. The formal definitions of the components are given in terms of schemas. Evolved nodes, home subscriber servers and surving gateways are assumed as power sets. Similar to above, all sets are assumed as finite collections of the components. All the four components are put in first part of the schema and invariants are defined in the second part of it.

Invariants: (i) Each evolved node in MME is contained in some of the home subscriber server. (ii) Each evolved node in home subscriber server is contained in MME. (iii) Each sim identifier in every equipment of every evolved node is in some of the home subscriber server. (iv) Every evolved node in MME is in some of the surving gateway. (v) Every evolved node in every surving gateway is in MME.

4 Model Analysis

In this section, formal analysis of the specification is provided. As we know there does not exist any real computer tool which may assure about complete correctness of a formal model, therefore, even the specification is well-written it may contain potential errors. It means an art of writing specification does not provide guarantee about correctness of a system. However, if the specification is analyzed and validated with the rigorous computer tools it certainly increases a confidence over the system by identifying errors, if exists.

Z/Eves is a powerful tool used here for analyzing the formal specification of long term evaluation system. A snapshot of the specification analysis is presented in Figure 2. The first column on the left of the Figure shows status of the syntax checking. The second column represents the proof correctness of the specification. The symbol 'Y' stands that specification is correct syntactically and proof is also correct. The symbol 'N' is used that errors exist which needs to be fixed. As formal specification of our systems is fully analyzed to be correct and hence 'N' does not appear in the Figure. All the schemas are checked to be correct in syntax and has a correct proof.

| Ta Ca Great Vide | | |
|------------------|--|----------------|
| YY | [Nodeld] | |
| Y Y | eNB | |
| | enb: NodeId | |
| | frequencies: F Frequency | |
| | cfrequencies: F Frequency | |
| | vfrequencies: F Frequency | |
| | uequipments: F UserEquipment | |
| | calloc: Frequency \rightarrow EquipmentId | |
| | valloc: Frequency \rightarrow EquipmentId | |
| | frequencies $\pm \emptyset \land$ cfrequencies $\pm \emptyset$ | |
| e e 🕤 | o 😽 🕊 🕱 🐨 | - 76 10 4 BOAR |

Figure 2. Snapshot of the Model Analysis.

In Table 1 given below, summary of the results is presented. In first column of the table, name of schema is given. The symbol "Y" in column 2 indicates that all schemas are wellwritten and proved automatically. Similarly, domain checking, reduction and proof by reduction are represented in columns 3, 4 & 5 respectively. The "NA" in column 4 is used that reduction was not required on the predicates and, hence, the specification is fully meaningful.

TABLE I.RESULTS OF MODEL ANALYSIS

| Schema Name | Syntax Type Check | Domain Check | Reduction | Proof |
|---------------|----------------------|-----------------|-----------|-------|
| Message | Y | Y | NA | Y |
| UserEquipment | Y | Y | NA | Y |
| eNB | Y | Y | NA | Y |
| HSS | Y | Y | NA | Y |
| SGW | Y | Y | NA | Y |
| MME | Y | Y | NA | Y |

5 Conclusion and Future Work

The Long Term Evaluation (LTE) is a complex mobile communication system which provides many real time multimedia applications. There exists some research work for analysis and design of the system and it needs much improvements for its modeling, specification and optimization due to its complexity and distributed nature [13], [14]. The telecom industry and academia are discovering network approaches and integration techniques for the system planning to propose scientific theories and solutions and, hence, it is assumed as an open research problem [15], [16], [17], [18].

In this paper, formal methods [19] are used for modeling, specification and analysis of the LTE system. It is observed that the use of formal methods was effective to model and verify properties of the system because of abstraction and rigorous computer tool support addressing complexity of the system. At first, formal specification of components of the LTE system was presented using Z notation. Then formal description of system was described after composition of the components. Formal specification is analyzed using Z/Eves toolset. Our experience has shown that Z notation was very useful at requirements analysis for further modeling and development of the system. The Z/Eves toolset increased our confidence for analyzing and validating the model because of its specification checking and analysis facilities.

For a moment, an abstract analysis of the LTE system and its components is provided at higher level of details. A detailed model will be presented for addressing other components, the interaction protocols and functionality at software level. The other challenges, for example, an ability to provide seamless and adaptive quality of service and optimization issues in such a heterogeneous environment will also be addressed.

References

[1] D. Singhal, M. Kunapareddy, V. Chetlapalli, V. B. James and N. Akhtar, "LTE-Advanced: Handover Interruption Time Analysis for IMT-A Evaluation", Proceedings of International Conference on Signal Processing, Communication, Computing and Networking Technologies, 2011.

- [2] Z. Bai, C. Spiegel, G. H. Bruck, P. Jung, M. Horvat, J. Berkmann, C. Drewes and B. Gunzelmann, "System Performance of UTRA LTE and LTE-Advanced", International Conference on Communications, Computing and Control Applications (CCCA), 2011.
- [3] Y. Wang, K. I. Pedersen, P. E. Mogensen, and T. B. Sørensen, "Resource Allocation Considerations for Multi-Carrier LTE-Advanced Systems Operating in Backward Compatible Mode", IEEE 20th International Symposium on Personal, Indoor, and Mobile Radio Communications, 2009.
- [4] R. Milner, "Communicating and Mobile Systems: The π -Calculus", Cambridge University Press, 1999.
- [5] S. Acharya, C. George and H. Mohanty, "Domain Consistency in Requirements Specification", Proceedings of the Fifth International Conference on Quality Software, 2005.
- [6] K. Taguchi, J. S. Dong, "An Overview of Mobile Object-Z", Springer-Verlag LNCS 2495: 144-155, 2002.
- [7] M. Bettaz, M. Maouche, "Towards Mobile Z Schemas", Department of Computer Science, International Journal of Computer Science and Application, 11(11), pp.101-17, 2005.
- [8] R. Duke, P. King, G. Rose, and G. Smith, "The Object-Z Specification Language", Department of Computer Science, University of Queensland 4072, Australia, 1991.
- [9] M. Bettaz, M. Maouche, "Mobile Z Notation", Research Report, Department of Computer Science, Philadelphia University, March 2002.
- [10] 3GPP TR 25.913. Requirements for Evolved UTRA and Evolved UTRAN. Available at http://www.3gpp.org.
- [11] 3GPP TS 23.401. GPRS Enhancements on EUTRAN Access. Available at http://www.3gpp.org.
- [12] Motorola Technical White Paper. "Long Term Evolution (LTE): A Technical Overview".
- [13] M. Mac and A. and M. Butler, "Service Specification using Z", Technical Report, Faculty of Physical and Applied Science, University of Southampton, UK, 1993.
- [14] R. Cam and S.Vuong, "A Formal Specification in LOTOS, of a Simplified Cellular Mobile Communication System", Formal Description Techniques-II (FORTE'89), pp. 485-99, 1990.
- [15] L. Chen et al., "System Level Simulation Methodology and Platform for Mobile Cellular Systems", IEEE Communication Magzine, 49(7), pp. 148-55, 2011
- [16] K. Doppler, "Device-to-Device Communication as an Underlay to LTE Advanced Networks", IEEE Communication Magzine, 47(12), pp. 42-49.
- [17] Y. Sun, "Editorial Advanced Signal Processing for Wireless and Mobile Communications", IET Signal Processing, 3(6), pp. 431–32, 2009.
- [18] J. M. Wing, "A Specifier's Introduction to Formal Methods", IEEE Computer, 23(9), pp.8-24, 1990.
- [19] H. Luo, et al., "Quality Driven Cross Layer Optimized Video Delivery over LTE", IEEE Communication Magzine, 48(2), pp. 102–109, 2010.

Dynamics Analysis of Air-filled Play Equipment under Condition of Multiple Repetitive Impacts

Yuki Tokoro*1,*2, Yoshifumi Nishida*2, Ilwoong Kim*2, Hiroshi Mizoguchi*1,*2

*1 Tokyo University of Science

*² Digital Human Research Center,

National Institute of Advanced Industrial Science and Technology (AIST)

Abstract—Recently, unintentional childhood injuries due to popular play equipment consisting of an air-inflated membrane structure have been increasing. However, we lack a basic understanding of such play equipment, and so qualitative standards for safe products and safe use have not been established. This study conducts experiments using a newly developed device, a child's femur impactor, to develop a kinetic model. The experimental results clarify the dynamics of an air-inflated membrane structure not only when a child falls on the structure but also when both a child and an adult jump repetitively. In particular, the experimental results reveal that the shock absorption performance of air-filled play equipment is changed by multiple repetitive external impacts. This paper also describes our development of a finite element model of an air-inflated membrane structure to analyze the newly found phenomenon of the propagation of internal air pressure and membrane tension.

Key Words: Biomechanics, Accident analysis, Play equipment, Air-inflated membrane structure, Risk assessment

1 Introduction

In recent years, the number of accidents of children playing on equipment consisting of an air-inflated membrane structure has been rapidly increasing [1]. Figure 1 shows an example of this type of play equipment. Typically, the membrane is made of material with high durability and airtightness. The inside of the sealed membrane is inflated with air. The internal air pressure is higher than the atmospheric pressure. This internal pressure generates repulsive forces against a child jumping and also works as shock absorption when a child falls. It has been assumed that this air-inflated membrane structure is safe enough to prevent serious injuries. However, bone fractures occur due to falling on the shock-absorbing air cushion of the equipment. Because this air-inflated membrane structure is a new type of play equipment, safety standards and operating manuals for the equipment are not based on scientific evidence. Research to understand the dynamics of these air-inflated membrane structures is strongly required.

This paper considers the following case as one of situations that causes serious injuries and and have been investigated in detail.

Case A child fractures his femur when he is jumping on the play equipment and an adult is jumping around him.

In this study, we adopt two methods to investigate and model the characteristics of the play equipment consisting of the airinflated membrane structure. First, we develop a child's femur impactor to measure the impact force on the femur when a child falls on the play equipment and examine the characteristics of the play equipment. Second, we develop a finite element model for explaining the resulting measurements.

2 Development of child's femur impactor

We developed a child's femur impactor to measure the impact force in a fall. Figure 2 shows the configuration of the

Size:3.45W×2.45H×3.45D(m) Fig. 1. A sample of equipment consisting of an air-inflated membrane

developed impactor. Figure 3 shows the weights attached to the impactor. By changing the attached weights, we can investigate the effect of body weight on the femur. The weight is attached to the top of the impactor. The weights used are 5 kg, 10 kg, and 15 kg weights. The specifications of the femur-impactor are listed in Table I. The size of the femur impactor is the same as the average size of a 3-year-old Japanese child. Force sensors are installed on the top and the bottom of the impactor and an acceleration sensor is installed on the central part of the impactor.

To measure the characteristics of the equipment in the situation that a child jumps repetitively, we need to generate a repetitive external force. In addition to the impactor, we also developed a special guide. Figure 4 shows the guide and the impactor. This guide keeps the movement of the impactor in a direction perpendicular to the ground (i.e., membrane). We can thus generate a repetitive impact by repeatedly raising and dropping the impactor manually. The guide and the impactor used for simulating the situation that a child is jumping repetitively. As for adult's jump, which is dealt with later, we utilze a real adult and the adult's impact is not measured.

3 Experiment on air-inflated membrane structure

Using the developed impactor, we carried out experiments on an air-inflated membrane structure. Typically, the membranes maintain the structure inside the equipment. Each membrane has a hole for transmitting air. We measured the inside air pressure. The pressure is 9.80 hPa (gauge pressure) in a stationary state.





Fig. 2. Developed femur impactor



Fig. 3. Weights attached to the femur impactor

3.1 Drop impact force in case of single impact

First we conducted drop experiments to study the characteristics of the equipment in the situation that a single child fall from a high position to the cushion portion of the play equipment. The developed femur impactor was used to measure the impact force applied to the femur when landing. To study the relation between the impact force, the total weight of the impactor, and the collision velocity, we changed the conditions of the weight attached to the impactor and the fall height. In this paper, the collision velocity is the speed immediately prior to contacting the membrane. We set the fall height to 0.4 m, 0.7 m, and 1.0 m and the total weight of

 TABLE I

 FEMUR-IMPACTOR SPECIFICATIONS

| Gross weight | 2.35 kg |
|--------------------------------------|--------------------------|
| Range of acceleration | $\pm 1960 \text{ m/s}^2$ |
| Detection axis of acceleratio sensor | 3 axes (x, y, z) |
| Resolution of acceleratio sensor | 1.60 m/s^2 |
| Range of load | \pm 20 kN |
| Detection axis of force sensor | vertical axis |
| Resolution of force sensor | 5.000 N |
| Sampling period | l kHz |



Fig. 4. Installed femur impactor and guide

the impactor was set to 7.35 kg, 12.37 kg, and 17.35 kg, as summarized in Table II.

Figure 5 shows the relation between the impact and the collision velocity and Fig.6 shows the relation between the impact and the weight. These graphs indicate that the impact force in the case of a single fall is proportional to both the collision velocity and the weight.



Fig. 5. Relation between collision velocity and impact

Using these experimental results, we calculated the impact force applied to a child's body. Figure 7 shows the relation between the impact force, the total weight of the impactor, and the collision velocity. The color scale indicates the magnitude of the impact force. Figure 7 is used later for a risk assessment.

| TABLE II Experimental condition | | | | |
|------------------------------------|--|--|--|--|
| Fall height 0.4m, 0.7m, 1.0m | | | | |
| Weight 7.35kg, 12.35kg, 17.35kg | | | | |

3.2 Drop impact force in case of multiple repetitive impacts

We next conducted experiments to study the characteristics of the play equipment when two or more people are jumping.



Fig. 6. Relation between weight and impact



Fig. 7. Relation between single impact, weight, and collision velocity

One typical example is the situation that a child is jumping and an adult is also jumping near the child. In this case, two external impact sources are present and the two impacts are repeated. We conducted two kinds of experiments. In the first experiment, an adult rides by standing on the equipment without jumping, shown in Fig.8. Simultaneously, we repeated the raising and dropping of the developed impactor and measured the impact force. In the second experiment, an adult whose weight is 60 kg jumps repetitively, shown in Fig.8. Simultaneously, we repeated the raising and dropping of the developed impactor. In the two experiments, the repetitive dropping of the impactor corresponded to the child's jump.

Figure 9 shows the result of the first experiment. The figure compares the case when the adult rides and the impactor falls repetitively and the case when no adult rides and the impactor falls repetitively. Figure 9 indicates that the impact force does not change when an adult rides at a distance from the impact fall point.

Figure 10 shows the result of the second experiment. The figure compares the case when no adult jumps and the impactor falls repetitively and the case when an adult jumps repetitively and the impactor falls repetitively. The max and min of each graph indicate the maximum and minimum value of the impact, respectively. From Fig.10, we find that the jumping adult can increase the child's impact by 60% in the maximum case and that, similarly, the jumping adult can decrease the child's impact by 35%. Namely, the case of the jumping adult is much different from the case of no adult. This phenomenon has not yet been investigated conventionally.



Fig. 8. Measured point



Fig. 9. Relation of impact and fall height in the case that an adult also rides on the play equipment without repetitive jumps

3.3 Investigation of the effect of the internal pressure change

We investigated the effect of the change of the internal steady-state air pressure of the play equipment. We changed the internal steady-state air pressure by reducing the amount of air transmitted by blowing air. In this experiment, the fall height of the impactor was 0.7 m and the attached weight was 13.25 kg. The pressure was changed from approximately 10 to 4 hPa.

Similarly, we measured the impact force when the internal steady-state pressure was reduced. In this experiment, the fall height of the impactor was also 0.7 m and the attached weight was 13.25 kg. The pressure was changed from approximately 10 to 4 hPa. Figure 11 shows the relation between the measured impact force and the reduced pressure. It is clear from Fig. 11 that the impact force is not significantly associated with the decreasing internal pressure of the play equipment.

3.4 Calculation of buckling load for child femur

We calculated the child's femur buckling load. Because we lack the material properties of living tissue such as muscle, we







Fig. 11. Impact force when the internal pressure is reduced (the comparison of the single fall condition and multiple repetitive impacts condition)

calculated the material strength of the femur by assuming the femur was bone without muscle around the bone. Despite this limitation, we can conduct the worst-case analysis by using the calculated value.

In the calculation of the buckling load, a femur can be approximated as a hollow cylinder. We calculate the Euler's buckling load by assuming that the both ends of the hollow cylinder are hinged ends. The length of the hollow cylinder indicates the length of a femur and the radius of the hollow cylinder indicates the radius of the femur.

Table III shows the size of femur of adults.

| TABLE III SIZE OF FEMUR OF ADULTS[2] | | | | |
|--|--------|--|--|--|
| Femur length | 455 mm | | | |
| External radius | 14 mm | | | |
| Internal radius | 6 mm | | | |

The femur length of the injured child is 325.4mm. Using this value and external and internal radius of the adults, we calculate the the external and internal radius of the child by scaling the adult's values[3]. Table IV shows the sizes of femur of the injured child.

TABLE IV SIZE OF FEMUR OF INJURED CHILD

| 0 | 525.1 mm |
|-----------------|----------|
| External radius | 10.0 mm |
| Internal radius | 4.29 mm |

According to the previous study[4], Young 's modulus of 6 year-old is 6.6 GPa. The area moment of inertia for a circular section can be expressed by

$$I = \frac{\pi}{4}(R^4 - r^4),$$
 (1)

where r and R indicate the internal radius and the external radius of the hollow cylinder. Euler's buckling load can be calculated by the following equation.

$$P = \frac{\pi^2}{l^2} EI.$$
 (2)

By substituting numerical values into the above equations, we obtained 4691.4 N as the buckling load.

Let us consider an example of risk estimation. The 95% percentile value of a 6-year-old boy's weight is 24.9 kg. When a child falls from 1.5 m height, the collision velocity becomes 5.5 m/s². Using Fig.7, we know that the impact force applied to a child whose weight is 24.9 kg is 2794.3 N. This calculation suggests that the probability of femur fracture is low when a single child is jumping. However It also suggests that the probability of femur fracture becomes high when a child is jumping and an adult is also jumping around the child since this situation can increase impact by 60 % according to the experimental results described in the previous section. Thus, by measuring the characteristics of the play equipment and the developed impactor, we can estimate the risk of bone fracture. This contributes to the risk assessment of the play equipment in the design phase since the fall height can be changed by improving the design.

4 Development of finite element model of play equipment

Development of a finite element model of the play equipment consisting of an air-inflated membrane structure is one of the most effective approaches to supporting the design of this equipment. As stated above, the characteristics of this type of play equipment are very complex and the absorption property significantly changes depending on the play situation. For example, when a person jumps on the air-inflated membrane structure, the inside air is pushed away from the landing point and the membrane around the landing point is stretched. This leads to an increase of the air pressure and the membrane tension in a local region. The effect of the pressure and tension increase propagates to a wide area inside the play equipment very rapidly. These phenomenon increases the impact force applied to another child. This situation is very complex for a designer estimating the risk. So, we developed a finite element model (FEM) of the air-inflated membrane structure.

4.1 Development of FEM of the air-inflated play equipment

To develop the FEM of the air-inflated play equipment, we need the material properties and the shape of the play equipment. First, we describe the material properties required for the simulation. A tensile test was carried out to investigate the film material of the air-inflated membrane structure. Figure 12 shows a photo of the tensile testing machine. We created a test piece of 10 mm width from the film material. The stressstrain diagram of the film obtained from the test is shown in Fig. 13. Young's modulus obtained by the tensile test was 0.918GPa. The material properties of the membrane are shown in Table V. The air pressure inside the air-inflated membrane structure was 15.91 hPa.

We used the above values for the material properties in the FEM. However, we simplified the shape and the structure of the equipment. The inside of the actual play equipment has many sheets of membrane, which we did not model. In this study, we created a base-level finite element model that approximates a rectangular air-inflated membrane structure.



Fig. 12. Tensile testing machine



Fig. 13. Membrane stress-strain diagram. The right side shows the fitted curve used in the analysis

TABLE V Material properties of the membrane

| Thickness | 0.451mm |
|-----------------|--------------------------|
| Young's modulus | 0.918Gpa |
| Poisson's ratio | 0.3 |
| Density | 1171.4 kg/m ³ |

4.2 The simulation using the developed model

To prove the effectiveness of the developed model, we conducted simulations for reproducing the propagation of the internal air pressure and the membrane tension on the conditions of single impact and multiple impacts by using the developed FEM. Figure 14 shows the model for the single-impact condition and Fig.15 shows the model for the multiple-impact condition. The weight of the falling object is 7.24 kg and the fall height is 1.5 m. Figure 16 shows propagation of the internal air pressure under the single-impact condition. Figure 17 shows the propagation of the membrane tension under the dual impacts condition. These figures indicate that the simulation can reproduce the propagation of the internal air pressure and the membrane tension.

Next we confirmed that the simulation explains the phenomenon that the impact generated by one falling object can increase the impact of another falling object. This is important for explaining the new finding that the propagation of the internal air pressure and the membrane tension generated by an jumping adult can cause the fracture of the femur bone of a child at a distance away from the adult. To confirm this, we compared the result of the single-impact condition and the result of the multiple-impact condition (two-impact condition). We compared the maximum values under these two conditions. Figure 18 compares the maximum values. The figure proves that the simulation explains that the multiple impact can cause a 57% increase of impact.



Fig. 14. The FEM under the single-impact condition



Fig. 15. The FEM under the multiple-impact condition



Fig. 18. Comparison between single-impact and multiple-impact conditions



Fig. 16. Simulation of the propagation of the internal air pressure



Fig. 17. Simulation of the propagation of the membrane tension

5 Conclusions

In this study, to understand the dynamics of play equipment consisting of an air-inflated membrane structure that has not yet been studied in detail, we developed an impact force measuring device and a simulation model. First, we developed a child's femoral impactor to measure the impact force for examining the characteristics of the play equipment. The experiment using the developed impactor revealed the new finding that the propagation of the internal air pressure and the membrane tension generated by a jumping adult can cause a fracture of the femur bone of a child at a distance away from the adult. We also presented a risk assessment using the results of the developed impactor measurements. Second, we developed a finite element model of the play equipment consisting of an air-inflated membrane structure and carried

out a simulation. The simulation results confirmed that the model was able to simulate the propagation of the internal air pressure and the membrane tension.

References

- [1] Quality of Life Council Comprehensive Planning Committee, Documents
- [2]
- Quality of Life Council Comprehensive Planning Committee, Documents on risk of play equipment. Handout 4.2.1, 2007.
 A. Krauze, M. Kaczmarek, J. Marciniak, "Numerical analysis of femur in living and dead phase ", Journal of Achievements in Materials and Manufacturing Engineering, Vol. 26, (2008), pp. 163-166.
 K. Furukawa, I. Watanabe, K. Mikik, " A Development of 6-year-old Child FE Model ", The Japan Society of Mechanical Engineers, No.05-2, (2005), pp. 229-230.
 A. Irwin, H. J. Mertz, "Biomechanical Basis for the CRABI and Hybrid III Child Dummies ", SAE Paper 973317, (1997). [3]
- [4]

A System Dynamics Approach for Complex Government Policies Design. Application in ICT Diffusion

Yannis Charalabidis¹, Euripidis Loukis¹, and Aggeliki Androutsopoulou¹

¹Department of Information and communication Systems Engineering, Information Systems Laboratory, University of the Aegean, Karlovassi, Samos, Greece

Abstract - In order to achieve e-governance, we are in need of new and more advanced tools, specifically designed towards supporting the policy making procedure. The purpose of this paper is to investigate the perspectives, provided by the development of decision support tools, to confront complex egovernment phenomena. The analysis is performed using a System Dynamics simulation model that enables policy makers to investigate the estimated impact of planned government initiatives. Simulation applies on the diffusion of Internet and Communication Technology. The development of the model, made in collaboration with the Observatory for the Greek Information Society, addresses the digital divide in Greece. Data from the i2010 initiative indicators have been used for the simulation. The results, arising from the execution of alternative scenarios, indicate the parameters to be changed through the implementation of actions to have the best impact on society.

Keywords: e-Inclusion, modelling, system dynamics, policy making

1 Introduction

Nowadays, more than ever, the need for improvements in public administration is evident. Whereas e-Governance is establishing as the most important public sector reform strategy, the potential of ICT in administrative processes should be exploited. For the transition to e-Governance is necessary to utilize the power ICT on public services to promote evidence-based, massively participative, justified and finally effective decision making.

On the other hand, the smooth operation of e-Governance functions presupposes a high level of e-Inclusion in the society, meaning the participation of all individuals in all aspects of information society, including dealing with public services. To foster e-Inclusion, EU has planned policies that aim at reducing gaps in ICT usage and promoting the use of ICT to overcome exclusion, and improve economic performance, employment opportunities, quality of life, social participation and cohesion. This perspective shifts the "public policy problem" of the digital divide from a matter of pure social inequality to a strategic issue in a global race for competitiveness [6]. In such complex phenomena, national and regional governments might have to deal with unprecedented challenges requiring a profound rethinking of the policy and strategic approaches to be implemented. Thus, the e-Inclusion issue calls for a deep understanding of the problem dynamics and dimensions, necessary to identify timely policy responses and to minimize errors.

In order to face the aforementioned problem and other similar phenomena, elected representatives and politicians should be empowered with ICT tools that will enhance the policy formulation process. For example, decision support tools allowing testing different policies and forecasting their impact, may provide significant policy assistance to identify efficient policy actions. These tools offer a dynamic assessment of policies that can result in more informed political choices and can be achieved through modelling approaches [20]. However, in order to confront complex government issues the limits of traditional analytical models should be overcame. In this direction, special interest has arisen in the domain known as Policy Modelling. A state of the art review of existing modelling and simulation approaches [3] has indicated the potential of their application in societal problems depending on their characteristics. From the aforementioned investigation of the above approaches it was concluded that System Dynamics and Agent Based Modelling are the most promising to represent complex social phenomena under the policy making perspective.

Herein, System Dynamics was chosen to demonstrate that models can change the way policy making is done. The main purpose of the study is therefore to investigate how and to what extent external factors can influence the ICT diffusion process, seeking to acquire knowledge about the society and citizens behavior. The ICT diffusion process is analyzed herein by means of a System Dynamics simulation model. Several scenarios are applied, representing alternative policy actions in order their impact to be forecasted. The development of the model is part of the preliminary work conducted in the context of the PADGETS, a research project in the domain of "ICT for Governance and Policy modelling", supported by the Seventh Framework Programme of the European Commission.

This paper is structured as follows; initially a literature review on previous studies regarding the system dynamics method implementation is proposed. Section 3 presents the system dynamics model that was developed, scope of the model, assumptions and limitations and the linkages between each factor are described. Section 4 provides details about and the scenarios implemented for each policy (simulation parameters) and simulation results. Discussion on the results follows in section 5. Finally, the last section examines the perspectives for further research and the conclusions that are dawn are discussed.

2 System Dynamics in related studies

To investigate the current state of the art (in terms of practice and research) in the domain of system dynamics policy modeling, a study about the social simulation models particularly in matters relating to public administration was conducted. This section aims to provide an overview of indicative good practices that have already been implemented in the public sector and address specific needs associated with each research area in an innovative way. The literature review identified many examples of the System Dynamics methodology dealing with social phenomena [13].

The first model ever built, Urban Dynamics [7] remains a classic example of system dynamics successfully applied to an important public policy problem. The latter emulates the evolution of an imaginary city beginning from its growth, stagnation, and then decay. At the core of Urban Dynamics are the interactions between housing, business, and population sectors of an urban system. In addition to generating insight into the causes of urban decay, the urban model can also help policymakers design policies to improve decaying cities or prevent stagnation and decay in urban areas that are still growing.

But there are many contemporary examples of System Dynamics modeling tools as decision support from government agencies in various countries in the world as well. As technology is the main source of national growth a strategic choice was proposed in Turkey entitled as "National Science & Technology Policy" in order to enhance her ability in science and technology, and to get the capability of transforming them to economic and social benefit. To this direction a system dynamics model was constructed for policy analysis with respect to technology improvement as a tool to determine the various technology improvement policies [5]. The purpose of modeling was to understand technology improvement system, to identify the related entities with their effects on national technology improvement policy and to see the trend of the technology improvement in Turkey with respect to 15 years time. Through simulating some possible scenarios, decision makers could understand why some behavior pattern of the technology improvement policy system is occurring and to see what might be done to alter the pattern. System Dynamics appeared to be a potential tool for this field, although it contains a number of interactive and conflicting variables and parameters.

System dynamics methodology has been used in United States to capture the influence of tax policy on market competition between traditional telephone market and VOIP market [14]. Since, Voice Over Internet Protocol (VOIP) constituted the fastest-growing market in the United States providing telephone-like service without the restrictions of telecommunication regulations, state governments feared the implications on the heavily-taxed fixed line phone service, that was less tax revenue to support crucial public services. To help policy makers decide how to impose tax on VOIP services and reduce the impact of VOIP development, a system dynamics model was built to gain insight into interactions between the VOIP market, traditional phone market, and tax policy. The target audience of the model, mainly government officers involved in the development of telecommunications regulations and taxation policy, used it to test various policy settings to determine to what extent they should tax the new technologies to collect maximum tax revenue with less impact on the market. Through the tax policy tests made, it was concluded that tax policy has little influence on market competition, traditional phone markets will continue to decline no matter what the tax policy is and finally, a fixed fee for tax can ensure stable tax revenue.

System dynamics modeling has been applied to issues of population health since the 1970s. One such case was detected in the United States, concerning chronic disease prevention [11]. A relatively simple model was built, exploring how a hypothetical chronic disease population may be affected by 2 types of prevention: upstream prevention of disease onset, and downstream prevention of disease complications. The model was not intended for actual policy making but for exploration of the methodology as a means of modeling multiple interacting diseased populations, and matters of national and state policy. Within three different policy scenarios tested, the conclusion drawn was that often prevention of occurrence can be neglected by the prevention of complications resulting in utilizing ineffectively available resources

Another implementation of the method in taxation policy, which combines the approach to public health, is related with the tobacco industry in New Zealand [16]. The study, carried out under the cooperation of the New Zealand Customs Service (NZCS) and Ministry of Health (MOH), analyzed the relationship of Customs Service outputs to desired government outcomes, in relation to the collection of tobacco excise duties and cigarette smoking. The study was done to demonstrate the utility of system dynamics in answering some questions of a common type in the public sector during policy development and review. Policy experiments with the model examined the effects of changes in excise duties and the effects of a price change on tobacco related behaviors and provided some very useful insights into the customs and health-related activities associated with the supply and consumption of tobacco products in New Zealand.

The paper digitization problem preoccupied the Italian Public Administration during the transition to an all-digital society. The need to cut down the production of paper through the digitization process, in fact, implies a deep change in many aspects that surround the world of documents, from the techniques and instruments needed to accomplish all the activities, to the definition of roles and competences in the documents' management context. To tackle these aspects, a System Dynamics analysis study has been conducted with the support of the Italian National Centre for the Information Technologies into Public Administrations (CNIPA) in order to show how social and psychological factors may in the end determine policy resistance and great obstacles to organizational change [2]. The model showed to what extent a "complete" digitization process would be positive and profitable for the administrations and in which ways the digitization process would spread. The finding that needs to be noticed is that such a great and demanding change would have enormous positive effects on the public administration's activities in terms of efficiency. System Dynamics predicted that the consequences of digitization will be perceived over time and will be relevant to cost savings and time, but in much less environmental impact.

Dynamic simulation approach was also used in a research project originated within the 'Swiss Priority Programme Environment' (SPPE) of the Swiss National Science Foundation (SNSF) focused on ecological issues [17][19]. In this study, a model simulating the management of solid waste at the local level, exemplified the dynamic interaction between public policies and environmentally relevant behavior of citizens. The objective was to identify solutions for environmental management. Therefore, System Dynamics methodology was the core device for integrating concepts from different disciplines in a simulation model enabling the formulation of a dynamic theory in the context of environmental management and analysing issues of policy resistance and compliance.

The case study of modelling the situation of Iran cell phone market indicates that the methodology of System Dynamics could prevent an incorrect implemented government policy [22]. Towards the effort to decrease the dependency of the country from imported goods, government of Iran increased the rate of tariff of imported cell phones suddenly in 2005. This strategy had adverse effects on the development of the domestic mobile phone industry, a fact that could have been avoided if the model that was built later would have been used. After the execution of three alternative policy scenarios, it was concluded that the best decision for the government would have been to increase the tariff rate gradually in order to enable domestic manufacturers of mobile phones compete with their foreign rivals improving their quality at the same time.

Finally, the case of E-Mexico constitutes another recent implementation of System Dynamics methodology for an egovernment project [15]. This program is defined as the strategy to create web-based content to the citizen in the areas of education, health, economy and government. Using the same technological infrastructure and under the leadership of the same Federal Ministry, four different networks of government and nongovernment organizations engaged in the creation of Internet portals to create relevant content in these areas. In this case System Dynamics was used as an integrated and comprehensive approach to understand complex egovernment phenomena. The model that has been, represents a theory of how institutional, organizational and technology elements interact among them to produce different technology enactments.

Additional examples of System Dynamics implementations have been identified on the field of occupational safety, the energy sector, national security (military organization), in the adoption of automation technologies, in changes in sociopolitical structure of a country and generally in all areas related to the development of a country.

3 The ICT Diffusion Model

The model addresses the problem of the "Digital gap in Greece". Digital divide is a term coined to characterize the inequality in the relationship between groups of individuals and their relationship with formation communication technologies (ICTs). The application of the model can be considered as an attempt to foster a higher level of e-Inclusion and increase the diffusion of complementary activities such as eGovernment and eParticipation.

The model tries to conceptualise the multidimensional phenomenon of the ICT Diffusion on vulnerable social groups, monitoring the aspects of the concrete problem and the correlations between them. Thus, the model examines the factors that constitute the digital divide and captures quantitative data obtained from relative studies [4][9][10] with the view to predict what impact will result from implemented actions to them. Focusing on social groups at risk of digital exclusion, it models the process of their e-Inclusion, by simulating the evolution of any parameters of the issue.

To assess the digital divide an indicator system has been identified by the Greek Observatory [25] for IS that utilizes the annual surveys of the i2010 initiative conducted for households in Greece [24]. i2010 strategy is the EU general policy framework for the information society and media for the period 2005-2010 that provides the framework for eInclusion European action. The proposed model incorporates the indicators that are measured by the i2010 strategy concerning the following factors of digital divide [23]:

- Internet access and usage
- Availability of broadband infrastructure nationwide
- Internet literacy and digital skills

The affect on those indicators resulting by any policy actions are estimated by testing alternative policy scenarios. By altering one or more parameters of the model, the user is able to execute "what if" scenarios in order to compare different sets of policies and finally observe the overall impact on society.

The casual loop diagram of the model depicts all the system's parameters and the linkages between them. The arrows in the following schema illustrate the relationship between each of the factors involved in the behavior of the system. Plus sign symbolizes the positive effect of a variable to another and minus the negative one.



Figure 1. Casual loop diagram

The final structure of the model is presented with the following stock and flow diagram. Based on the previously presented casual loop diagram the stocks and flows, the basic building blocks of the System Dynamics methodology, were identified. A stock is the term for any entity that accumulates or depletes over time. A flow is the rate of change in a stock. The rest variables are auxiliary and contribute to the calculation of the flow rates. For example, in our case the percentage of broadband users increases according to the broadband take up rate reducing in parallel the stock of internet potential users. Then, the flow of broadband take up rate is determined by the cost of broadband access, the technology factor and the rate of broadband coverage.



Figure 2. Stock and flow diagram for ICT diffusion

4 Simulation of the Model

The model is designed with a time horizon of 15 years, from 2005 to 2020. Due to the rapid development of ICT and the ongoing actions towards the digital convergence, we consider the above time interval long enough to provide data on which we model the diffusion phenomenon. The data from 2005 to 2008, which are available through the measurements of

indicators by the "Greek Observatory" should be enough to see historical behaviors in order to set the data constants and formulate the equations between the variables. The model generates values of the model variables on an annual basis from 2008 to 2020, but can estimate values in shorter basis. The initial values of the indicators consist of the most recent
measures available. All percentages reflect the proportion of each cluster in relation with the total population in Greece.

In order to examine the impact of possible policy actions on the indicators four output variables have been selected and are being observed during the whole simulation lifecycle: the number of broadband users, the volume of Internet access at home, the percentage of Internet users (either broadband or not) and finally the percentage of digitally skilled people. There were 5 scenarios implemented for the analysis of different policy actions to reinforce the ICT diffusion. The outputs of the runs were evaluated accordingly and compared with the output corresponding to the base run simulation. Scenarios applied to the simulation model are given below.

4.1 Base case run of the model

The model output for the selected variables for the base case run of the model is provided in Figure 3 and Table 2. The base case run is the output from the simulation run from 2008 to 2020 with the initial set of model assumptions. As illustrated, all four variables are characterized by an upward trend and reach very high levels at the end of the simulation. The percentage of digitally skilled population presents the greatest growth rate, reflecting the accumulation of the increases of the rest variables that are affecting it.



Figure 3. Evolution of indicators in the base run

4.2 Scenario 1: Increase in broadband coverage rate

First scenario represents what will happen if government decided to intensify the broadband coverage projects. However we examine only the coverage rate, since the number of annual coverage projects depend on factors that are beyond the system boundaries. The figure concern only the variables that are affected in the specific plot and depict the policy result that is to achieve full broadband coverage in the country faster. As the broadband coverage is already high in Greece alternative actions should be considered to increase the number of broadband users.



Figure 4. Evolution of broadband coverage

4.3 Scenario 2: Reduction of access cost

Simulation second scenario assumes that a policy that reduces the cost of Internet access is planned, for example a government subsidy for the monthly fee for broadband access services. Another example of such an action is the programme "Diodos" implemented by the "General Secretariat for Research and Technology" that foresees special rates to purchase a broadband connection for students. The implementation of such a strategy will lead to significant growth in broadband and overall Internet access, which is expected since the smaller the cost will be more people will gain access. A small part of this increase will be transferred in the use of Internet but the impact on digital literacy will be minimal.



Figure 5. Increase of broadband users in Scenario 2

4.4 Scenario 3: Increase of access points

A similar policy with the previous one could be to increase the available public access points. The specific policy already applied within projects involving municipalities and regions in Greece and will be probably amplified through the installation of free Wi-Fi access points in transportation stations, an action that is currently under discussion. The aforementioned policy actions could have been tested with the proposed simulation tool before being taken decisions. As shown in the figure, the usage growth rate could present an initial increase, but will be declined after 2012 due to the remaining potential internet users, reflecting the balancing loop of the model. However this increase is not projected in the overall internet users, due to the fact that the majority of the population is online mainly from home and make use of access points just supplementary.



Figure 6. Evolution of usage rate

4.5 Scenario 4: Notification about ICT benefits

Another candidate policy for the dissemination of ICT could be targeted to inform citizens about the benefits offered by ICT. For example, a campaign to make citizens aware of the opportunities and services provided by the Internet will increase the overall perceived usefulness. The current scenario's output showed that a movement like this would entail a slight increase in usage, some of which will shift and increase digital skills. However, the change caused is so small that such policy action would not have a significant impact.

4.6 Scenario 5: Population training

The last scenario indicates the results would be expected by reinforcing citizens' digital skills. Educational training programs for various social groups can contribute to digital inclusion.

5 Simulation Results

The most recent surveys of the Greek Observatory validate the model as the simulation results of previous years (2008-2011) coincide with the real metrics at a large extent. Specifically, in 2010 already 46% of Greek households have internet access which converges with the prediction of the base run concerning the household internet access which was at 47,45%.

Simulation results suggest that an approach that combines more than one policy may be more effective, reflecting the reinforcing feedback that System Dynamics is based on. The above comparisons allowed us to understand that the most influencing factors are the cost of access and the population level of ICT skills, as the diffusion processes for the other scenarios did not show any significant difference. Therefore, an alternative scenario combining the second and the fifth policy choices could affect all the variables very significant as seen from the following charts. Figure 9 represents the output of simulations referred to the policies of access cost reduction and citizens training. As it can be observed, the percentage of digital literacy and the usage rate projected into the future reach very high levels in the last year of simulation.





Figure 7.Output of the simulation approach combining scenario 2 and 5

The above plot that depicts a desired evolution in the diffusion process can provide precious information for policy makers, highlighting the concept that they should design and imply specific policies. Overall, the implementation of the model aims to prove that similar techniques such as other modelling approaches and policy simulation experiments can also be used, as a grounded decision support system that informs and enhance the debate on policy design.

6 Conclusions

Previously the process of ICT diffusion in Greek society has been projected in the future through a System Dynamics simulation. The results obtained show that the factors, which exert a significant influence on ICT diffusion are related with the digital literacy of people and the cost of Internet access. However the model is a "prototype" only and hence does not contain all the variables and relationships that would be necessary to answer the research questions in depth. Although, the model is regarded as sufficiently "robust" to provide indicative results to the policy questions it can be further developed to include additional parameters of the problem and can be used the basis for a more comprehensive policy tool that could be further developed either for Greek and or any other country trying to struggle with the implications of policy relating to the digital gap. In addition it can be customized in order to overcome the drawback of not accounting external factors for, e.g., the socio-economic context of each country.

However the current paper does not only intend to present the results of a simulation of a real policy problem, but also to demonstrate the value of System Dynamics methodology and other policy modeling approaches in better understanding of social phenomena and in visualising their internal processes. Under this perspective, a significant aspect in model building is the engagement of policy makers. Civil servants, governmental actors and other stakeholders should be strongly involved in the preparatory work needed to be done before creating the model. Finally, future work should be performed to build a network of peers from the public and private sector interested in policy modelling, visualisation, mass collaborative platforms, large-scale societal simulations and in the use of these advanced tools by governments.

7 References

[1] Andersen, D. F., J. A. M. Vennix, et al. (2007). Group Model Building: Problem Structuring, Policy Simulation and Decision Support. The Journal of the Operational Research Society 58(5): 691-695.

[2] Armenia, S., L. Roma, et al. (2008). A new system dynamics model for the analysis of the paper digitization process in the Italian Public Administration. Proceedings of the 26th International Conference of the System Dynamics Society. Athens, Greece.

[3] Charalabidis, Y., Loukis, E., Androutsopoulou, A., (2011). Enhancing Participative Policy Making through Modelling and Simulation: A State of The Art Review. European Mediterranean Conference on Information Systems (EMCIS) 2011. 30-31 May 2011, Athens, Greece

[4] Becker, J., B. Niehaves, et al. (2008). Digital Divide in eGovernment: The eInclusion Gap Model. Electronic Government: 231-242.

[5] Durgun, M. S. (2003). A System Dynamics Approach for Technology Improvement Policy Analysis: The Case for Turkey, STPS - Science and Technology Policy Studies Center, Middle East Technical University.

[6] Ferro, E., J. Gil-Garcia, et al. (2007). The Digital Divide Metaphor: Understanding Paths to IT Literacy. Electronic Government: 265-280.

[7] Forrester, J. W. (1971). Counterintuitive behavior of social systems. Theory and Decision 2(2): 109-140.

[8] Ghaffarzadegan, N., J. Lyneis, et al. (2009). Why and How Small System Dynamics Models Can Help Policymakers: A Review of Two Public Policy Models. Proceedings of the 27th International Conference of the System Dynamics Society. Albuquerque, New Mexico, USA.

[9] Gil-Garcia, J., N. Helbig, et al. (2006). Is It Only About Internet Access? An Empirical Test of a Multi-dimensional Digital Divide. Electronic Government: 139-149. [10] Helbig, N., J. Ramon Gil-Garcva, et al. (2009). Understanding the complexity of electronic government: Implications from the digital divide literature. Government Information Quarterly 26(1): 89-97.

[11] Homer, J. B. and G. B. Hirsch (2006). System Dynamics Modeling for Public Health: Background and Opportunities. American Journal of Public Health 96(3): 452–458.

[12] Jiang, W. and H. Bin (2007). Modeling and simulation of group behavior in e-government implementation. Proceedings of the 39th conference on Winter simulation: 40 years! The best is yet to come. Washington D.C., IEEE Press.

[13] Kovacic, A. and B. Pecek (2007). Use of Simulation in a Public Administration Process. SIMULATION 83: 851-861.

[14] Liu, C.-Y. and W.-T. Wang (2005). System Dynamics Approach to Simulation of Tax Policy for Traditional and Internet Phone Services. Proceedings of the 23rd International Conference of the System Dynamics Society. Boston.

[15] Luna-Reyes, L. F. and J. R. Gil-García (2009). Using Institutional Theory and Dynamic Simulation to Understand Complex E-Government Phenomena. Proceedings of the 27th International Conference of the System Dynamics Society. Albuquerque, New Mexico, USA.

[16] Robert, Y. C. and V. C. Leslie (2006). Demonstrating the utility of system dynamics for public policy analysis in New Zealand: the case of excise tax policy on tobacco. System Dynamics Review 22(4): 321-348.

[17] Schwaninger, M., S. Ulli-Beer, et al. (2008). Policy Analysis and Design in Local Public Management A System Dynamics Approach. Handbook of Transdisciplinary Research: 205-221.

[18] Teekasap, P. (2009). Cluster Formation and Government Policy: System Dynamics Approach. Proceedings of the 27th International Conference of the System Dynamics Society. Albuquerque, New Mexico, USA.

[19] Ulli-Beer, S. (2003). Dynamic Interactions Between Citizen Choice and Preferences and Public Policy Initiatives -A System Dynamics Model of Recycling Dynamics in a Typical Swiss Locality. Proceedings of the 2003 International Conference of the System Dynamics Society. New York City, U.S.A.

[20] Volkery, A. and T. Ribeiro (2009). Scenario planning in public policy: Understanding use, impacts and the role of institutional context factors. Technological Forecasting and Social Change In Press, Corrected Proof.

[21] Zamanipour, M. (2009). A System Dynamics Model for Analyzing the Effects of Government Policies: A Case Study of Iran's Cell Phone Market. Proceedings of the 27th International Conference of the System Dynamics Society. Albuquerque, New Mexico, USA.

[22] Kountzeris, A. (2008). e-Inclusion and measurement of the Digital Divide. Observatory for Digital Greece.

[23] Kountzeris, A. and M. Konstantatos (2009). e-Inclusion and Digital Literacy in Greece. Observatory for Digital Greece

BEM for design inception Harnessing the power of clients' design intuition

A. Simondetti1, S. Roberts¹, and D. Birch² 1Arup, London, United Kingdom 2Department of Computing, Imperial College, London, United Kingdom

Abstract

Recentbestpracticedemonstrateshowbuiltenvironmentmodelling (BEM) enhances collaboration between the designer and other consultants. The designer's mission, however, is to "shape a better world", which involves more than design optimisation; it requires designers, including engineers, first exploring many solutions, as well as those that may not immediately be obvious. This mission is a radical one, and achieving it depends on building creative, collaborative partnerships with clients. To make the most of our clients' aspirations, ideas and insights, we need to empower them with an intuitive representation of relevant parameters. BEM offers excellent opportunities to enhance clients' involvement at design inception meetings by making it simple and cost effective to explore a wide variety of design options. This vision paper is informed by a review of digital design software that took place in 2010 and 2011.



Figure 1 Screenshot of HierSynth by David Birch. This computational framework uses the Unity3D games engine to publish interactive building massing models. CityEngine is used to procedurally generate the models; Radiance software generates daylight factors; Arup's Integrated Resource Management (IRM) generates carbon consumption data; and Imperial College's SynCity generates utility grids for gas and water

Context

BEM allows multiple aspects of a building, infrastructure or urban design project to be presented as a unified virtual 3D model. BEM facilitates holistic design, virtual testing and the optimisation of operations. Its benefits include better documentation, faster construction, built projects that operate more efficiently, and improved end-user familiarisation and training. BEM also has a social dimension; it unites in a shared digital world the people who commission, fund, design, construct and inhabit the built environment. Everyone – from the client to the general public – can better understand what's in it for them.

Commercial suppliers of BEM software have many visions for their products. These visions, however, do not always fully embrace the needs of built environment designers and their clients. This is because most of the software was originally developed for the automotive and aerospace industries, but the process of designing and constructing the built environment is different from the process of designing and manufacturing a car. In particular, built environment clients may have to deal with a very wide range of stakeholders as well as communities that are affected by their projects. And built environment designers do not own their supply chain and, therefore, often have to communicate with different audiences at each stage of a project.

To help us make full use of and shape the development of the most useful BEM tools, Arup's Foresight, Innovation & Incubation team is engaged in an ongoing review of new design software and innovative applications of existing software. More than simple desk research, the review involves testing and using a range of applications in collaboration with academics, as well as in-depth discussions with suppliers, talking to other users and participation in conferences and ideas workshops.

The vision of the **Future tools** programme is to move beyond the now proven benefits of component-based software (coordinated and up-to-date design information) and to harness the benefits of digital tools that can involve clients and other stakeholders in the design process.



Figure 2 This diagram of the design inception process is used throughout the paper to anchor different aspects of the vision. It details the sequence of steps that follow project inception. The "review" quadrant on the right represents the face-to-face discussions that allow the designer to harness clients' intuition.



Figure 3 Stages of project development. This vision focuses on design inception.

The challenge

The success of Arup's built environment projects is due in part to the creative partnerships it forges with clients and other members of the project team. A key part of Arup's role is to stimulate, understand and harness the client's expertise, ideas and design intuition. It is important to establish this collaborative relationship at the earliest stages of the design process because this is when the most significant decisions are made – the ones that offer the biggest opportunities for the client.

Ideally, design inception should be about exploring a variety of ideas to identify the best design solutions rather than focussing on getting a single solution right. The brief for a mixed use urban regeneration scheme, for example, might be to provide office, residential, retail and leisure space. But what is the best way to accommodate those uses on the site? Or in the case of a museum extension, where should it go? Beside, behind or above the existing building? Or could it be built underground?

Identifying the best design options involves balancing a wide set of factors including the developer's financial and commercial constraints, the user's requirements, the architect's aspirations and the engineer's parameters. To address this challenge, the vision here of BEM empowering collaborative design inception is inspired by Professor Sevil Sarivildiz and Dr Mi-Design 1



Figure 4 Diagram showing the performance of the two designs in terms of form, functionality and energy efficiency. The Pareto theoretical front of constant total performance shows that design 1 is best overall. However, design 2 is nearer the diagonal line that represents a balanced performance.



Figure 5 This plot shows the contrast between optimisation for one design to a local maximum of performance versus exploring different designs. Optimisation can proceed in small step changes through the design space, eventually finding the "peak". Searching for better designs, such as the ones to the right, require large steps, and harnessing design intuition can help make these. chael Bittermann's research at Delft University of Technology. The illustrations at the top of the facing page show output from their Intelligent Design Objects toolkit, which they created to help architects explore design solutions with developers. The output includes 3D representations of the various solutions that maximise design performances. A Pareto front diagram shows the performance of each solution in terms of three aspects – in this case form, functionality and energy use. However, these aspects of design do not address "hard" engineering parameters.

Each of the two designs could be optimised to a range of performance parameters, shown by isolated "performance" peaks in the lower diagram. However design inception should be about large steps across the design space to new regions, such as to the right currently hidden by the figurative cloud. Exploring through a process of intuition is about making these large steps.

The vision

Our vision is of a future where similar BEM tools will allow clients, engineers and other project team members collaboratively to explore a whole range of design solutions in the space of hours. The tools will use client-focussed parameters such as comfort as opposed to engineering-focus sed parameters such as steel tonnage. Analysing multiple parameters in parallel, the tools will unite and interrelate the expertise and intuition of the whole project team to identify the solution that most effectively balances the key requirements of all stakeholders. To fulfil this vision, BEM specialists will need to collaborate with clients to establish the relevant parameters.



Figure 6 Four stages of the exploratory process of design inception. Designs are iterated after each loop

This process of design exploration with clients has four stages (see the design inception diagrams on the facing page): setting the control parameters, computing these into a design solution, publishing the results and reviewing the design. Free-flowing discussion during the review stage may result in a new design solution that starts the next iteration using a revised set of parameters and parameter values.

At the outset, the client and designer decide which control pa-



Figure 7 Our vision of new BEM tools that speed up the design exploration process, allowing several design iterations during a single client meeting

rameters are most important and agree their relative values. In the case of a new speculative office tower, the key parameters might be the net lettable area and the comfort of circulation measured in vertical transportation passenger waiting time. Moving on to the control stage of the process, the designer then selects a solution from a database of generic ones: a tall, slender tower, a shorter tower with an atrium and so on. Using an intuitive interface, the designer also inputs the control parameters for the selected design solution. Next the toolkit automatically calculates, for example, how many lifts and how many service floors a tower of a certain level of quality needs. Finally each solution is published as a photorealistic 3D visualisation together with a plan, a table of quantitative data and false colour representations of physical data, such as comfort. (An example is shown on page 12.) The easy-to-understand output is reviewed collaboratively by all parties, who can decide there and then to fine tune the control parameters and/or add new ones to create the next option.

For this process to work, the BEM tools will need to provide a high quality of visual output to allow clients intuitively to understand the merits and drawbacks of each design option. Even more importantly, the process must be fast \neg allowing a number of iterations to be explored during the course of a meeting or workshop. To achieve the necessary speed, the tools will use generic data (gleaned from a database of the engineer's experience of designing similar buildings) and consequently the solutions they generate will be accurate but not totally precise. It is important that all parties understand this and realise – for example – that cost and energy use data is indicative rather than absolute.

These collaborative BEM tools will not replace the engineer or any other member of the design team. Instead, the tools will enable designers to provide computational leadership, using their professional intuition to quickly broaden the design space in search of novel aspects while drawing together everyone's expertise and ideas.

How it's done now

Many of the existing BEM tools that enhance collaboration between architects and engineers during design development are simply too slow and cumbersome to use in a design inception meeting; the designer needs to spend hours just inputting the precise parameters. Designers, therefore, continue to use long-established methods to demonstrate their credentials and explore ideas at initial client meetings and design workshops. The designer usually produces a few options before the initial meeting. These are based on the key parameters of the client's initial brief and are influenced by the designer's previous projects of the same type, which may be used as case studies to explain the thinking behind the design options. Presented either as a slide show or as a series of presentation boards, the proposals illustrate the designer's expertise while stimulating debate and helping the client to pinpoint their key parameters and design preferences.

Slide shows and presentation boards, however, are not interactive media; they can't provide instant responses to the client's questions and ideas. In most situations, the designer goes away to produce further options based on the new parameters, sequentially taking on board comments from, say, the architect, the developer and the developer's advisors. As a result, the client may have to wait two weeks or more before the new options are published and ready for review. And, therefore, only a few iterations are possible within the project's time constraints.

Because this process is so cumbersome, important strategic design decisions are sometimes left unmade – or are made by default based on past experience rather than the specific opportunities offered by the project. Less important decisions may take priority while more important decisions may drop off the agenda. And sometimes an inappropriate amount of time is spent making decisions.

Alternatively, at the initial meeting the designer might quickly sketch some new options on the spot. But while these intuitive "back of the envelope" drawings allow the client to feel involved in the design process, these new design solutions are untested and unquantified, and may send the design process in the wrong direction.





First steps towards the vision

New BEM tools are available today to help us move closer to the vision of rapid, intuitive design collaboration. During collaboration with an architect on a competition scheme or exploration of a project brief with a developer or end user client, these tools help to stimulate debate and allow rapid production of design solutions. While not yet offering parallel analysis of multiple parameters or rapid publication of both high quality pictures and numbers, the tools are helping designers to move beyond business-as-usual towards uniting engineering and client constraints to find the best solution.

Arizona State University's (ASU) Decision Theater supports round-table discussions among large groups of stakeholders. The Theater is lined with seven rear-projection screens, these can be configured to each display a live single-disciplinary model that simulates 'what if' design scenarios. Because the visual output is abstract and the models are loosely coupled , the process can be very rapid. This is important when a discussion involves a large group of diverse stakeholders, most of whom have little interest in many of the design parameters. ASU's researchers have observed empirically that seven seconds is the longest time that a group of this kind can wait for feedback before the flow of their discussion is disrupted. Longer processing times are only accepted by design teams, who are usually more engaged in the totality of the design.

Arup in New York has developed an energy toolkit that - in as little as half an hour - can produce accurate data on the energy implications of a particular design option using a database of generic solutions. Providing information on, say, the energy saved by installing solar panels in various locations or by using different cooling systems, this toolkit helps the designer and client to understand the energy implications of various possible solutions.

The toolkit has already contributed to two competition-winning design proposals - for a local government building in Italy and an extension doubling the size of a landmark museum building. In the case of the museum, the toolkit empowered the winning team to do something that none of their competitors had tried: locating the extension in front of the existing main facade but mainly underground. The toolkit allowed the team to demonstrate that this innovative solution was not just possible but also, counter-intuitively, affordable - offering vital reassurance to the commissioning client. A lesson learned is that the client now wrongly expects the energy performance of the detailed design to match the design inception results. When BEM tools are used at design inception, it is important to manage the client's expectations. It needs to be made clear that the solutions are accurate (resulting from a robust model) but are not precise (being based on generic rather than specific data).

To aid better specification of building services for laboratories, Arup has created a toolkit based on a database of power consumption per unit floor area of various general and specialist types of lab equipment. (In a laboratory building, research equipment is responsible for most of the energy requirements.) The database allows rapid calculation of a building's power and cooling requirements to meet the client's goals of comfort (important if they want to attract and retain world-class researchers) and efficiency.

Computational design optimisation (CDO) – championed by Arup Lighting in London – allows design teams and developers to balance key parameters through the exploration of numerous design options. For example, Arup used CDO for the design of a residential, commercial, retail and leisure development in Lewisham, south London. Two conflicting project aims were to maximise the total volume of the development while minimising the negative effects of the scheme on local residents. With CDO, Arup Lighting was able to increase the built volume of the scheme by 20% compared to usual standards while maintaining good natural lighting in the development area and for the surrounding buildings.

At present, however, CDO is an offline tool and cannot be used to provide on-the-spot answers as part of an interactive design discussion. And, while it is able to analyse client parameters such as development volume, it can't analyse multiple parameters from different disciplines.

In the following pages, three BEM tools are outlined that enhance collaboration at design inception. The tall building simulation tool streamlines the stages of control, compute and publish in the design exploration process. The bowl geometry simulation tool allows the architect and client collaboratively to appraise multiple design iterations. And the facade optioneering toolkit allows the engineer and architect to work together to find the best option.



Figure 9 Consultant and designer collaborate at design inception

Consultant and designer collaboration: tall building simulation tool

This case study examines part of the design exploration process - from control to publish - as shown (with one iteration) in the diagram.

The tall building simulation tool is a prototype to streamline the control, compute and publish stages of the design process. It was developed by Arup in Amsterdam in collaboration with Aedas, Davis Langdon and Hilson Moran. The tool allows the design team, the cost consultant and the developer's real estate advisors to understand the commercial and operational implications of the various options when setting the design brief for a new tower. The factors to be considered may include the size of the floorplates, the number of storeys, and its sustainability. Using parametric and associative modelling, the tool offers multi-dimensional simulation to demonstrate how changes to the plan, height and geometry of the tower affect performance, marketability, lifecycle costs and energy use. The easy-to-interpret output is published within a few seconds as a threedimensional geometric model, a plan and sets of engineering and cost data on a dashboard.

The tool empowers collaboration between the design and consultant teams, allowing everyone to understand the design opportunities and constraints (such as the desirability of the floorplate and the potential for natural lighting and ventilation) while offering proof of concept for the chosen solution.



Figure 10 The tall building simulation tool can help establish the "mix" of a building at initial design meetings. The diagram shows one iteration of the design exploration process, but many iterations may be developed during the course of a meeting

enue objectives, spectators' desire to feel part of the action and broadcasters' technical needs, as well as a range of life safety, building and football regulations.

The bowl geometry simulation tool (developed by J Parrish, formerly at Arup Sport in London) allows the architect and client to review subtly different architectural design iterations to compare, for example, lines of sight.

Getting the geometry of the stadium bowl right is fundamental to creating a good experience for spectators. The aim is to bring them as close as possible to the action while maintaining good sightlines. Achieving this requires a careful balancing act because these aims can conflict both with each other and with other aspects of the brief. For example, increasing the space between seating rows creates better sightlines but draws spectators further away from the field, while resulting in a larger stadium and increased construction costs.

The bowl geometry simulation tool allows quick and cost effective production of many subtly different bowl geometry models – based on default stadium data – at the design inception stage. As well as seating arrangements and sightlines, the models include circulation spaces and exit routes, concessions, and bowl and roof structures, enabling the client intuitively to explore and understand all the options.



Figure 11 Architect and client develop designs collaboratively

Willish Certenticity different design Vision 31 of 31 of a rectangular statium design Overenticity different design Vision 31 of 31 of a rectangular statium design Overenticity different design Vision 31 of 31 of a rectangular statium design

Figure 12 Designing a stadium bowl involves a lot of fine tuning. The bowl geometry simulation tool can compute and publish 3D illustrations of slightly different options within minutes.

geometry simulation tool This case study examines the publishing stage of the design exploration process, as shown (with several iterations) in the

diagram. A new soccer stadium has to meet a complex set of requirements. These include the client's construction budget and rev-

Architect and client collaboration: bowl



Figure 13 The facade optioneering tool unites parametric modelling and analysis software to allow multidisciplinary exploration of facade design options. The results are published as 3D geometric models alongside a range of building data.



Figure 14 Engineer and consultant enjoy the collaborative review of comparable options

Engineer and consultant collaboration: facade optioneering toolkit

This case study shows how the published output of several design iterations can be compared side by side. Providing all the information necessary for cross-disciplinary discussion, the facade optioneering toolkit allows project teams to make informed decisions during the early stages of the facade design by comparing design iterations side by side.

The toolkit was developed by Arup in Sydney in collaboration with RMIT University in Melbourne. Using the DesignLink software development kit as a platform, the tool unites parametric modelling and analysis software to allow easy exploration and optimisation of facade design options. The toolkit helps the engineer and consultant to work creatively together. This contrasts with the bowl geometry simulation tool outlined in the previous case study, which enhances collaboration between the client and architect.

As the geometric parameters are changed, the toolkit updates the geometric model and provides multidisciplinary analyses. A graphical interface displays three different design options as 3D geometric models, together with structural, thermal, lighting and carbon emissions data and net lettable area and cost figures for each option.



Figure 15 An example of BEM used in the design of a work of art. The screenshot is from a project bespoke tool developed by Tristan Simmonds (formerly at Arup) for sculptor Antony Gormley to express his artistic intuition for Space Station, a massive sculpture created for his 2007 show at London's Hayward Gallery.

Conclusion

The case studies show how BEM enhances collaboration between the designer and other consultants and – in the case of the bowl geometry simulation tool – the architect and the client.

Arup's mission, however, is to "shape a better world", which involves more than design optimisation; it requires designers, including engineers, first exploring many solutions, as well as those that may not immediately be obvious. This mission is a radical one, and achieving it depends on building creative, collaborative partnerships with clients.

To make the most of our clients' aspirations, ideas and insights, we need to empower them with an intuitive representation of relevant parameters. BEM offers excellent opportunities to enhance clients' involvement at design inception meetings by making it simple and cost effective to explore a wide variety of design options.

REFRENCES

[1] Bittermann, M.S. "Intelligent Design Objects: a cognitive approach for performance-based design (IDO)", doctoral dissertation, TUDelft, (2009)

[2] | Luebkeman, C., Shea, K., "CDO: Computational Design + optimization in building practice". The Arup Journal, (2005)

[3] Luebkeman, C., Simondetti, A., "Practice 2006: toolkit 2020" in Intelligent

Computing in Engineering and Architecture, 13th EG-ICE Workshop 2006, Ascona, Springer. (2006)

[4] Shea, K., Sedgwick, A., and Antonutto, G., "Multicriteria Optimization of

Paneled Building Envelopes Using Ant Colony Optimization", Proceedings of the 13th EGICE

Workshop on Intelligent Computing in Engineering and Architecture, Edited by I. F. C.

Smith, Ascona, Switzerland, June 25-30, pp 627-636. (2006).

[5] Roberts, C "Spatially-linked Integrated Resource Management (IRM)" Conference presentation, New Partners for Smart Growth 2010, Seattle, W, (2010)

[6] Shea, K., Aish, R., and Gourtovaia, M. "Towards Integrated Performance-driven Generative Design Tools." In Automation in Construction 14(2), pp. 253-264. (2005)

[7] Flager, F. and J. Haymaker, A Comparison of Multidisciplinary Design, Analysis and Optimization

Processes in the Building Construction and Aerospace Industries, in 24th International

Conference on Information Technology in Construction, I. Smith, Editor. Maribor, Slovenia.

p. 625-630 (2007)

Design a medical application for Android platform using model-driven development approach

J. Yepes¹, L. Cobaleda², J. Villa D¹, J. Aedo¹ ¹ARTICA, Microelectronic and Control Research Group ²ARTICA, Software Engineering Research Group University of Antioquia Medellín, Colombia

Abstract - Since the complexity of embedded systems has grown significantly, it has been necessary increase the abstraction level. For that reason, we propose a method to design a system for medical device interoperability, called SIMMIT (System Integration Medical Monitoring and Interoperability for Tele-care) based on a Model-Driven Development approach. This work presents a strategy and the tools to design a software application for embedded systems supported by functional and non-functional requirements. This approach starts with a system specification report, which describes both the structure and the functionality of the system. Thus, with this document functional system model is creating whose application specification is independent from implementation details. On the next development stage, the model evolves by adding to it features of a specific platform. Finally, in this phase the Java code generation for Android platform is automatically done from Unified Modeling Language diagrams. This application is running on Android operating system in the OMAP3530 processor.

Keywords: Android, Embedded System, Methodology, Model-Driven Development.

1 Introduction

The rise of complexity in embedded systems applications at the present time has been considerable. The new applications require accomplishing demands of quality factors such as: good performance, reliability, security, portability, interoperability, robustness, scalability and low power. Thus, these factors introduce an additional effort in order to fulfill the development of reliable systems, with development time and cost reasonable.

In order to face this problem, Model-Driven Development (MDD) approach is introduced as design software methodology. This approach focuses on the modeling of functionality of the system without consider the technology in which it will be implemented in order to manage the complexity.

Typically MDD methodology defines four models types during the development process: Computation Independent Model (CIM), Platform Independent Model (PIM), Platform Specific Model (PSM) and Platform Specific Implementation (PSI). The first model defines the system requirements. The second one defines the functional system model. The third one defines the software model with details of the platform and the last model correspond to the application generated in a target language, that implement the functionality defined in the previous models [1].

On the other hand, Unified Modeling Language (UML) has been adopted as the de facto standard modeling language in the MDD approach. For Embedded System (ES) design, UML provides a rich set of constructs to support the modeling of system functionality, behavior and structure. Some limitations of language related to ES design such as the modeling of realtime requirement have been addressed with the use of UML profiles that permit extend the expressiveness of the such language [2].

In this work, we applied a software development methodology for embedded systems. We used a MDD approach to design a medical device for health monitoring in critical situations, called SIMMIT. The application is running on Android [3] operating system and the OMAP3530 platform.

This article is organized as follows. Chapter 2 presents some related work. Chapter 3 describes the OMAP3530 platform and tools, Chapter 4 provides an overview of the methodology used. Chapter 5 describes the case study. Finally the conclusions and future work are presented in Chapter 6.

2 RELATED WORK

In [4], proposed a methodology based on MDD approach to the development of mobile applications implemented in a cell phone. This work has developed a graphical modeling language specific to mobile applications, and coming up with a generic algorithm for the conversion of this graphical model into code. The main effort was putted in the design of interaction techniques, which will allow creating mobile applications easily. Unlike this work, authors presents a list of guidelines for modeling mobile applications, regardless platforms or the context.

Another work [5] proposes model-driven development of mobile personal healthcare applications. The authors

developed an approach in order to modeling care plans for chronic disease, using two domain-specific visualizing languages (DSVLs). The first allows healthcare providers to model complex care plans, health activities, performance measurements and sub-care plans. The second DSVL describes a mobile device interface for the care plan. A code generator synthesizes mobile device implementation of this care plan application. Unlike these works, in our application communicates with a medicals devices which is managed from the communication and processing of physiological variables.

3 PLATFORM AND TOOLS DESCRIPTION

Android mobile platform is supported by Texas Instrument's OMAP3530 processor. The OMAP generation of high-performance, applications processors are based on the enhanced device architecture and are integrated on TI's advanced 45-nm process technology. This architecture is designed to provide best in class ARM and Graphics performance while delivering low power consumption. This balance of performance and power allow the device to support a huge variety of multimedia applications [6].

The OMAP3530 integrates a GPP (General Purpose Processor) ARM Cortex TM-A8 @600MHz, a DSP (digital signal processor) TMS320C64x @430MHz plus a graphics accelerator 2D and 3D PowerVR SGX 530. The GPP controls all hardware resources using a generic operating system like Linux, Windows CE or, in this case, Android. The DSP acts as coprocessor of GPP. It also integrates various peripherals and interfaces to connect the different types of external devices [7].

The tools using are the following:

3.1 Eclipse

Eclipse is a software development kit (SDK) [8] consists of the Eclipse Platform, and Java development tools and a multi-language software environment composed by an integrated development environment (IDE) with an extensible plug-in system. Eclipse can be used to develop applications in various programming languages including Ada, C, C++, Java, Perl and Python. Development environments include the Eclipse Java development tools (JDT) for Java, Eclipse CDT for C/C++, among others.

3.2 Android Software Development Kit

The SDK (Software Development Kit) for Android is officially supported with Eclipse Plug-ADT (Android Development Tools plugin) and includes a set of development tools. This SDK including: debugger, libraries, phone virtual machine, documentation, sample code and tutorials [9].

3.3 IBM Rational Rhapsody

Rational Rhapsody is a modeling environment based on UML [10]. It was primarily designed to accelerate development and reduce costs, improve quality and managing complexity, through visualization of the models. Moreover IBM Rational Rhapsody helps to maintain consistency across the development life cycle to facilitate agility in response to requirements.

4 METHODOLOGY

The current trend is to use the methodologies in high levels of abstraction in the early stages of design, in order that the description of the system can be quickly and completely [11]. MDD (Model Driven Development) approaches have been proposed as a clear methodology for developing embedded systems [12], because provides the ability to streamline, standardize and replicate design practices by allowing a completely independent functional specification to implementation.

MDA (Model Driven Architecture) is one of the most popular MDD approaches proposed by OMG (Object Management Group) for software development, based on models at different levels of abstraction. In the MDA approach, a system is modeled using a platform independent model (PIM), which is transformed into a platform specific model (PSM) using design patterns. The languages used to express these models are defined by means of meta-models that are able to define abstract and concrete syntax and operational semantics of modeling languages.

The main objective of MDA is to separate the functionality of the architecture to build flexible systems, which do not depend on a hardware platform, or a specific software architecture, i.e. the application PSM can easily generated for different platforms, using the same PIM.

This work defines a process for generating a concurrent application Code which will run on the OMAP3530 hardware platform. From a PIM application, where define its functionality, is necessary generate a PSM model that considers the low level details necessary for the application to run properly on the target platform. Such details are closely tied to the hardware platform, in our case it is necessary to consider: Android operating system. The process consists of 4 stages:

A. Development of requirements: In this stage it is necessary to select the requirements that will guide the subsequent phases; i.e. the actors, their responsibilities, some functional features and constraints.

In the case study, a set of requirements is identified in the first iteration. Some of them were selected to establish the base architecture. These are related to the critical or main functions of the system. Although this application development could require more iteration, so the design is projected to be scalable.

B. Application Model: This stage consists in develop a PIM from the requirements document. In this phase is necessary identify basic functional features of the system, in that order the uses case and general classes are identifying from requirements document and generated using a modeling environment based on UML Rational Rhapsody.

C. Application Specific Model: This stage describes the functionality of the application, along with non-functional features of the system is obtained by enriching the PIM model by means of a transformation process. The transformation PIM to PSM is probably the most common focus in MDA. This is traditionally the mapping from the essential analysis model to the platform-specific design model. In this part the PIM model is transformed semi-automatically using meta-modeling transformation techniques to generate output for the design phase according to the platform to implement. Moreover, design patterns are applied manually in a PIM which transform in PSM.

D. Code Generation: With the PSM fully developed, you can feed a tool that transforms the PSM model described in UML to code c, c + + or java. The code is generated with Rational Rhapsody tool and subsequently the particular modifications are introduced in the code directly. It is important to emphasize that it is possible to accomplish the iterative development process due to the Rational Rhapsody tool capability for keeping the consistency between models and code.

E. Running the application: Once the code generation process is complete, proceed to test the application on either a virtual or on the OMAP3530 platform using Android operating system.

5 **CASE STUDY**

The proposed methodology is evaluated by developing a prototype called SIMMIT (System Integration Medical Monitoring and Interoperability for Telecare). The objective is to integrate and transmit towards a medical center the medical record information of a patient in emergency state within a medical assistance vehicle or in a remote station.

The signals derived from monitoring equipment (ECG, heart rate, respiratory rate, oxygen saturation and blood pressure) and recording the findings in the patient should be integrated and appropriate to the patient's electronic medical records in a standard format for then sent to a remote location via a

communication network when the medical staff required. In Figure 1 illustrates the SIMMIT and its environment.



Figure 1. SIMMIT Diagram.

In order to implement the case study was carried out the methodology proposed in paragraph 4. First they drew up a requirements document, from there took place a model of the application. In Figure 2 illustrate a use-case diagram of SIMMT.



With the use cases and after an object analysis performed a conceptual class diagram in Figure 3 illustrates a portion of the class diagram.

UiMonitorina ጓ ecgDatos:byte[] UiMonitoring() iniciar():void parar():void salir():void evStart() ACT-002 APH evStop() evGuardar() procesar(): void CtrlCapture Processing ecgDatosProc:int[] fc:int hrPort ecgDatos:byte[] Processing() caPort CtrlCapture() escala():void capturaLista():boolean procesar(data:byte[]):int[] apturar():boolean anturarECG0.void

Figure 3. Class diagram

In order to verify the PIM, was made a diagram of state machine, as shown in Figure 4.

Once the PIM is checked, proceed to perform the PSM. This model must take into account the nature of the application described in the PIM, and the details to implement it on a specific platform. In our case we use the OMAP3530 platform and Android operating system. The PSM should include design patterns, as these provide unified solutions to recurring software problems. In Figure 5 illustrates how the class diagram has been modified in order to add the features of Android and the observer pattern [13].



Figure 4. Statechart diagram

The next phase is to generate code from the PSM with the help of the tools IBM Rational Rhapsody and Eclipse, as shown in figure 6. The final step is to run the application designed with the MDD methodology on the OMAP3530 platform as is showing in Figure 7. If everything is correct it has completed the design process, otherwise is necessary proceeds with an iterative process that consists of refining the model PSM and generates code again to fulfill all the system requirements.



Figure 5. PSM class diagram

```
//## class UiMonitoring
public class UiMonitoring extends Activity implements Observer {
   double[] ecgData = new double[30];
                                            //## attribute ecgData
   // Constructors
   //## operation UiMonitoring()
   public UiMonitoring() {
       //#[ operation UiMonitoring()
       enableUi(false);
        //#]
   }
   /**
    * @param savedInstanceState
   */
   //## operation onCreate(Bundle)
   public void onCreate(Bundle savedInstanceState) {
        //#[ operation onCreate(Bundle)
        super.onCreate(savedInstanceState);
        setContentView(R.layout.monitoring);
       addButtonsListeners();
```

Figure 6. Code snippet



6 Conclusions

In this work a software application for medical health care called SIMMIT has been developed using a MDD approach with success. This methodology starts with UML modeling application functionality regardless of implementation details, via an intermediate model that adds these details up to the executable code generation for OMAP3530 platform with Android 2.2 operative system.

The methodology aims to improve the productivity of embedded system design, rescuing the benefits of reusability, scalability, maintainability and modularity of system components and it provides the ability to separate the functionality of the application of the implementation details by means of models.

ACKNOWLEDGMENT

We would like to express our thanks to the Excellence research Center, ARTICA and to the members of Microelectronic research group and Software Engineering Research Group from Antioquia University.

7 References

- [1] Document. MDA Guide. 2003. http://www.omg.org/cgibin/doc?omg/03-06-01.pdf [Online. Cited: Octuber 22-10-2010].
- [2] P. Green, "Uml as a framework for combining different models of computation," in UML for SOC Design, G. Martin and W. Müller, Eds. Springer US, 2005, pp. 37– 62.
- [3] "What is Android?, "Available: http://developer.android.com/guide/basics/what-isandroid.html [Online].
- [4] F. Balagas, H. Husmann, "Model-Driven Development of Mobile Applications ", In Proceedings of the 2008 23rd IEEE/ACM International Conference on Automated

Software Engineering (ASE '08). IEEE Computer Society, Washington, DC, USA, 509-512.

- [5] A. Khambati, J. Grundy, J. Warren, and J. Hosking. "Model-Driven Development of Mobile Personal Health Care Applications". In Proceedings of the 2008 23rd IEEE/ACM International Conference on Automated Software Engineering (ASE '08). IEEE Computer Society, Washington, DC, USA, 467-470.
- [6] OMAP3530 Architecture. http://focus.ti.com/docs/prod/folders/print/omap3530.ht ml (Available: February, 2012).
- [7] OMAP35x Applications Processor Texas Instruments OMAP Family of Products Technical Reference Manual, October 2009.
- [8] Eclipse software development kit. Available: http://wiki.eclipse.org/Main_Page [Online].
- [9] Android Software Development Kit. Available: http://developer.android.com/guide/developing/tools/ind ex.html [Online].
- [10] IBM Rational Rhapsody. Available: www.telelogic.com/products/rhapsody/index.cfm. [Online].
- [11] A. Sangiovanni-Vincentelli. Quo Vadis SLD: Reasoning about Trends and Challenges of System-Level Design. Proceedings of the IEEE, 95(3):467-506, March 2007.
- [12]L. Bondé, C. Dumoulin, J. Dekeyser. "Metamodels and MDA Transformations for Embedded Systems". In: Forum on Design Languages (FDL'04), Lille, 2004.
- [13]B. Powel Douglas. (2002) Real-Time Design Patterns: Robust Scalable Architecture for Real-Time Systems. Boston U.S.A Addison-Wesley.

Muscle Force Exploration through Simulation for Passenger Seat Design

Sangho Park, Piao Jinghu, Murali Subramaniyam, Sun Junfeng, Taesu Yim

Department of Mechanical Design Engineering, Chungnam National University, 79, Daehangno, Yuseonggu, Daejeon 305-764, South Korea

Abstract - Along with recent advancement, many automobile industries have strongly encouraged research on their product design, structure, consumer comfort and ergonomics and so on. In particular, consumer satisfaction or human comfort is playing an important role including seat comfort, driving posture, visibility, interior space, etc., In general, sitting in an automobile seat for prolonged period can cause back pain or worsen an existing back or neck problem. This study performed the simulation to explore muscle force during prolonged period of sitting for passenger seat design. The trapezius muscle forces have been measured from the simulation and compared for the different backrest inclination angle and sitting period.

Keywords: — sitting posture, back pain, muscle force, ADAMS/LifeMOD

1 Introduction

The automobile is not only as a transportation for humans, but also have a very close relationship with humans as the entertainment or work tools, recently, depend on rapidly development of the sciences, the automobile's performance has been greatly improved, and with the development of economic the consumer culture also has a qualitative change, the consumers pay more attentions to the comfort of driving, therefore, how to improve the comfort of the automobile become a hot topic in the field of automobile industry.

The comfort and fatigue of driving depend on the road, driving speed and time, driving posture or sitting position and the cushion device of the seat, as the related research on the driving posture, Rebiffe,1969[1] had a undertook study about when the each joint of our body in some angle is the best driving posture, obtained a relationship of the driving seat and pedal that applies to various height driver, although the results of the seat angle range is too big, and the angle is actually a 2D parameters rather than the 3D parameters that actual seat design to be used, but the results are very valuable for the seat designers as some design basis. Verrirst, 1986[2] introduced a kind of adjustable experimental device that can measure the variable parameters driving of posture, and

Schneider's,1979[3] studies tell us, driving posture, the position between the seat and steering wheel is a complex interactional influence for the comfort. In South Korea, for the automobile seat that suitable for Koreans shape also have some related research Se Jin Park,2000[4],using statistical method to derive the standard of the driving posture Sung Jun Park,2006[5], as described above, although there are many researches for the driving posture and fatigue, But this field is currently in a chaotic state, therefore, is very necessary to study the factors that affect the riding or driving comfort , especially in the premise of ensuring the driver or passengers' safety, how can to provide the comfortable riding feeling to the driver and passengers, and lifting the fatigue of riding the automobile. It is still need to be very large and complex research to solve these problems.

Recently, many software applications have been developed for impact simulation, biomechanical analysis, movement simulation and surgical planning. The software enables users to perform human body modelling and interaction with environment where the human motion and muscle forces can be simulated.

The created human body may be combined with any type of physical environment or system for full dynamic interaction. The results of the simulation are the human motion, forces exerted by the muscles, and the stresses or strains at the desired location of the human body [6].

2 Research method

The presented model can be applied to understand the complex spine biomechanics and clinically important analysis such as contact forces between each vertebra and wheelchair model, load acting on the intervertebral disc joints, corresponding angles between vertebrae in the seated position and tension in the spine muscles. These results aid clinicians to develop mechanical design of back support, such as placing conventional pillows and towels at appropriate positions which can be an effective and convenient alternative to expensive special seating [7].

While studying the muscle force when driving, a manautomobile system simulation model is needed. In this study, a dynamic model created which contains passenger model and a passenger's seat. The human modeled with South Korean characteristics, which contains skeletons, muscle, ligaments and joints. With this dynamic model simulated the process of passenger's driving. The simulation has been performed to explore the muscle force, when the backrest at different angles.

3 Discussion

With the simulation results, comparative analyses have been performed. Through analyzes, the relationship between the comfort level and the backrest inclination angle have been found. We wish the values and results will be useful for the bus passenger's seat design.

4 Acknowledgment

This research was supported by Basic Science Research Program through the National Research Foundation of Korea (NRF) funded by the Ministry of Education, Science and Technology (2012-0003090).

5 References

[1] Rebiffed, R. "The driving seat. Its adaptation to functional and anthropolatric requirements",1969.

[2] VeriTest, J.P. "Driving posture and comfort drivers of road vehicles (passenger cars, commercial vehicles, and heavy lorries)", Recherche Transports Securitas, Special Issue, Institute National de Recherche sur les Transports et leur Securite, France, pp. 38–44,1986.

[3] Schneider, L.W ,Anderson, C.K. and Olson, P.L. "Driver anthmpometry and vehicle design characteristics related to seat positions selected under driving and nondriving conditions", SAE paper No.790384.

[4] Se Jin Park, M.K Jeoung, K.S Kwon, S.W Kim. 2000, "A Study on Evaluation of Fatigue in Vehicle", Journal of the Ergonomics Society of Korea, pp.27-30.

[5] Sungjoon Park. "Estimation of Driver's Standard Postures by a Multivariate Analysis Method", Journal of the Ergonomics Society of Korea, pp.27-33,2006.

[6] Lee, K.W. CAD System for Human-Centered DesignComputer-Aided Design & Applications, Vol. 3, No. 5,2006, pp 615-628.

[7] Tay Shih Kwang. Ian Gibson, and Bhat Nikhil Jagdish "Detailed Spine Modeling with LifeMODTM".

SESSION ALGORITHMS AND NOVEL APPLICATIONS

Chair(s)

TBA

Generating an Informed Virtual Geographic Environment through Cell Merging in order to Geosimulate the Propagation of Zoonoses

Mondher Bouden¹ and Bernard Moulin²

Department of Computer Sciences and Software Engineering, Laval University, Quebec, Canada ¹Mondher.Bouden@ift.ulaval.ca, ²Bernard.Moulin@ift.ulaval.ca

Abstract - Zoonoses (infectious diseases transmitted from insects to animal and to humans) such as Lyme disease is a concern for public health authorities, especially after their proliferation due to global warming. Notwithstanding the obvious influence of geographic features (as for example the suitability of geographic areas on the species' biological processes such as survival, feeding and breeding), they are not integrated in currently available zoonosis simulations based on mathematical compartment models which specify the characteristics of the transitions between the different stages (e.g. eggs, larvae, adults, inflected, etc.) of the involved species. Hence, there is a need for a software to simulate such a propagation taking into account the geographic characteristics of the landscape and the spatial behaviors (and interactions) of the involved species. In this context, we propose a novel approach using an informed Virtual Geographic Environment (VGE) composed of a set of cells in which the transitions of the different biological stages of the involved populations, as well as their interactions can be simulated plausibly. In order to reflect in the VGE the landscape's influence on the biological phenomena, we automatically generate the VGE cells from vector-based land-cover data provided by a Geographic Information System. The cells are obtained by a stepwise aggregation of landcover polygons selected according to biologically relevant qualitative criteria such as the survival suitability of a cell for a given species. For geosimulation purposes one critical issue, discussed in the paper, is to significantly reduce the huge number of land-cover polygons that are associated with the studied geographic areas. We propose a threshold-based merging algorithm which iteratively reduces the number of polygons while generating a spatial subdivision composed of cells with the maximal size and suitability for a given biological phenomenon. Our approach also enhances each cell with qualitative information about the relative geographic orientation of its neighbors. The resulting informed VGE provides the foundation for the simulation of a variety of spatial-temporal phenomena such as the migration of birds importing infected insects (i.e. juvenile ticks in the case of Lyme disease) and the survival of tick colonies in suitable areas.

Keywords: Informed Virtual Geographic Environment, Merging of Cells, Multi-Level GeoSimulation, Spread of Infectious Diseases.

1 Introduction

The expansion of some zoonoses such as the West Nile Virus (WNV) or Lyme disease led public health authorities to develop monitoring systems [1] in order to better understand the epidemiology of the disease and the level of risk it can represent for human populations. However, these monitoring activities cannot be used to forecast the probable propagation of a zoonosis over a territory. There is a need for tools to simulate such a propagation taking into account the geographic characteristics of the landscape and the spatial behaviors and interactions of the involved species. In this context, we are interested in using modeling and computer simulation to develop a decision support tool which can help public health policy makers to better understand the zoonosis propagation phenomena and to explore the possible effects of intervention scenarios at appropriate time and places and at the appropriate level of expected risk.

Besides, several approaches have been proposed to model and simulate the spread of zoonoses. However, these approaches such as mathematical modeling, cellular automata and traditional multi-agent systems have some weakness when trying to model and simulate the influence of geographic and climatic features on the disease spread and the spatialtemporal interactions of various kinds of actors (e.g. mosquitoes, ticks, birds, mammals, etc.). Indeed, the simulation based on mathematical models that generally uses compartment models and differential equations [2] does not take into account the characteristics of the geographical space in which populations operate, except in limited cases such as patchy models [3] which use an abstracted representation of space that is not based on data provided by Geographic Information Systems (GIS). In spite of the fact that a simulation based on cellular automata models the evolution of the spatial characteristics of a geographic area involved in the disease, it does not represent individuals and their mobility [4]. On the other hand, agent-based simulations of epidemics represent the disease vectors (e.g. animals) as agents, but usually do not take advantage of data provided by GIS in order to properly locate agents in the geographic space [5]. Moreover, given the limitations of computational resources of computers and the lack of data, we cannot plausibly represent each individual by an agent, especially if we have to simulate a population composed of millions or even billions of individuals involved in zoonoses, as for example mosquito populations and tick populations transmitting the WNV and Lyme disease respectively.

In this context, we propose to use a multi-level population-based geosimulation approach [6] to remedy the shortcomings of current methods. We acquired some experience with the development of WNV-MAGS System [7], a tool allowing public health decision makers to assess several intervention scenarios in order to understand and estimate the magnitude of the evolution of the WNV in a large territory. Furthermore, we are currently developing a generic solution (Zoonosis-MAGS) to be applied to a variety of zoonoses such as Lyme disease, with the aim to create realistic simulations at different levels of granularity. To develop such populationbased geosimulations we introduced a new theoretical model (called MASTIM: Multi-Actor Spatio-Temporal Interaction Model) which can be used to simulate the interactions of various types of actors, including those representing populations containing a huge number of individuals [8].

In this paper we examine the critical issue of how to accurately represent and generate the VGE in which these geosimulations can be carried out, taking into account that the studied territories are huge and that we need to consider the geographic characteristics that influence the biological cycles of the involved species (i.e. areas that are un/suitable to the survival and proliferation of insects) and their behaviors (e.g. feeding, displacements, migration, etc.). A good idea is to divide the VGE into basic cells in which the different stages of the involved populations can be simulated. Such a space subdivision can be generated using either raster or vector data structures. In a raster-based system, the environment is divided into uniform-sized cells. Such raster-based VGE are used by cellular automata [5]. However, such grids of cells are usually artificial and not related to the spatial characteristics of the studied phenomena. In a vector-based system, the environment is represented using geometrical primitives such as polygons (regions of space which are well defined using GIS data). Several kinds of polygons can be used to make a spatial subdivision of the VGE. The choice of the polygon types depends on the spatial characteristics that are important for the observation and analysis of the zoonosis propagation. In the case of WNV in which mosquitoes are spread over the territory, we used polygons representing either municipalities or census tracts, depending on the area of interest (municipalities are used to cover large areas such as the southern part of the province of Quebec whereas census tracks are used to characterize smaller areas such as the Ottawa metropolitan area). This administrative division fits the surveillance data which were available for the simulation. In the case of the simulation of the establishment of tick colonies and the spread of Lyme disease, it does not make sense to use such an administrative subdivision because the cells representing municipalities or census tracks are too large and have no biological significance considering that tick populations can only survive in grass areas at the edge of (or in) forests. We therefore thought of using a spatial subdivision based on land cover characteristics. Land-cover data may be provided in vector-format in which polygonal cells cover large territories (as for example in the Geobase database). However, another problem arises when it comes to simulate the propagation of zoonoses. Indeed, we need to compute for each cell the interactions of the involved species using some kind of compartment model. Given the huge number of land-cover cells it is impossible to carry out such computation for each time step in each individual cell of the land-cover subdivision. Consequently, we got the idea to merge land-cover cells

having similar characteristics with respect to the phenomenon to be simulated. Hence, we could get the largest polygons possible in order to get plausible simulations while minimizing the needed computations. The merging method that we propose is based on criteria selected by the user who models the phenomenon. A very common criterion is related to the suitability of the habitat which may affect the survival, feeding, and/or breeding behaviors of species. In order to generate an efficient VGE, the cells resulting from such a merging process need to have the largest possible sizes with respect to the selected criteria. In this paper we propose an approach to create such a 'biologically informed VGE'.

In Section 2 we present the GIS data that we use to generate the VGE. In Section 3 we present our new approach including the technique used to merge cells and to create the informed VGE. In Section 4 we discuss the usefulness of the resulting VGE for the geosimulation of the zoonosis propagation. Finally, Section 5 concludes the paper and evokes some future works.

2 Presentation of the GIS data

We use the land-cover shape files provided by the Geobase database (<u>www.geobase.ca</u>). The land-cover information is the result of vectorization of raster thematic data originating from classified Landsat 5 and Landsat 7 ortho-images, for agricultural and forest areas of Canada and for the Northern Territories. The land-cover data covers the totality of the Canadian territories and is divided in different regions using the index maps of the National Topographic System of Canada (NTS). Each region is identified by a unique number (e.g. 22M, 30K, etc.).



Figure 1. The six NTS regions of interest.

Since we are interested in the southern part of Quebec and more specifically in the regions of *Montérégie* and *Estrie*, we consider six regions (i.e. 21E, 21L, 31G, 31H, 31I and 31J) represented in Figure 1. For instance, region 31H contains 133 780 distinct polygons and its dimension is nearly 156 x 110 km. The vector data are distributed as surface features (polygons) that have descriptive attributes such as the land-cover class which are based on the EOSD (Earth Observation for Sustainable Development) Land-Cover Classification. Indeed, the Coverage type attribute takes its value between 0 and 233 and represents different categories of land-covers (e.g. 20 for water, 34 for urban area, 50 for shrub land, 220 for deciduous forest, etc.). It is also worth mentioning that some polygons may have complex shapes as shown in Figure 2.



Figure 2. Two complex polygons representing an urban area (left) and a forest (right).

3 A new approach to generate a biologically informed VGE

We begin this section by presenting an overview of our new approach. Then, we present the criterion used to merge cells. After that, we discuss the different steps of our approach including the preprocessing of GIS data and deleting holes, the progressive merging of cells and the creation of an informed VGE.

3.1 Overview of our new approach

Figure 3 presents an overview of our approach which is based on the merging of cells representing polygons provided by GIS data. Due to the large number of cells to be processed and because it is difficult to anticipate when to stop the merging process to get the largest polygons which satisfy the selected criterion, we propose an approach that is carried out in several steps. Indeed, in a first step we preprocess the GIS data by selecting the region of interest and computing the suitability of all polygons with respect to the selected criterion (see Sections 3.2 and 3.3). Then, we create plain polygons by removing their holes in order to reduce the number of basic cells (see Section 3.3). The result is stored in a new database which is updated as we go through the different processes. Besides, the most important step of our approach is the progressive merging of suitable polygons. We begin this process by selecting polygons having the best suitability (100%) according to the selected criterion. Then, we sort these polygons (the biggest is first processed) and get their neighbors. These neighbors are also sorted (the biggest and the best one according to the selected criterion is first processed) in order to try to merge them to the original polygons. Note that we choose different thresholds to stop the aggregation process if merging an additional neighbor decreases the suitability of the original polygon of more than 10%. In this way we preserve the suitability of the resulting polygons and thus we allow for a progressive merging process. We then iterate the merging process on a new set of polygons with a suitability value in the interval [90%, 100%] that we

call the 'absorption threshold interval' (ATI). We therefore continue to apply this process several times by progressively increasing the ATI for the processed polygons and by simultaneously reducing the threshold used to stop merging. Finally, we do the opposite by selecting polygons having the worst suitability (0%) according to the selected criterion and trying to merge them with other unsuitable polygons (see Section 3.4). As it is the case for suitable polygons, this process allows for generating homogeneous unsuitable areas (aggregation of unsuitable polygons) with maximal size. In addition, it further reduces the number of polygons used to generate the VGE. All the processes mentioned above are carried out until reaching a satisfactory result according to the user's appreciation.



Figure 3. Overview of our new approach.

Indeed, our approach based on heuristics in the form of merging rules using different criteria and thresholds, gives acceptable results because we apply a descending sort to the different processed polygons using area and suitability as order criteria. This sorting process allows the system to first merge the biggest polygons with the best valuation according to the selected criterion. Therefore, we ensure that the polygons selected by the merging algorithm are among the best candidates. Moreover, the last step of our approach is the creation of an informed VGE in which the zoonosis propagation can be simulated with computational efficiency (see Section 3.5).

3.2 Using suitable habitats

The biological phenomena that we study are closely tied to the characteristics of the landscape that is simulated by the VGE. Indeed, a zoonosis propagation greatly depends on the survival of the populations involved in the transmission of pathogens (i.e. virus, bacteria). For example, the capacity of tick populations to survive in suitable habitats (as for example 'sparse deciduous forests' according to Geobase terminology) is an important factor that influences the spread of Lyme disease. In fact, a species can settle in an area only if it is suitable to its survival. Therefore, for each polygon of the GIS data presented in Section 2 we propose to determine what we called a 'suitability degree' in order to estimate the quality of suitable habitats for a given species. To this end, we add a new attribute to the database table associated with each polygon. This attribute which represents the suitability degree of the polygon for a given species takes its values between 0 and 1. The value 0 represents a place that is unsuitable for the species' individuals and the value 1 represents a very suitable (100 % of suitability) place (see Table 1).

Table 1. Suitability degrees (SD) of habitats for ticks inrelation to land cover.

| Code | Туре | SD | Code | Туре | SD |
|------|---|------|------|--|------|
| 0 | No data | 0 | 100 | Herb | 0.2 |
| 10 | Unclassified | 0 | 101 | Tussock graminoid tundra | 0.2 |
| 11 | Cloud | 0 | 102 | Wet sedge | 0.2 |
| 12 | Shadow | 0 | 103 | Moist to dry non-tussock graminoid/dwarf shrub tundra | 0.2 |
| 20 | Water | 0 | 104 | Dry graminoid prostrate dwarf shrub tundra | 0.2 |
| 30 | Barren/Non-vegetated | 0 | 110 | Grassland | 0.2 |
| 31 | Snow/Ice | 0 | 120 | Cultivated Agricultural Land | 0.2 |
| 32 | Rock/Rubble | 0 | 121 | Annual Cropland | 0.2 |
| 33 | Exposed land | 0 | 122 | Perennial Cropland and Pasture | 0.2 |
| 34 | Developed | 0 | 200 | Forest/Tree classes | 0.8 |
| 35 | Sparsely vegetated bedrock | 0 | 210 | Coniferous Forest | 0.8 |
| 36 | Sparsely vegetated till-colluvium | 0 | 211 | Coniferous Dense | 0.6 |
| 37 | Bare soil with cryptogam crust - frost boils | 0 | 212 | Coniferous Open | 0.9 |
| 40 | Bryoids | 0 | 213 | Coniferous Sparse | 0.9 |
| 50 | Shrubland | 0,8 | 220 | Deciduous Forest | 0.85 |
| 51 | Shrub tall | 0.75 | 221 | Broadleaf Dense | 0.7 |
| 52 | Shrub low | 0.75 | 222 | Broadleaf Open | 1 |
| 53 | Prostrate dwarf shrub | 0.2 | 223 | Broadleaf Sparse | 1 |
| 80 | Wetland | 0.2 | 230 | Mixed Forest | 0.8 |
| 81 | Wetland - Treed | 0.2 | 231 | Mixedwood Dense | 0.65 |
| 82 | Wetland - Shrub | 0.2 | 232 | Mixedwood Open | 0.95 |
| 83 | Wetland - Herb | 0.2 | 233 | Mixedwood Sparse | 0.95 |

3.3 Preprocessing of GIS data and deleting holes

We tested our approach using the shape file of Region 31H (See Section 2). This file contains 133 780 polygons representing different kinds of land covers. Using Geomedia [9], we converted this file to an Access spatial database (read/write) in order to have the ability to modify it. For this conversion, we used the reference coordinate system of Canada (NAD83-CSRS). Then, we added and computed the suitability attribute for each polygon as mentioned in the previous section. We also developed a new application which is able to query the Access database and modify its attributes.

We noticed that a large number of polygons have some holes. These holes are considered as polygons by Geobase with their own land cover attribute. To reduce the number of polygons, we decided to remove these holes and to create plain polygons. We assumed that the filling of holes will not greatly reduce the suitability degree (SD) of the polygon, since the size of most holes is a small portion of the polygon size. To explain our approach, let us take the example of the polygon shown in Figure 4. Indeed, this polygon represents dense deciduous forest (SD = 0.7) and contains 20 holes. Some holes contain themselves other holes which brings the total number of holes to 35.



Figure 4. Creating a solid polygon from a polygon that contains 20 holes and 15 sub-holes.

Table 2 shows how we compute the new SD after deleting these holes. For example, there are 16 holes representing a dense mixed wood forest (SD = 0.65). The sum of areas of these holes is equal to 890025.2 m^2 . We compute an area proportion (AP) which is equal to the ratio between the sum of areas of the holes and the total area of the polygon without holes (890025.2/21368279.7 = 0.041651701). Then, we compute the suitability degree proportion (SDP) which is equal to the multiplication between AP and SD (0.041651701 x 0.65 = 0.027073606). The sum of all SDP represents the new suitability of the solid polygon. We should mention that the total number of polygons listed in Table 2 is 36, including the polygon containing holes. Moreover, to delete holes, the system uses a spatial query which finds all the polygons located in the interior of the polygon containing holes. Then, we merge them to obtain a filled polygon which will have the new calculated SD (0.69542 instead of 0.7).

Table 2. Computing the new suitability after deleting holes.

| Land Cover | Nb | Code | Area (m ²) | SD | AP | SDP |
|-----------------------------------|----|------|------------------------|------|-------------|-------------|
| Water | 1 | 20 | 96894.1 | 0 | 0.004534483 | 0 |
| Shrub tall | 3 | 51 | 35420.8 | 0.75 | 0.001657635 | 0.001243226 |
| Herb | 1 | 100 | 11284 | 0.2 | 0.000528072 | 0.000105614 |
| Annual Cropland | 1 | 121 | 60536.5 | 0.2 | 0.002833008 | 0.000566602 |
| Perennial Cropland and Pasture | 1 | 122 | 21789.3 | 0.2 | 0.001019703 | 0.000203941 |
| Coniferous Dense | 1 | 211 | 73049.7 | 0.6 | 0.003418605 | 0.002051163 |
| Broadleaf Dense | 2 | 221 | 19914816.3 | 0.7 | 0.931980327 | 0.652386229 |
| Broadleaf Open | 1 | 222 | 13574.2 | 1 | 0.00063525 | 0.00063525 |
| Mixedwood Dense | 16 | 231 | 890025.2 | 0.65 | 0.041651701 | 0.027073606 |
| Mixedwood Open | 9 | 232 | 250889.6 | 0.95 | 0.011741217 | 0.011154156 |
| | 36 | | 21368279.7 | | 1 | 0,695419786 |



Figure 5. The result of deleting holes of polygons of Region 31H.

Moreover, we added a constraint to our application in order to block the removal of holes that exceed a threshold area and a threshold SD. This allows for a parameterization of our approach that can be customized to specific needs. After deleting the holes of polygons belonging to Region 31H (Figure 5), we succeeded in removing 34 108 holes (25.5% of the initial number of polygons). The remaining number of polygons is therefore 99 672 (74.5%). Besides, it remains only one hole considering the thresholds that we set for area (10 km2) and SD (0.6).

3.4 Progressive merging of Cells

After deleting holes, we store the result in a new Access spatial database. Then, we start the cells merging process (by updating this database) in different steps using different intervals of suitability. In a fist step, we query the database to get polygons having SD = 1 (100 % of suitability). Then, we apply to the 6 261 obtained polygons a descending sort using their area as an order criterion. Thus, the first processed polygon will be the one having the best SD and the biggest area. For each of these polygons, we apply a spatial query in order to find their immediate neighbors. We sort these neighbors in descending order according to their SD and area. Then, we try to merge to the selected polygon its neighbors one after the other, until we obtain a merged polygon that has a SD that does not drop below a chosen stopping-threshold of 0.9. After this first merging pass, 7 137 polygons have been deleted from the database. Thus, until now we succeeded in removing 41 245 (30.8%) polygons and the remaining number of polygons is 92 535 (69.2 % of the initial Geobase polygons). Figure 6 shows (on the left) a polygon (in blue) with three neighbors (red). This polygon is merged to one of its neighbors (on the right) and the resulting polygon keeps a SD higher or equal to the chosen threshold.

Then, we apply the merging process in successive steps by increasing the ATI (used to select the initial processed polygons) and decreasing the threshold of SD used to stop merging. In fact, we used in the next process the polygons that have a SD belonging to the interval [1, 0.9] and a stopping-

threshold of 0.8. For the next merging iterations, we have used respectively the following intervals and thresholds: ([1, 0.8], 0.7), ([1, 0.7], 0.6), ([1, 0.6], 0.5).



Figure 6. Merging of cells using threshold of SD.

After completion of the merging of suitable polygons, we apply another process in order to merge the unsuitable polygons as mentioned in Section 3.1. We select and sort by decreasing sizes polygons with a SD = 0 and we try to merge them with their neighbors which have a SD less than 0.5. We use this value since it represents the last stopping-threshold which is used by the merging process of suitable polygons. The proportions of cells removed after each of the applied merging processes are shown in Figure 7 (the first process is the deletion of holes and the last process is the merging of unsuitable polygons).



Figure 7. Proportions of cells removed after the different merging iterations.

3.5 Creating an informed VGE

After merging the Geobase cells using the different processes, it remains only 17 % (22 698) of the initial polygons belonging to Region 31H. For the purposes of the geosimulation of zoonosis propagation (see Section 4), we propose to add additional data to each polygon such as the identifiers of its neighbors which is very useful for a variety of spatial functions. To this end, we add to each polygon a new attribute which contains the list of the IDs of its neighbors. Then, we apply a spatial query to the Access database in order to get the immediate neighbors of each polygon. We also add to each polygon its neighbors' orientation which is very useful for some processes that we want to simulate such as bird migrations in which the direction of displacement is very important (see Section 4). Thus, we associate two kinds of orientation data with each neighbor ID of a given polygon. The first one is a quantitative information representing an angle between a virtual line (see Figure 8a) and the North axis. These lines are drawn using the centroid of a given polygon and the centroid of each of its neighbors. To this end, we use a Geomedia function that returns the forward azimuth (FA) of a line segment. The returned value is normalized to fall between zero and 2π and is measured clockwise from North. We use this quantitative information to also store a qualitative information which represents either one of the four cardinal directions (i.e. North, East, South, West) or one of the four ordinal directions (i.e. North-East, South-East, South-West, North-West) or one of the eight further divisions represented in Figure 8b.



Figure 8. Quantitative (a) and qualitative (b) information used for the neighbor's orientation.

4 The usefulness of the obtained VGE for the geosimulation

Our approach (generating an informed VGE through cell merging) provides a lot of benefits if we compare it to approaches which have been used up to now. For example, the raster-based VGE which are used by cellular automata [4] does not take into account some important factors related to disease spread such as the population's survival in a specific geographical space. Another example is the GDBSCAN [10] (generalizing density-based clustering algorithm) which may be used to cluster polygons, but is based on extracting densityconnected sets of neighboring objects within a circular region. Besides, the geosimulation of the zoonosis propagation needs to deal with huge populations of various species which are located in large territories. The use of the Geobase land-cover cells to generate the VGE representing such large territories is almost impossible, especially because the number of cells may reach almost one million (number of cells for the six regions shown in Figure 1). Fortunately, our approach can significantly reduce the number of cells used to create the VGE while respecting the user's needs (choice of criteria and thresholds). The reduction in the number of cells in which the populations of interest have to evolve and interact will improve the effectiveness of the simulation engine used to geosimulate the zoonosis propagation.

Moreover, we have to model the individuals' mobility in order to simulate the disease spread. Our approach facilitates this task because we do not need to use agent-based approaches which try to explicitly take into account the trajectories of each individual (i.e. agent) or group of individuals located in the VGE [5]. In the case of zoonoses, it is not feasible to use such approaches since we have to handle huge numbers of individuals. Besides, our informed VGE allows for modeling a variety of processes influencing the geosimulation of the zoonosis propagation. As an example, we are currently modeling the migration of birds that import juvenile ticks (some of them being infected) in Quebec from the US. We think of modeling the Spring migrations as waves which distribute across cells birds of various species carrying ticks. Using a geo-referenced database of birds crossing the border on migration corridors at different periods of Spring, our simulation is initialized by associating the incoming birds to selected cells located at the border of Quebec and US [11]. Then, bird groups spread to neighboring cells with respect to their attractiveness. This cell attractiveness for birds is computed using a qualitative attribute and associated with polygons (cells) in the same way as the suitability attribute presented in Section 3. The process runs until all individuals are distributed. Indeed, such a process is facilitated by the identification of each neighbor of each cell which is available in our informed VGE. Moreover, the distribution of individuals over cells should agree with the location of migration corridors. This is why it is important to know the orientation of neighboring cells with respect to a given cell (see Section 3.5) in order to properly distribute individuals. This is possible thanks to our informed VGE which stores in each cell the quantitative and qualitative orientation of its neighbors.

5 Conclusion and future works

In this paper, we presented a new approach to generate an informed VGE used to geosimulate the propagation of an infectious disease, taking into account the spatial-temporal characteristics of this phenomenon. We exploited vector-based land-cover data to progressively merge polygons according to their degree of suitability for selected biological processes. The resulting polygons are used as basic cells which are the fundamental elements of our informed VGE. We believe that this innovative approach can be used not only to simulate the propagation of zoonoses, but also that it can be adapted to various other phenomena that do not necessarily relate to the spread of infectious diseases.

As future works, we plan to refine our approach and we are particularly interested in improving the efficiency of the progressive merging process. Indeed, we currently handle only the immediate neighbors of each polygon. It can be worth to take into account in the merging process the neighbors of neighbors and so on (considering a breadth-first or depth-first search) until reaching a satisfactory result according to the user's appreciation which is based on the selected criteria and thresholds. Besides, we investigate an algorithm that will be able to get the intersections between polygons of different informed VGE. Indeed, we need to create a VGE for each species involved in the phenomenon and then to compute the intersections between the corresponding polygons in order to obtain a VGE composed of cells which are qualified by the suitability parameters of all the species of interest.

6 Acknowledgments

Many thanks to GEOIDE, the Canadian network of centers of excellence in geomatics (CODIGEOSIM Project), INSPQ (*Institut national de santé publique du Québec*) and the Saint-Hyacinthe Division of the Public Health Agency of Canada (PHAC) for their support (finance, expertise and data).

7 References

- P. Gosselin, G. Lebel, S. Rivest, and M. Douville-Fradet. "The integrated system for public health monitoring of West Nile virus (ISPHM-WNV): a real-time GIS for surveillance and decision-making." *International Journal* of Health Geographics. 2005, 4:21.
- [2] M.J. Wonham, T. De-Camino-Beck, and M.A. Lewis M.A. "An epidemiological model for West Nile virus:

invasion analysis and control applications." In Proceeding of Royal Society of London. Series B, Biological Sciences, 2004 Mar 7, 271(1538):501-507.

- [3] R. Liu, J. Shuai, H. Zhu, and J. Wu. "Modeling Spatial Spread of West Nile Virus and Impact of Directional Dispersal of Birds." *Mathematical Bioscience Engineering*, 3, 2006, pp. 145–160.
- [4] S.H. White, A. Martin del Rey, and G. Rodriguez Sanchez. "Using Cellular Automata to Simulate Epidemic Diseases." *Applied Mathematical Sciences*, 3(20), 2009, pp. 959-968.
- [5] S. Emrich, S. Suslov, and F. Judex. "Fully Agent Based. Modellings of Epidemic Spread Using Anylogic." In Proceeding of EUROSIM, September 2007, Ljubljana, Slovenia, pp. 9-13.
- [6] M. Bouden, and B. Moulin. "Zoonosis-MAGS: A generic multi-level geosimulation tool for zoonosis propagation". In Proceeding of Global Geospatial Conference 2012. Spatially Enabling Government, Industry and Citizens. Quebec City, Canada, 14-17 May 2012.
- [7] M. Bouden, B. Moulin, and P. Gosselin. "The Geosimulation of West Nile Virus Propagation: A Tool for Risk Management in Public Health." *International Journal of Health Geographics*, 7 :35, 2008, pp. 1-19.
- [8] M. Bouden, and B. Moulin. "A Spatio-Temporal Interaction Model used to Geosimulate the Zoonosis Propagation". In Proceeding of Symposium on Theory of Modeling and Simulation (TMS'12). Spring Simulation Multi-Conference. The Society for Modeling & Simulation International. March 26-29 2012, Orlando, FL, USA.
- [9] R. Lisichenko. GIS Using Geomedia Professional V6, OnWord Press, 1 edition, 2008, 480 pages.
- [10] J. Sander, M. Ester, H.P. Kriegel, and X. Xu. "Density-Based Clustering in Spatial Databases: The Algorithm GDBSCAN and Its Applications". *Data Mining Knowledge Discovery*. 2(2): 169-194, 1998.
- [11] C.M. Francis, and D.J.T Hussell, 1998. Changes in numbers of land birds counted in migration at Long Point Bird Observatory, 1961–1997. *Bird Population*, 4:37–66, 1998.

Games for Children with Cerebral Palsy

Leonidas Deligiannidis Wentworth Institute of Technology 550 Huntington Av. Boston, MA 02115, USA 001-617-989-4142

deligiannidisl@wit.edu

ABSTRACT

The objective of this study is to determine the suitability and feasibility of novel multimodal computer games utilizing Virtual Reality technology and gaming for children with Cerebral Palsy (CP). The goal of the games is to provide fun experiences and sensations to the user where he/she will become motivated to engage in physical activity, and thus, we may provide a medium for motor, speech, and memory rehabilitation.

A common problem with children with CP is a reduction in motivation. Immersive Active Gaming (IAG) may allow these children to engage in physical activities despite their negative motivation and other psychological disorders such as depression. With our games we hope to not only encourage and motivate the participants to engage in physical activity but also promote equality. We can promote equality by designing and enabling games to be played by children with and without CP; utilizing appropriate interaction devices based on a child's abilities. Sisters, brothers, parents, relative, and friends can all play the same games together with a child with CP. After all, we are all equal with our own strengths and weaknesses.

Keywords

3D multimodal Games, Cerebral Palsy.

1. INTRODUCTION

One of the most important problems in designing games for children with Cerebral Palsy (CP), is that each case is unique because the abilities of each child are often different. Children with CP see the difference between themselves and other children without disabilities. For example, playing at school recess, children with CP cannot run and jump on monkey bars as their classmates. Gymnastics classes look and feel different for them, going to the restroom is not an activity, depending on the severity, that can be performed easily and may require physical assistance. Even playing video games becomes a difficult task, since they are required to manipulate a dozen buttons and a couple of mini-joysticks present on most gaming console controllers. This could lead to depression and sadness. The rate of depression is three to four times higher in people with disabilities such as cerebral palsy and it seems to depend not on the disability itself but rather on how well they cope with the disability. Emotional support, self-esteem, and stress are all factors that impact one's mental health [1]. Questions they could ask to themselves are: "Why cannot I play like the rest of my friends at recess?", "Why do I have to go to physical therapy, which is not such a fun activity anyways?", "Why am I different? I didn't do anything wrong to deserve this". Instead of focusing on ones weaknesses, we should focus on strengthening one's abilities. This can be done if we have the right tools and mindset. As a starting approach, we first need to boost their self-esteem for two major reasons: a) to avoid psychological problems these children may develop, and b) to improve their self-esteem so that they feel less disadvantaged which will improve their lives. So, the games should be designed to utilize new devices and interaction techniques that could negate the physical differences in the ability of each player. For example, a child with no physical disabilities can play a game using a mouse, a keyboard, or a joystick. A child with difficulty moving his/her fingers, could play games where moving the whole arm indicates an event. Or if a child has severe hand movement difficulty, he/she could use the neck or legs to indicate turns and direction of movement in the game. A child with no movement abilities could use voice commands to achieve the same goal. Others have used Sony's PlayStation 2 EyeToy camera to capture mouse motion [2][3], and others a web camera [4][5].

When designing games for children with CP, it is important to take special care not to exacerbate the children's disabilities. Games that would not be well suited are those that are based on timing, speed, or dexterity since these games may frustrate the children and not allow them to enjoy the actual story of the game. In case of competitive games (children competing with each other), we should investigate carefully the potential difficulties in designing games that can adapt to physical disabilities. Then the question becomes: "How many differences in a player's ability can be accommodated and when real competition is compromised by such adaptations?"

To promote equality, games should be designed so that children with and without disabilities can play the games together. This will make the children with CP feel more equal since they can play the same game with other children without disabilities. However, a devise that works well for children with disabilities could impair a child without disabilities. Thus, the game should allow for different types of devices to perform the same tasks so that we don't take away the fun of the game play.

The intriguing component of our study is to "level the playing field" for children with CP, allowing them to play and compete on more equal terms with other players. This focus is rather different from the majority of other games for such audiences that tend to focus more on remediating the disability itself.

Current gaming interfaces draw their strengths from exploiting the user's pre-existing abilities and expectations rather than trained behaviors. For example, navigating through a conventional computer system requires a set of learned, unnatural commands, such as keyboard keys to be typed in, mini-joysticks on game pads to be manipulated, etc. On the other hand, navigating through Immersive Virtual Environments exploits the user's existing realworld "navigation commands" such as positioning the head and eyes, turning the body, walking toward something of interest, and pointing at the direction of interest. This naturalness can reduce the mental effort required to operate the system and thus enables one to focus on enjoying the game-play instead of the mechanics of how to play the game. As a result, this reduces the time needed to learn how to play a game because it avoids information overload, time pressure, and stress which directly affects the outcome of the gaming experience. This, most likely, takes away the fun of the game and possibly turns game-playing into an activity that provides no enjoyment, discomfort, and intimidation.

2. CEREBRAL PALSY

Cerebral Palsy (CP) is a group of lifelong neurological disorders cased of irreversible damage to brain cells. It affects movement, learning, hearing, seeing, and thinking. It occurs due to problems with brain development as early as while the baby grown in the womb. But, it can occur at the time of birth and during the first 2 to 3 years of age. CP is a non-contagious and non-progressive brain disorder. The severity of CP varies from child to child and this makes every case unique. The United Cerebral Palsy Foundation [6] reports that about 750,000 children and adults in the United States have symptoms of cerebral palsy. An additional estimate of 8,000 babies and infants will be diagnosed every year.

The motor area of the brain's outer layer (cerebral cortex) is responsible for directing muscle movement. This is the affected area of the brain that leads to impairment of motor functions; the muscles continually receive signals to contract disabling a person to regulate their muscle tone, and this causes tightness and stiffness of their muscles. In most cases, the muscles and the nerves are healthy. But because the component of the brain that controls the

muscles is injured, children with CP face difficulties controlling the movement of their muscles. Many children with CP have increase muscle spasticity in one or both arms and / or legs. As their muscles continuously contract, their muscles become stiff and tight and this interferes with walking, moving, and even speech. There are cases where there are bone and joint deformities where the muscles become permanently fixed and tight and may required surgical intervention. Additionally, CP often is accompanied by disturbances of sensation, perception, cognition, communication, and epilepsy.

Unfortunately, there is no known cure for CP yet. Therapy normally is performed by an occupational, physical, or speech therapist. Since there is no cure for CP, the goal of the treatment is to make people with CP as independed as possible to live an effective life. Because of muscle movement problems, which restrict people with CP to walk, run and in general exercise their muscles, it is possible that they end up with decreased muscle mass. By performing a specific set of exercises, physical therapy focuses in increasing the performance of the muscles by preventing the muscles to become atrophic and fixed in an On the other hand, occupational abnormal position. therapy focuses on enabling people with CP to master to perform daily activities such as dressing, eating, etc. Speech therapy may be required in many cases where the people with CP are unable to produce intelligible speech, or have problems in other areas of communications such as facial expressions and hand gestures and postures. In certain circumstances, a speech therapist can teach a one to use special communication devices such as a voice synthesizer. Recreational therapists can also help children with CP to improve gross motor skills by using, for example, horseback riding techniques.

3. BACKGROUND

Most of the recent Virtual Reality (VR) research is focused on the assessment of selected cognitive, functional and behavioral functions and processes. This can complement traditional in-person rehabilitation techniques, which makes VR technology a great tool for medical practitioners for both the assessment and cognitive/functional rehabilitation of people with CP. VR has been used as a medium for the assessment and rehabilitation of a variety of clinical populations. This includes populations with cognitive [7-10] and meta-cognitive [11-16] deficits, such as visual perception, attention [17], and memory [15-16] [18-23]. Other applications are directed at the rehabilitation of motor movement difficulties [24-26] to provide recreational opportunities for people with disabilities.

Virtual Reality has been adopted by many researchers as a suitable technology to be used as an assessment and a treatment tool in physical and occupational therapy [27-35]. The ability to provide simulated and meaningful virtual environments, Virtual Reality offers us opportunities to engage in purposeful tasks related to real-life events.

While VR is capable to provide fully immersed environment experiences, at times it seems more appropriate to reduce the level of immersion because of physical capabilities people possess. This reduces the sense of presence and realism in the Immersive Virtual Environment; however, it allows a subject to experience the environment in a less restricted form. For example, in a desktop-VR environment, users are not attached to heavy and relatively large head mounted displays (HMDs), or to be placed within a physical projection room such as the CAVE, or wear cumbersome gloves and other interaction devices to feel the texture and weight of virtual objects. VR is an attractive technology because of its advantage to provide synthetic environments and encourage participants to perform tasks that are difficult to achieve and deliver using conventional neuropsychological methods [36-38].

Since there is no cure for CP, major goals of therapy include enhancement of functional ability by improving sensory, motor, and cognitive functions. Using VR technology, this is achieved by enabling the participants to interact with synthetic environments, engage in activities required by the game or synthetic environment, manipulate objects in a way that immerses them within the simulated environment, which helps produce a feeling of presence in the virtual world [39][40]

There are some limitations of the VR technology, not because of the technology itself, but rather because of the one's ability to utilize awkward VR devices in order to interact with a synthetic environment. For example, muscle spasticity limits and makes uncomfortable a subject's ability to use equipment such as a glove, or wearing a relatively bulky and heavy-head-mounted display (HMD), etc [41], which could also have a side effect of nausea; something that is not present in desktop-VR. Low vision people are limited in using a HMD simply because of the subject's limited field of view and decrease peripheral vision [42], or in a case where the subject is a child, by overloading the neck muscles. In addition, it has been reported that low frequency flashing lights used as visual effects can trigger seizures [43][44].

4. GAMING ENVIRONMENT

Our gaming laboratory, where we develop the games and perform user studies, is located on the ground level at Williston Hall. The laboratory is handicapped accessible and it is located near the main parking lot. The dimensions of the lab are 24x16 feet and there are two high end PCs (Dell XPS 730X) running Windows Vista Ultimate with 6 G of memory, two drives in RAIN 0 configuration, a sound card and a dual NVIDIA GeForce GTX280 SLI enabled video card. Both PCs are equipped with 2 Intel i7 quad CPUs running at 3.2GHz.

Both PCs can direct their output, using Genfen (www.gefen.com) DVI switchers, to our main display which is a SHARP PN-S525 1920x1080 52 inches LCD monitor. The PCs are also connected to a 1000W ONKYO

TX-SR707, 7.2 channel receiver. Below the display we installed a fan with adjustable speed control, which enables us to render wind. Its speed is controlled by the game using a Phidgets (www.phidgets.com) Advanced Servo controller. The fan is placed to aim at the player's feet and not the face to avoid breathing problems the fan might cause.



Figure 1. Servo controller connected to the adjustable speed fan, and the robotic arm.

Its idle speed is around 100 feet per minute while its maximum speed is around 500 feet per minute. The servo controller is connected via USB to a WiFi Phidgets SBC single board computer which is the gateway to the computer running the game. We also built a custom robotic arm for the final component of the game, which is also controlled by the servo controller as shown in figure 1.



Figure 2. Diagram of the Gaming Laboratory.

navigation, we use the Polhemus Latus For (www.polhemus.com) 3D tracker with two receptors and one marker. The Latus system provides six degrees of freedom (6DOF) tracking wirelessly. Multiple receptors connect to the Latus to form a network of tracking units that extend the envelope of tracked space. The Latus tracks the light-weight wireless markers in space which contain their own battery and transmition frequency. The unique transmition frequency identifies each 6DOF marker. A diagram of the lab is shown in figure 2 and a snapshot of the lab setup is shown in figure 3.



Figure 3. A snapshot of the set up in the Gaming laboratory

We placed all the furniture and the hardware equipment around the lab leaving the center area of the lab available for the players. There is enough room for two players to stay side-by-side to play a game. We removed all metallic furniture away from the Latus receptors and we secured the receptors on custom made 5/4 inches PVC pipe stands using plastic screws and glue to minimize signal interference as the Latus tracks the markers in space using a magnetic field.

5. BLOB VILLAGE

The game we developed is the Blob Village. It is played by two players. Blobs are the creatures living on a beautiful little village by the snowed mountains where there is a little lake, a castle, and a swimming pool. Figure 4 shows a bird's eye-view of the Blob Village.



Figure 4. Bird's eye-view of the Blob Village.

The main characters of the game are the blobs that have their own personality. They wander around their village and provide verbal clues to the player as to what to do next. The blobs are friendly creatures who love to interact with children. In fact, if they see you approaching them, they will leave whatever they do and come to talk to you. They even try to get your attention by waving their hands and calling you to go closer to them. A snapshot showing a blob waving his hand is shown in figure 5.

The user moves around the village trying to communicate with the blobs to find the "blue" blob who feels sad today and wants to be hugged by a happy child and play together. During the play, the player is able to walk around and even fly to get to different places in the Blob Village. As the player walks and flies in the game, he / she experiences different 3D sound effects; while soft background music is playing. In addition, depending on the speed of the walk and whether the player if flying or not, the player feels the intensity of the wind. To walk and fly, two players need to coordinate to perform the right action. The players could even use different input devices to accomplish the same task. This depends on the physical abilities of the children and that is why every setup of the game could require different input devices. The duration of a day in Blob Village is much shorter than what it is in real life. This is reflected in the game by moving through the stages of dawn to dusk. Using shader programming, the effects of the sun, water, and everything else in the game, change depending on the time of the day in the Blob Village.



Figure 5. A blob in the game waving his hand.

After the player finally finds the sad "blue" blob, the blob asks the player if he or she wants to touch, hug, and play (physically) with the blob. When the player agrees to play with the sad blob so that the blob would feel happy again, the blob starts spinning in the game until it disappears from the display.



Figure 6. A snapshot of the Blob Village while the player is flying.

At that moment, a hidden custom made blob toy becomes visible. One of the servo motors kept this toy blob hidden from the player's view behind a curtain right below the display and it is emerged into view at this point giving the feeling of the blob being materialized. While playing a sound effect, the robotic arm turns slowly to reveal the blob toy right below the screen and asks the player to give it a hug.

This moves a 3D multimodal game from being a totally synthetic environment, into the reality where we can grab, hug, kiss, and play with the main character of the game. Meanwhile, to find the sad blob required the cooperation of two players, one with physical disabilities and without. None players could play the game all by themselves, it required two players - two equal players, it just happened that some of us have some difficulty in moving some parts of our body. On the other hand, the child's self-esteem could be increased by playing the game because at the end, after we experience the beautiful environment of Blob Village, we managed to find the sad blob and play with it and make it feel happy again. Figure 6 is a snapshot of the Blob Village while flying, and figure 7 is a snapshot of the final stage where the player found the "blue" sad blob which is located by the swimming pool.



Figure 7. A snapshot of the sad blob located at the swimming pool.

5.1 Blob Behavior

The Blobs, the main characters of the game, are three dimensional blob-like characters who wander around the blob village. Each blob is designed to have 4 animations, shown in figure 8: a) "idle" where they wobble left and right, b) "left wave" where they wave their left hand, c) "right wave" where they wave their right hand, and d) "jump" where they perform a jump.

We can animate the blobs while they are moving as well as when they are standing still. They are all of different color and there are 4 different textures we apply on them to dynamically animate their facial expressions; which are performed independently of their body animation. The four textures that we animate on the blobs, by swapping them in and out using a shader, are shown in figure 9.



Figure 8. The four animations of the blobs (idle, left wave, right wave, and jump)



Figure 9. The four animated textures utilized to provide facial expressions.

As the blobs wander around in their village, they try to avoid bumping into each other. When they are in a course of collision, depending on their unique age they either stop moving and make a sound, or go around the stopped blob making another sound. The player is a special entity in the game. So when the player is approaching a blob, the blob turns to the direction of the player, makes an inviting sound (plays one of the 33 pre-recorder by children phrases such as "hello", "hi there", "hey look at me", etc), waves his left or right hand, and walks towards the player. When a blob gets close to the player, he stops moving and provides verbal hinds such as "now go to the island", or "now fly over the mountains to go to the castle", etc. Because the blobs can wander all over the village, we constrained them to stay in one of the 4 areas in the village. This way there are blobs everywhere. The critical component is to always have the blue sad blob by the pool, since we are supposed to find the blue sad blob and make him feel happy again. The four rectangle areas define the main land, the island, the castle, and the swimming pool as shown in the figure 10 below.



Figure 10. The four areas of the Blob Village

5.2 THE STORY

The story of the game is relatively simple, yet challenging for the players. The player starts off in the main land where she tries to communicate with the blobs to figure out what to do next. One of the blobs tells the player that there is a sad blue blob in the village that really wants to play with the player. To find the sad blob, the player needs to go to the island. This requires the player to fly over the lake around the island where she meets another group of blobs. These blobs instruct the player to go to the castle. At the castle, the blobs instruct the player to go to the swimming pool where the blue blob has been isolated there because he is sad and wants to play with the player. Finally, the player goes to the swimming pool where she meets the blue blob. The blue blob then becomes very happy and asks the player if she wants to play with him. Meanwhile, getting from place to place, the player experiences 3D sound effects, challenging navigation techniques, gained the ability to fly, etc. As the player flies, a fan renders the wind which depends on the player's speed and elevation. Now when a 3D character asks a child to play with him, our first reaction is "...but we are already playing the game". However, the game is about to turn from a synthetic 3D game into a game of reality. Upon the player's positive reply, the blue blob spins around faster and faster until it disappears from the display. At that moment, a robotic arm that is placed right below the display starts turning, revealing the blue blob toy right in front of the player. The player now is free to reach out a grab the real blue blob toy as shown in figure 11, and give him a hug.



Figure 11. The game reveals the toy blue blob that can be touched and hugged by the player.

6. RESULTS

We designed a large village with many components (3D terrain with collision enabled, water shaders, daylight effects, 3D models, etc), the main character of the game (the blobs with animations and facial expression animations), 3D sound effects (soothing music in the background, blob sounds and speech, which are recordings of real children), synthetic behavior (such as day, afternoon, and night effects), other environmental effects such as the utilization of a fan to render wind, as well as a robotic arm that reveals the blue toy blob. For navigation, we used the Polhemus Latus six degrees of freedom (6DOF) device.



Figure 12. The five postures for navigating in the Blob Village: a) stop, b)left turn, c) right turn, d) forward, and e) fly.

The Latus uses wireless tracking of markers making the entire game totally untethered; the player is not attached to any wiring. For navigating in the game, we used 5 gestures that implement: a) "stop" to stop motion, b) "turn left" to turn the player's view to the left, c) "turn right", to turn the player's view to the right, d) "forward" to move the player's position forward, and e) "fly" to change the elevation of the player as shown in figure 12. Some of the gestures can by performed by the second child playing the game in cooperation.

We plan to begin our formal evaluation with children with CP soon, using Robson's scale [46] for the measurements. During development and after we finished BlobVillage, 11 students taking our Introduction to Games Programming course, three faculty members and 4 children without CP tried the game and they gave us positive feedback. This is encouraging to us to continue our research with formal evaluations. They found the usage of the fan to render wind, the 3D sound effects, and the final stage of the game (when the blue blob toy is revealed using the robotic arm) very interesting and new. They found the scenery beautiful but they thought that the fan speed was a little bit high and also that the sound volume was a bit too high as well. We plan on adjusting the fan speed and the volume in the final release of the game.

7. FUTURE WORK

We plan on developing more interaction and navigation Polhemus techniques using the (http://www.polhemus.com/) Latus 6DOF device as well as other custom-built devices, the Wiimote controller from Nintendo, and the Kinect 3D camera from Microsoft. We have several families with children with Cerebral Palsy that cannot wait to come to our lab for the formal trials. We want to find out if using inexpensive off the shelf devices such as the Nintendo's Wiimote, or the Kinect camera, is as effective as using the relatively expensive Polhemus Latus tracker. What makes the Wiimote and the Kinect attractive is that they are wireless. Our game does support the Wiimote but its usage in the current version of the game as an interaction device is left out for future research.

8. CONCLUSION

We want to enable children, and adults, who may otherwise feel social isolation secondary to self-esteem issues to engage in game playing [45]. The importance and originality of our work is not only in the use of game-play to boost self-esteem, but also in the inclusive nature of game-play for people who might not otherwise benefit from standard commercially available games because of their standard interface models. The concept of inclusion is important for socio-cultural standpoints as well as practical standpoints. The bigger picture is not simply whether a person with physical disabilities can hold and use a game control device, but what the implications of excluding these people from the technology are.

In this paper we presented a sophisticated multimodal game that will become the basic platform for our future research. From the feedback we received from our students, faculty members, and children that played the game we have no doubt that we can use similar games to boost the selfesteem of children with CP. In addition, by playing these games in a multi-user environment we can promote equality among children with or without CP since these games are played cooperatively to achieve a common goal.

REFERENCES

- [1] National Institute of Health, National Institute of Neurological Disorders and Stroke http://www.ninds.nih.gov/disorders/cerebral_palsy/cerebral_palsy.ht m Retrieved Jan 21 2010
- [2] D. Rand., R. Kizony, PL Weiss., "Virtual reality rehabilitation for all: Vivid GX versus Sony PlayStation II EyeToy", Proc. 5th International Conference on Disability, Virtual Reality and Associated Technologies, Oxford, UK, 2004, pp. 87-94.
- [3] Patrice L Weiss, Debbie Rand, Noomi Katz, and Rachel Kizony, "Video capture virtual reality as a flexible and effective rehabilitation tool", Journal of NeuroEngineering and Rehabilitation Vol. 1(12), 2004
- [4] Margrit Betke, James Gips, and Peter Fleming, "The Camera Mouse: Visual Tracking of Body Features to Provide Computer Access for People With Severe Disabilities". IEEE Transactions on Neural Systems and Rehabilitation Engineering, Vol. 10 No. 1, March 2002 pp1-10.
- [5] Rick Kjeldsen, "Improvements in Vision-based Pointer Control". Proc. of the ACM SIGACCESS conference on Computer and Accessibility, Portland, Oregon, Sep. 29 2006 pp 189-196
- [6] United Cerebral Palsy Foundation, http://www.cerebralpalsysource.com/Resources/foundation_cp/index .html Retrieved Jan. 25 2010
- [7] Rizzo, A.A., Buckwalter, J.G., Humphrey, L., van der Zaag, C., Bowerly, T., Chua, C., Neumann, U., Kyriakakis, C., van Rooyen, A. & Sisemore, D. "The Virtual Classroom: A Virtual Environment for the Assessment And Rehabilitation Of Attention Deficits". CyberPsychology and Behavior, Vol. 3(3), 2000, pp.483-499.
- [8] Zhang L, Abreu BC, Masel B, Scheibel RS, Christiansen CH, Huddleston N, Ottenbacher KJ. "Virtual reality in the assessment of selected cognitive function after brain injury". Am J Phys Med Rehabil. Aug. 2001. Vol.80(8), pp597-604.
- [9] Grealy MA, Johnson DA, Rushton SK. "Improving cognitive function after brain injury: the use of exercise and virtual reality". Arch Phys Med Rehabil. Jun. 1999. Vol.80(6), pp661-667.
- [10] Weiss PL, Naveh Y, Katz N. "Design and testing of a virtual environment to train stroke patients with unilateral spatial neglect to cross a street safely". Occup Ther Int. 2003. Vol.10(1). pp39-55.
- [11] Lam YS, Tam SF, Man DWK, Weiss PL. "Evaluation of a computer-assisted 2D interactive virtual reality system in training street survival skills of people with stroke". Proceedings of the 5th International Conference on Disability, Virtual Reality & Associated Technology. Oxford, UK, 2004, pp27-32
- [12] Larson, P., Rizzo, A.A., Buckwalter, J.G., van Rooyen, A., Kratz, K., Neumann, U., Kesselman, C., Thiebaux, M., Van Der Zaag, C. "Gender issues in the use of virtual environments". CyberPsychology and Behavior, 1999, Vol. 2(2), pp113-123.
- [13] Robert S. Astur, Maria L. Ortiz, Robert J. Sutherland, "A characterization of performance by men and women in a virtual Morris water task: A large and reliable sex difference". Behavioural Brain Research, Jun. 1998, Vol.93(1-2), pp185–190.
- [14] Thomas, K.G., Laurance, H.E., Luczak, S.E., Jacobs, W.J. "Age related changes in a human cognitive mapping system: Data from a computer-generated environment". CyberPsychology and Behavior, 1999, Vol.2(6), pp545-566.

- [15] A Rizzo, J G Buckwalter, P Larson, A van Rooyen, K Kratz, U Neumann, C Kesselman, M Thiebaux, "Preliminary findings on a virtual environment targeting human mental rotation/spatial abilities". Proc. of the 2nd European Conference on Disability, Virtual Reality and Associated Techniques, 1998, Sköve, Sweden, pp 213–220.
- [16] McComas, Joan., Pivik, Jayne., Laflamme, Marc. (1998). "Children's transfer of spatial learning from virtual reality to real environments". CyberPsychology and Behavior 1998, Vol.1(2), pp121-128.
- [17] Wann, J.P., Rushton, S.K., Smyth, M., & Jones, D. "Virtual environments for the rehabilitation of disorders of attention and movement". Virtual reality in Neuro-Psycho -Physiology. Amsterdam Netherlands, 1997, Vol.44, pp. 157–164.
- [18] Luciano Gamberini, "Virtual Reality as a New Research Tool for the Study of Human Memory". CyberPsychology & Behavior., Vol.3(3), June 2000, pp337-342.
- [19] Paul N. Wilson, "Active Exploration of a Virtual Environment Does Not Promote Orientation or Memory for Objects", Environment and Behavior, Vol.31(6), 1999, pp752-763
- [20] Jaime Sánchez, Héctor Flores, "Memory enhancement through audio", ACM SIGACCESS Accessibility and Computing, Issue 77-88, 2003, pp24-31
- [21] Rose FD, Brooks BM, Attree EA, Parslow DM, Leadbetter AG, McNeil JE, Jayawardena S, Greenwood R, Potter J., "A preliminary investigation into the use of virtual environments in memory retraining after vascular brain injury: indications for future strategy?", Disabil. Rehabil., Vol.21(12), 1999, pp548-554
- [22] Dinh, H.Q.; Walker, N.; Hodges, L.F.; Chang Song; Kobayashi, A., "Evaluating the importance of multi-sensory input on memory and the sense of presence in virtual environments", Proc. of the IEEE Virtual Reality, 1999, Los Alamitos, CA, pp. 222-228.
- [23] Grealy, M.A., Johnson, D.A., Rushton, S.K., "Improving cognitive function after brain injury: The use of exercise and virtual reality", Archives of Physical Medicine and Rehabilitation, 1999, Vol.80(6), pp661-667.
- [24] Rachel Kizony, Noomi Katz, Patrice L. (Tamar) Weiss. Adapting an immersive virtual reality system for rehabilitation. The Journal of Visualization and Computer Animation, Special Issue: Virtual Reality in Mental Health and Rehabilitation, Nov. 2003 Vol.14(5), pp261-268.
- [25] Sveistrup H, McComas J, Thornton M, Marshall S, Finestone H, McCormick A, Babulic K, Mayhew. "A. Experimental studies of virtual reality-delivered compared to conventional exercise programs for rehabilitation". Cyberpsychol Behav. Jun. 2003 Vol.6(3)6, pp245-249.
- [26] Alma S Merians, David Jack, Rares Boian, Marilyn Tremaine, Grigore C Burdea, Sergei V Adamovich, Michael Recce and Howard Poizner. "Virtual reality-augmented rehabilitation for patients following stroke". Physical Therary. Sep. 2002, Vol.82(9), pp898-915.
- [27] Albert "Skip" Rizzo, "A SWOT Analysis of the Field of Virtual Reality Rehabilitation and Therapy", Presence: Teleoperators and Virtual Environments, Vol. 14(2), April 2005 pp 119-146
- [28] Patrice L. (Tamar) Weiss, Noomi Katz, "The potential of Virtual Reality for rehabilitation". Journal of Rehabilitation Research and Development Vol.41(5) Oct. 2004 pp:vii-x
- [29] Elkind, James S." Uses of virtual reality to diagnose and habilitate people with neurological dysfunctions". CyberPsychology and Behavior, Fall 1998 Vol. 1(3), pp263-274.

- [30] Pugnetti, Luigi, Mendozzi, Laura, Motta, Achille, Cattaneo, Annamaria, Barbieri, Elena, "Evaluation and retraining of adults' cognitive impairments: Which role for virtual reality technology?" Computers in Biology and Medicine, March 1995, Vol.25(2), pp213-227.
- [31] Rose, F.D., Attree, E.A., and Brooks, B.M." Virtual environments in neuropsychological assessment and rehabilitation". Virtual Reality in Neuro-Psycho-Physiology. 1997, Amsterdam, Netherlands, pp. 147-156.
- [32] Christiansen C, Abreu B, Ottenbacher K, Huffman K, Masel B, Culpepper R., "Task performance in virtual environments used for cognitive rehabilitation after traumatic brain injury", Archives of Physical Medicine and Rehabilitation, Vol.79(8), 1998, pp888-892.
- [33] Davies, R.C., Johansson, G., Boschian, K., Lindé, A., Minör, U., Sonesson, B., "A practical example using virtual reality in the assessment of brain injury". The International Journal of Virtual Reality, 1998, Vol.3(4), Sweden, pp. 1-7.
- [34] Strickland, D., "Virtual reality for the treatment of autism", Stud. Health Technol Inform, Vol.44, 1997, pp81-86.
- [35] Cromby, J.J., Standen, P.J., Newman, J., Tasker, H. "Successful transfer to the real world of skills practiced in a virtual environment by students with severe learning difficulties". Proc. of the First European Conference on Disability, Virtual Reality and Associated Technology, Maidenhead, UK, 1996, pp. 103–107.
- [36] Greenleaf, W.J., and Tovar, M.A. "Augmenting reality in rehabilitation medicine". Artificial Intelligence in Medicine, 1994, Vol.6, pp289-299.
- [37] Kuhlen, T., and Dohle, C. "Virtual reality for physically disabled people". Computers in Biology and Medicine, 1995, Vol.25(2), 205-211.
- [38] Wilson, P.N., Foreman, N., and Stanton, D. (1997). "Virtual reality, disability and rehabilitation". Disability and Rehabilitation, 1997, Vol.19(6), pp213-220.
- [39] Slater Mel, "Measuring presence: A response to the Witmer and Singer Presence Questionnaire", Presence: Teleoperators and Virtual Environments Vol. 8(5), 1999, pp560-565.
- [40] Eric B. Nash, Gregory W. Edwards, Jennifer A. Thompson, Woodrow Barfield, "A Review of Presence and Performance in Virtual Environments". Int. Journal of Human-Computer Interaction, Vol. 12(1), May 2000, pp1-41
- [41] Wikipedia "Spasticity" http://en.wikipedia.org/wiki/Spasticity. Retrieved Jan. 5 2010
- [42] Scholar pedia "Hemineglect" http://www.scholarpedia.org/article/Hemineglect Retrieved Jan. 5 2010
- [43] D Kasteleijn-Nolst Trenité DG, Martins da Silva A, Ricci S, Rubboli G, Tassinari CA, Lopes J, Bettencourt M, Oosting J, Segers JP, "Video games are exciting: a European study of video game-induced seizures and epilepsy", Epileptic Disord, vol. 4(2), pp. 121-8, June 2002.
- [44] Michelle Bureau, Edouard Hirsch, and Federico Vigevano, "Epilepsy and Videogames", Epilepsia Vol. 45(s1) pp24-26 Jan 2004
- [45] J. Beecham, T. O'Neill, and R. Goodman, "Supporting young adults with hemiplegia: services and costs", Health Soc. Care Community, vol.9(1), pp. 51-9, January 2001.
- [46] Robson, Philip. J. (1989). Development of a New Self Report Questionnaire to Measure Self Esteem. *Psychological Medicine*, 19, 513-518.

Hardware-Software Cosimulation of Feedback Controller for Synchronization of Inferior Olive Neurons

Desta Edosa¹, Keum W. Lee², Venkatesan Muthukumar¹, and Sahjendra N. Singh¹

¹ Department of Electrical and Computer Engineering, University of Nevada, Las Vegas, 4505 Maryland Parkway, NV 89154-4026, USA

² Division of Electronic Information and Communication, University of Kwandong, Gangwon, Korea

Abstract—The design of a control system for the synchronization of Inferior Olive Neurons (IONs) and hardwaresoftware cosimulation of the closed-loop system are presented in this paper. Each ION is described by a set of four nonlinear differential equations. These IONs exhibit limit cycle oscillations (LCO), but are not necessarily in phase. The objective is to control one of the IONs so that both oscillate in unison (with zero relative phase). A simple linear feedback control law for the synchronization of the IONs using an output variable is developed. In the closed-loop system, asymptotic convergence of the state vectors of the IONs is accomplished. Then the hardware-software cosimulation of the complete closed-loop system is considered. Hardware synthesis and HW-SW cosimulation tests are performed to examine the performance of the controller. The oscillatory properties of the Inferior Olive Neurons (IONs) can be used to provide timing signals for motors to mimic vertebrae like movements. These vertebrae like movements can be used in bio-inspired underwater and unmanned vehicles with oscillating fins.

Keywords: Inferior Olive Neuron, ION synchrony, HW-SW simulation

1. Introduction

The synchronous activity of olivo-cerebellar system, is one of the key neuronal circuits in the brain, provides motor control signals for the movement execution. This neuronal network is organized around clusters of inferior olive neurons (IONs) [6], [16], [15]. The inferior olive neurons (IONs) have various features including, the subthreshold activity in which the membrane potential has sustained fluctuations. This rhythmic activity has been termed as spontaneous subthreshold oscillations. The orbits of the IONs may have different shapes (sinusoidal, quasi-periodic, periodic waveform with spikes, and irregular) [12].

The dynamical behavior and functional significance of interconnection has been explored by Armstrong [1]. Benardo and Foster [3] have examined the oscillatory pattern of IONs. A structural study of inferior olivary nucleus of cat has been performed by Sotelo et al. [19]. For low amplitude oscillations, a model based on electrical coupling of neurons with heterogeneous channel densities has been considered

by Manor et al. [17]. A variety of mathematical models capturing the important characteristics of IONs have been proposed in literature. Velarde et al. [20] and Llinas et al. [16] developed ION models using Vander Pol Oscillator and FitzHugh-Nagumo (FN) systems. A relatively simple ION model has been also proposed by Kazantsev et al. [11], [10]. The dynamical behavior of ION depends on its parameters. The structure of the orbit of an ION undergoes drastic changes when its parameters vary. The study of qualitative changes in the orbit structure termed as, bifurcation of neurons has been performed by Guckenheimer and Labouriau [5] and Izhkevich [7], [8]. For inferior olive neurons, a bifurcation analysis has been treated by Katori et al. [9]; Lee and Singh [14]. Authors have also studied the synchronous activity of interconnected neurons. Neurons are connected via gap junction to form electrical coupling. Katori et al. [9] have studied the spatio-temporal dynamics of IONs using conductance-based model. Neuronal synchronization in the mammalian brain have been examined by Bennett and Zukin [4].

In view of the important role of the IONs in motor control, researchers have shown interest in the models of these IONs for the control of autonomous air and underwater bio-robotic vehicles. Robust oscillatory patterns of a variety of shapes and sizes of the IONs are suitable for executing different kinds of maneuvers. For the control of robotic vehicles, it is essential to control the relative phase angles of the cluster of IONs. As such the problem of synchronization of IONs is important [22]. In literature, synchronization of IONs has been explored and linear or nonlinear or adaptive control systems have been developed by Bandyopadhya et al. [2] and Lee and Singh [13], [14]. Bennet and Zukin [4] have shown that a model of gap junction in conjunction with postsynaptic capacitance as a low pass filter achieves neuronal synchronization. In view of the work of Bennet and Zukin [4], it is of interest to develop control systems which are simple in form because nonlinear and adaptive control systems developed for synchronization of IONs are not attractive from the viewpoint of implementation.

Contribution of this paper lies in the development of a simple robust control system for the synchronization of IONs. It is assumed that a single output of each ION is measured for feedback. One of the IONs is treated as a
reference ION, and the other ION is controlled by the application of an extracellular stimulus. A control law is developed so that the trajectories of the controlled ION asymptotically follows the the trajectories of the reference ION. The control input is the product of a suitably chosen gain and the scalar output error of the IONs. It is shown that synchronization is accomplished for a set of values for the feedback gain. The control system is robust to variations in the key parameters of the IONs and the controller gain. Then the hardware implementation of the complete closedloop system, including a pair of IONs and the feedback control system, is developed. Laboratory tests as well as numerical simulation results are obtained. It is seen that the hardware circuit produces the waveforms observed in digital simulation closely. The hardware circuit is especially useful for real-time control of air and underwater vehicles using flapping wings and fins.

The organization of the paper is as follows: Section 2 presents the mathematical model of the IONs. A control law for synchronization is developed in Section 3. The results of digital simulation is given in Section 4. Section 5 considers the HW-SW implementation and simulation of the closed-loop system, and laboratory test results and comparison of these with the simulated results are presented in Section 6.

2. IO neuron model

In this paper, the ION model described in Kazantsev et al. [11], [10] are considered for synchronization. Let the state vector of the *i*th neuron be $(u_i, v_i, z_i, w_i)^T$, i = 1, 2. (*T* denotes matrix transposition.) The model has polynomial nonlinearities in variables u_i and z_i of degree three. For simplicity in notation, often the arguments of various functions will be suppressed. The nonlinear equations describing the ION1 (ION₁) are

$$\dot{u}_{1} = k\epsilon_{Na}^{-1}[u_{1}^{2} - u_{1}^{3} + (u_{1}^{2} - u_{1})a - v_{1}]$$

$$\dot{v}_{1} = k(u_{1} - z_{1} + I_{Ca} - I_{Na})$$

$$\dot{z}_{1} = [z_{1}^{2} - z_{1}^{3} + (z_{1}^{2} - z_{1})a - w_{1}]$$

$$\dot{w}_{1} = \epsilon_{Ca}(z_{1} - I_{Ca} - \mu^{*} - I_{ext1})$$
(1)

The variables z_1 and w_1 are responsible for subthreshold oscillations and low-threshold (Ca^{2+} -dependent) spiking, and the variables u_1 and v_1 describe the higher-threshold (Na^+ - dependent) spiking. The oscillation time scales are controlled by the parameters ϵ_{Ca} and ϵ_{Na} ; and I_{Ca} and I_{Na} drive the depolarization level (equilibrium point) of the system.

The parameter k sets the relative time scale of the two systems. The parameters a_i 's (appearing in the nonlinear functions) play an important role in shaping the trajectories of the IONs. I_{ext1} denotes the extracellular excitation used here as the control input. The bias term μ^* provides flexibility in getting different kinds of waveforms. The ION_1 is treated as the slave ION.

The reference ION is described by

$$\dot{u}_{2} = k\epsilon_{Na}^{-1}[u_{2}^{2} - u_{2}^{3} + (u_{2}^{2} - u_{1})a - v_{2}]$$

$$\dot{v}_{2} = k(u_{2} - z_{2} + I_{Ca} - I_{Na})$$

$$\dot{z}_{2} = [z_{2}^{2} - z_{2}^{3} + (z_{2}^{2} - z_{2})a - w_{2}]$$

$$\dot{w}_{2} = \epsilon_{Ca}(z_{2} - I_{Ca} - \mu^{*})$$
(2)

Note that ION₂ has no input. These IONs (with I_{ext1} =0) exhibit limit cycle oscillations as well as bursting phenomenon for a set of values of μ^* and a [13]. However these oscillations are not necessarily in phase. We are interested in designing a simple linear feedback control law such that in the closed-loop system, the state vector of ION₁ asymptotically tracks the state vector of the reference ION. Furthermore, it is assumed that only the output error signal $(z_1 - z_2)$ is measured for feedback.

3. Synchronizing Control System

In this section, for the synchronization of the IONs, a linear feedback control law is designed. Define state vectors $x_1 = (u_1, v_1, z_1, w_1)^T \in \mathbb{R}^4$ and $x_2 = (u_2, v_2, z_2, w_2)^T \in \mathbb{R}^4$. Then (1) and (2), can be compactly written as

$$\dot{x}_1 = f(x_1) + BI_{ext1}$$
$$\dot{x}_2 = f(x_2) \tag{3}$$

where the nonlinear vector function $f(x_1) \in R^4$ and $f(x_2) \in R^4$ are easily obtained from (1) and (2), and one has $B = [0, 0, 0, -\epsilon_{Ca}]^T$.

Let $e = x_1 - x_2$ be the state vector error of the two IONs. Then using (3), the dynamics of the error are given by

$$\dot{e} = f(x_1) - f(x_2) + BI_{ext1} \tag{4}$$

Expanding $f(x_1) = f(x_2 + e)$ about x_2 gives

$$\dot{e} = f(x_2) + \frac{\partial f(x_2)}{\partial x_1} e + BI_{ext1} - f(x_2) + h.o.t$$
 (5)

where h.o.t. denotes higher-order terms in e. For small e, (5) can be approximated by the variational equation of the form

$$\dot{e} = A(t)e + BI_{ext1} \tag{6}$$

where

$$A(t) = \frac{\partial f(x_2(t))}{\partial x_1}$$

is the Jacobian matrix, evaluated along the trajectory $x_2(t)$ of ION₂. It easily follows that the matrix A(t) is

$$\begin{bmatrix} \alpha & -\frac{k}{\epsilon_{Na}} & 0 & 0\\ k & 0 & -k & 0\\ 0 & 0 & \beta & -1\\ 0 & 0 & \epsilon_{Ca} & 0 \end{bmatrix}$$
(7)

where $\alpha = \frac{k}{\epsilon_{Na}} [2u_2 - 3u_2^2 + (2u_2 - 1)a]$ and $\beta = 2z_2 - 3z_2^2 + (2z_2 - 1)a.$

Let us assume that the reference ION is undergoing a limit cycle oscillation. As such $x_2(t)$ is a periodic trajectory. Suppose that the period of $x_2(t)$ is T_p . Then the matrix A(t) is also periodic and its period is T_p . For the synchrony of the two IONs, it is essential to design a control system such that the state vector error e(t) converges to zero. Let us select a control signal of the form

$$I_{ext1} = g(z_1 - z_2)$$
(8)

where g is a feedback gain (yet to be determined). The closed-loop error system is

$$\dot{e} = f(x_1) - f(x_2) + gBc(x_1 - x_2) \doteq f_e(e, t)$$
(9)

where c = [0, 0, 1, 0], and the argument t indicates the dependence of f_e on $x_2(t)$. Substituting the control input (8) in (6), gives the variational equation

$$\dot{e} = (A(t) + gBc)e \doteq A_c(t)e \tag{10}$$

The matrix $A_c(t)$ is also periodic.

First consider stability of the equilibrium point e = 0of the variational equation (10). For this purpose, let us compute the transition matrix of $A_c(t)$ by solving the matrix differential equation

$$\dot{\Phi}(t,t_0) = A_c(t)\Phi(t,t_0); \Phi(t_0,t_0) = I_{4\times 4}$$
(11)

where I denotes an identity matrix of indicated dimension. The growth property of $\Phi(t, t_0)$ depends on the characteristic multipliers (eigenvalues of $\Phi(T_p, 0)$). For the asymptotic stability of the origin of (10), the characteristic multipliers must be strictly within the unit disc [18]. Note that $A_c(t)$ is a function of the gain g, and its characteristic multipliers depend on it. It will be seen in the next section that there exists a set of values of the gain for which asymptotic stability of (10) is assured. Of course, asymptotic stability of the variational equation implies asymptotic stability of e = 0 of the nonlinear time varying system (9). In the next section, a set of values of g are obtained and the performance of the controller is examined.

4. Digital Simulation

The ION parameters given in [13] are used for numerical computation. These parameters are: a = 0.01 and $\mu^* = 0$. The initial conditions of ION₁ and ION₂ are are $x_1(0) = (-0.1, -0.1, 0.1, 0.1)^T$ and $x_2(0) = (-0.2, 0.1, -0.1, 0.3)^T$. Note that the initial conditions of the IONs are not equal. For simulation, the derivatives of \dot{x}_i have been scaled by a factor of 60. The trajectories of the reference ION are shown in Figure 1. We observe that ION₂ is exhibiting LCO, and the limit cycle is orbitally stable. The period of oscillation is $T_p = 1.4948$. In fact, the existence of

orbitally stable limit cycle can be confirmed by computing the characteristic multipliers of $\Phi_o(T_p, 0)$, where $\Phi_o(t, t_0)$ is the transition matrix associated with A(t). The computed values of the eigenvalues of $\Phi_o(T_p, 0)$ are 0, 0.0000, 0.0283 and 1.0053 (\approx 1). Actually, the fourth eigenvalue is exactly one; this small error is due to numerical computation which cannot be avoided. Since one of the characteristic multipliers for the reference ION is unity and remaining are within the unit disk, the periodic orbit is orbitally stable. Note that e = 0 of the variational equation $\dot{e} = A(t)e$ cannot be asymptotically stable. Therefore, for the synchronization of the two IONs, feedback signal I_{ext1} is essential.

For obtaining the values of the gain for the stabilization of the error dynamics, the transition matrix $\Phi(T_p, 0)$ associated with the closed-loop matrix $A_c(t)$ and its eigenvalues are computed. Figure 2 shows the magnitude plot of four eigenvalues. It is seen that for the gain $g \in [-20, -5]$, all the eigenvalues of $\Phi(T_p, 0)$ are within the unit disk; and therefore, for any value of the gain in this set, the equilibrium point e = 0 of the closed-loop variational equation is exponentially stable.

For the selected value of the gain g = -10, the closedloop system including (1) and the control law (8), and the ION₂ model (2) are simulated. The trajectories of both the IONs and the time history of the eigenvalues of $\Phi(T_p, 0)$ are shown in Figure 3. It is observed that for g = -10, the eigenvalues of $\Phi(T_p, 0)$ have magnitude less than 1. Therefore, e = 0 of the nonlinear error dynamics (9) with g = -10 is exponentially stable. As predicted, we observe that the error vector e(t) asymptotically converges to zero; and therefore, the controller quickly accomplishes synchrony of the IONs. In the steady-state control input vanishes.

5. Hardware Implementation

This work also presents implementation of the proposed ION synchronization and the ION (reference ION) in hardware using a Hardware-Software (HW-SW) co-simulation methodology. The Xilinx System Generator (XSG) for DSP is an add-on module for the Mathworks Simulink software. The ION model and the proposed loopback control are designed and implemented using Xilinx System Generator for DSP using predefined device optimized DSP blocksets and custom HDL codes. The tool generates synthesizable HDL code that is mapped in the Xilinx reconfigurable chip (FPGA). The developed system represents a bit-accurate and cycle-accurate model of the software simulation model. The advantages of using System Generator for DSP can be summarized as follows: 1) Rapid prototyping of complex and high-performance DSP systems from high-level abstraction, 2) Bit and cycle accurate floating and fixed point implementation of the DSP algorithm, 3) Automatic generation of HDL code for synthesis to reconfigurable devices, 4) Support for hardware-software cosimulation, and 5) Support of Xilinx tools to explore power, latency and device utilization of the developed implementation.

In our implementation the Xilinx System Generator for DSP 12.1 was used to generated HDL codes that where synthesized using Xilinx ISE 12.1 Design Suite. A XSG implementation of the ION and the top-level implementation of the reference ION and the controller are shown in Figures 6 and 7 respectively. The design was implemented on the XtremeDSP development platform that contains a Spartan-3A DSP 3400A Xilinx FPGA. The Spartan-3A DSP Development Platform provides a great environment for developing signal processing designs. The Xilinx partan-3A DSP 3400A FPGA is based on the XtremeDSP DSP48A Slice. The 250 MHz DSP48A Slice provides an 18-bit x 18-bit multiplier, 18-bit pre-adder, 48-bit post-adder/accumulator, and cascade capabilities for various DSP applications. The DSP48A slice also support a wide math functions, DSP filters, and complex arithmetic without the use of general FPGA fabric.

6. Results

The paper accomplishes the following research objectives: 1) An ION model has been discussed, simulated and implemented in reconfigurable hardware. The Matlab/Simulink simulation results and HW-SW cosimulation results are shown in Figures 1 and 4, 2) A linear feedback control law to synchronize two IONs has been developed, simulated and implemented in reconfigurable hardware. The Matlab/Simulink simulation results and HW-SW cosimulation results are shown in Figures 3 and 5, 3) Table 1, shows a detailed hardware resource usage for our implementations, and 4) Table 1, also shows the maximum latency of the hardware implementation and the maximum frequency of operation of the implemented hardware.

Table 1: Hardware resource usage for reference ION and the ION synchronizer

| # | ION | ION Synchronizer |
|------------|------------|------------------|
| Slices | 734 | 1467 |
| Flip-Flops | 217 | 990 |
| LUTs | 922 | 1885 |
| IOB | 65 | 145 |
| DSP Slices | 10 | 20 |
| BUFGMUX | 1 | 1 |
| L_max | 39.179 ns | 41.739 ns |
| f_max | 25.524 MHz | 23.958 MHz |

7. Conclusions

This paper implements a hardware closed-loop control system for synchronization of inferior olive neurons (IONs). Each ION is modeled by a set of four non-linear differential equation. The model of the ION has be simulated and implemented in reconfigurable hardware and the results verified



Fig. 1: Trajectories of the reference ION - Simulation



Fig. 2: Magnitude plot of four eigenvalues - Simulation

using HW-SW cosimulation. The software simulation and HW-SW co-simulation results are identical. The closed-loop control system containing two IONs and a simple linear feedback control law has been implemented. The object of controlling one ION so that both oscillate in unison has been achieved. The closed-loop control system has been implemented in reconfigurable hardware and the software simulation and HW-SW cosimulation results are identical. The performance of the control system for hardware latency and device utilization has been evaluated and discussed. Typically an array or network of such IONs are required and there exists a need to develop controllers for ION synchronicity. In future work, the problem of synchronizing an array of IONs and Hardware in Loop (HIL) testing of the ION with motors will be addressed.



Fig. 3: Trajectories of both the IONs and the time history of the eigenvalues of $\Phi(T_p, 0)$ -Simulation



Fig. 4: Trajectories of the reference ION - HW-SW Co-simulation



Fig. 5: Trajectories of both the IONs - HW-SW Co-simulation



Fig. 6: XSG implementation of the reference ION



Fig. 7: XSG implementation of the synchronous ION system

References

- [1] Armstrong DM (1974) Functional significance of connections of the interior olive. Physiological Reviews 54(2):358-417
- [2] Bandyopadhyay PR, Singh SN, Thivierge DP, Annaswamy AM, Leinhos HA, Fredette AR, Beal DN (2008) Synchronization of animalinspired multiple high-lift fins in an underwater vehicle using olivocerebellar dynamics. IEEE Journal of Oceanic Engineering 33(4): 563-578
- [3] Benardo LS, Foster RE (1986) Oscillatory behavior in inferior olive neurons: Mechanism, modulation, cell aggregates. Brain Research Bulletin 17(6):773-784
- [4] Bennett MV, Zukin RS (2004) Electrical coupling and neuronal synchronization in the mammalian brain. Neuron 41(4): 495-511
- [5] Guckenheimer J, Labouriau IS (1993) Bifurcation of the Hodgkin and Huxley equations: a new twist. Bulletin of Mathematical Biology 55(5):937-952.

- [6] Ito M (1984) Cerebellum and Neural control. Raven, New York
- [7] Izhikevich EM (2000) Neural excitability, spiking, and bursting. International Journal of Bifurcation and Chaos 10(6):1171-1266.
- [8] Izhikevich EM (2007) Dynamical systems in neuroscience: the geometry of excitability and bursting. MIT Press, MA
- [9] Katori Y, Lang EJ, Onizuka M, Kawato M, Aihara K (2010) Quantitative modeling of spatio-temporal dynamics of inferior olive neurons with a simple conductance-based model. International Journal of Bifurcation and Chaos 20(3):583-603.
- [10] Kazantsev VB, Nekorkin VI, Makarenko VI, Llinas R (2004) Selfreferential phase reset based on inferior olive oscillator dynamics. PNAS 101(52):18183-18188
- [11] Kazantsev VB, Nekorkin VI, Makarenko VI, Llinas R (2003) Olivo-cerebellar cluster-based universal control system. PNAS 100(22):13064-13068
- [12] Lampl I (1994) Characterization of subthreshold oscillations in inferior olivary neurons: mechanisms of generation, mode of operation and their functional role. PhD thesis: Hebrew University, Israel.
- [13] Lee KW, Singh SN (2009) Adaptive global synchrony of inferior olive neurons. Bioinspiration and Biomimetics 4(3):(13 pages).
- [14] Lee KW, Singh SN (2012) Bifurcation of orbits and synchrony in inferior olive neurons. Journal of Mathematical Biology (To appear).
- [15] Llinas RR (2009) Inferior olive oscillation as the temporal basis for motricity and oscillatory reset as the basis for motor error correction. Neuroscience 162(3) : 797-804
- [16] Llinas RR, Leznik E, Makarenko VI (2004) The olivo-cerebellar circuit as a universal motor control system. IEEE Journal of Oceanic Engineering 29(3):631-639
- [17] Manor Y, Rinzel J, Segev I, Yarom Y (1997) Low-amplitude oscillations in the inferior olive: a model based on electrical coupling of neurons with heterogeneous channel densities. Journal of Neurophysiology 77(5):2736-2752
- [18] Rugh, WJ (1996) Linear System Theory. Prentice-Hall, NJ
- [19] Sotelo C, Llinas RR, Baker R (1974) Structural study of inferior olivary nucleus of the cat : morphological correlates of electronic coupling. J. Neurophysiology 37(3):541-559
- [20] Velarde MG, Nekorkin VI, Kazantsev VB, Makarenko VI, Llinas R (2002) Modeling inferior olive neuron dynamics. Neural Networks 15(1):5-10
- [21] Welsh JP, Llinas RR (1997) Some organizing principles for the control of movement based on olivocerebellar physiology. Prog. Brain Res. 114:449-461
- [22] Wang, W., Slotine, JE (2005) On partial contraction analysis for coupled nonlinear oscillators. Biological Cybernetics. 92: 38-53. Prentice-Hall, NJ

A Novel Temporal Framework for Relational Event Representation

Sadi Evren SEKER

Department of Computer Engineering, Istanbul University, Istanbul, Turkey academic@sadievrenseker.com

Abstract: Temporal logics are widely used in many study types, such as question answering, ontology, natural language processing, search engines, text summarization, or even visual tools like Gantt charts or UML diagrams.

Computable temporal languages are the logical systems, built on temporal logic, that can be computed to find a result. They can also be defined as computer-computable languages, built on temporal logics. All timeline drawing or planning software uses temporal logic in order to visualize or process cases.

Also, semantic web studies are one of implementation areas where the temporal modeling and reasoning is massively needed. Relation between events or event types and event of subjects can be modeled by using temporal logic.

This study introduces the temporal logics and the computable temporal languages in the current literature. For the first time, some of temporal logic problems are pointed and solved during this study.

Also a novel temporal framework is implemented with JAVA and published on the web which covers the solutions of temporal logic problems.

1. Introduction

The aim of this study is to achieve a machine-computable temporal logic and apply this novel logic into a software implementation. The main motivation behind the new temporal logic comes from the natural language modeling of the temporal statements.

For some natural languages, some studies already declare temporal differences as affected by their culture and therefore their languages.

In this study, some temporal difficulties and their solutions are discussed and an ensemble of solutions is modeled in a reasonable logical model. This paper is organized to start with some background information about two temporal logics, Reichenbach and Allen's Temporal logics and than points some difficulties in the modeling of temporal relations. Finally these problems are solved and a visualization tool is introduced which is implemented over the new suggestions as the proof of solutions are possible to be coded in programming languages.

1.1 Reichenbach Temporal Logic

Reichenbach temporal logic is built on simple three temporal anchors:

- Speech time (symbolized by S)
- Reference time (symbolized by R)
- Event time (symbolized by E)

Most of Reichenbach's study was focused on the natural languages. Thus, he formulated the order of these times.

For example, a sentence like "I read the book" can be formalized as R=E<S. On the other hand, a sentence like "I have read the book" can be formalized as E<R=S. Please note that in the former model, the event takes place before the speech time and the speech refers to the event time, so the event and reference times are equal and smaller than speech time on the model. For the latter example, the event again takes place before the speech time, but the speech is referring to the current time, so the speech time and reference times are equal and greater than the event time.

By a simple probability calculation, we can end up with 13 possible orders for the above temporal anchors. Obviously, not all of these probabilities are meaningful in a natural language. Reichenbach has named these possibilities, using Anterior, Simple and Posterior aspects, and Past, Present and Future tenses. In Reichenbach's opinions, there can only be 9 possible meaningful times in English or in any natural language. Table 1 covers these possibilities and samples for each of the case:

| Perm | RTN | English Tense | Sample |
|--|----------------------|--------------------|---------------------------|
| E <r<s< td=""><td>Anterior past</td><td>Past perfect</td><td>I had slept</td></r<s<> | Anterior past | Past perfect | I had slept |
| E=R <s< td=""><td>Simple past</td><td>Simple past</td><td>I slept</td></s<> | Simple past | Simple past | I slept |
| R <e<s< td=""><td></td><td></td><td></td></e<s<> | | | |
| R <s=e< td=""><td>Posterior past</td><td></td><td>I would sleep</td></s=e<> | Posterior past | | I would sleep |
| R <s<e< td=""><td></td><td></td><td></td></s<e<> | | | |
| E <s=r< td=""><td>Anterior present</td><td>Present perfect</td><td>I have slept</td></s=r<> | Anterior present | Present perfect | I have slept |
| S=R=E | Simple present | Simple present | I sleep |
| S=R <e< td=""><td>Posterior present</td><td>Simple future</td><td>I will sleep</td></e<> | Posterior present | Simple future | I will sleep |
| S <e<r< td=""><td></td><td></td><td></td></e<r<> | | | |
| S=E <r< td=""><td>Anterior future</td><td>Future perfect</td><td>I will have slept</td></r<> | Anterior future | Future perfect | I will have slept |
| E <s<r< td=""><td></td><td></td><td></td></s<r<> | | | |
| S <r=e< td=""><td>Simple future</td><td>Simple future</td><td>I will sleep</td></r=e<> | Simple future | Simple future | I will sleep |
| S <r<e< td=""><td>Posterior future</td><td></td><td>I shall be going to sleep</td></r<e<> | Posterior future | | I shall be going to sleep |

• TABLE 1. ALL POSSIBLE 13 PERMUTATIONS OF Reichenbach's temporal logic, their English tenses/aspects, and a sample of each case.

Please note that in Table 1, blank lines represents meaningless cases of the permutations in English and RTN stands for the Reichenbach Tense Name.

2. ALLEN'S INTERVAL LOGIC

Allen's Interval Logic (AIL) or Allen's Temporal Logic (ATL) deals with orders of events. The representation of event orders like "event A is before event B" or "event A is at the same time with event B" are the operators of this logic [4].

The basic variables in AIL are the intervals, and Allen has built his logic over binary operators working on those intervals. In AIL, 13 basic binary operators connect intervals by constraints. These intervals can be considered to be running threads or any operations on the timeline.



Fig. 1. Linearity of timeline

Linearity of time can vary between different temporal logics. Figure 1 displays four different types of timeline linearity. The figure 1-1 is linear in the past and in the future. Figure 1-2 is only linear in the past and is non-linear in the future. This type of linearity can be classified as semi-linear. Figure 1-3 is the opposite form of figure 1-2, where the future is linear and the past is non-linear. Figure 1-4 is non-linear in both the past and future.

AIL supports all types of linearity in figure 1. AIL can be demonstrated by a box diagram, where the boxes represent the events and the arrows represent the relations between them.

For example, in an example sentence like "John ate an apple at the table after he entered the room" we have the events "eat" and "enter." There are also hidden events, in which John goes to the table and takes the apple, in order to eat an apple from the table after he has entered the room. The timeline of the example is given in figure 2.



Fig. 2. Sequence of events in the sample sentence

If the states of the events are considered, we know John was outside of the room before he entered the room. He was also away from the table before he approached the table, and he had no apples before taking the apple.





room, or that John had no apple while he was away from the table or while he was outside of the room.

In AIL, there are 13 types of possible relations. The list below includes the possible operators of Allen Interval Logic:

- Before (x,y) or After (y,x)
- Overlaps (x,y) or Overlapped (y,x)
- Meets (x,y) or MetBy (y,x)
- Contains (x,y) or During (x,y)
- Starts (x,y) or StartedBy(y,x)
- Ends (x,y) or EndedBy (y,x)
- Equals (x,y)

The above sample can be modeled using AIL. Let's say entering room (ER) requires us to be outside of the room (OR); after entering the room, the state is inside the room (IR) and similarly, approaching the table (AT) changes the state of being away from the table (SAT) to the state of being close to the table (SCT). Taking the apple (TA) is a transformation of state from not having the apple (NHA) to having the apple (HA). All these states are pre-requirements in the case of eating the apple (EA).

The above sentence can be modeled in Allen Temporal Logic in the following formulation:

Meets(OR,ER) \land Meets(ER,IR) \land During (ER,SAT) \land During (AT,IR) \land Meets(SAT,AT) \land Meets(AT,SCT) \land During(AT,NHA) \land During(TA,IR) \land During (TA,SCT) \land Meets(NHA,TA) \land Meets(TA,HA) \land During(EA,HA) \land During (EA,CT) \land During(EA,IR) \land Meets(TA,EA)

The formulation above demonstrates all the temporal states and events in the sample sentence. On the other hand, a reader can interpret the above sentence and can add more states, which can still be modeled by AIL. For example, if John follows the order of events above when he is hungry, then this state can be added to the model of AIL. In this case the model would be:

Occurs (hungry, NHA) \land Holds (hungry, TA) \land Meets (hungry, EA)

So, from the AIL model, John eats an apple when he gets hungry and does not have an apple; he takes an apple while he is hungry; his state of hunger ends when eating the apple.

3. Missing Temporal Models in AIL and Reichenbach Temporal Logic

Unfortunately, Allen's Interval Logic is not sufficient for a representation of Turkish temporal logic [1]. One of the specific problems for Turkish temporal logic is the positive/negative verb repetition. These terms in Turkish represent a continuous event by using two verbs with opposite meanings. For example, in English a single verb like "blink" is represented in Turkish with two separate verbs "yanıp sönmek" (to light and to fade). This concept can also be represented by "to flash" or "to twinkle" in English – both single verbs. Another example is the translation of the term "pacing up and down" or to "pace back and forth" in Turkish.

Again this term is represented by two separate verbs "gelip gitmek" (to come and to go). Or another example: "restart" in Turkish is "kapatip açmak" (to close and to open).

The examples above show a group of terms in Turkish where ATL is insufficient, because the precise order of words in Turkish temporal logic is not important. For example, the term for "restart" in Turkish, "kapatıp açmak" (to close and to open), can also be represented as "açıp kapatmak" (to open and to close), which is exactly the reverse order of the former term. On the other hand, the semantic representation of these terms is only one event in the ATL, which creates a problem in the case of trying to represent a single event with two separate verbs.

Allen's Temporal Logic is a linear logic that is suitable for representing events in a linear manner. Unfortunately, the temporal logic behind Turkish natural language is not exactly linear. Although there are some studies that model time in a non-linear domain [6], TimeML has been implemented linearly using ATL. For example, let's try to represent the Turkish sentence below in ATL.

"The life signal on the safety buoy was blinking while the divers were under water."

The above English sentence can easily be represented in ATL, as shown in Figure 4.

| A | Diver's out water | Diving | Diver's under w | ater | |
|---|-----------------------|--------|-----------------|-----------|---------------------|
| 8 | Diver's out water | | | Diving | Diver's under water |
| c | Life signal is not bl | inking | Start Blinking | Life Sign | al is blinking |

Fig. 4. Sequence of events in sample sentence

Since we have no idea of the starting time of the blinking of the life signal, either A-C or B-C or any time combination of "diving" and "start blinking" between these times is considered as correct from the above input sentence.

The ATL representation of the above case would be shown as below:

Meets(SB,DUW) ^ During (DUW,LSB) ^ During (DOW,LSNB)

where, "sb: signal blinking", "duw: divers under water", "lsb: life signal blinking", "dow: divers out water", "lsnb: life signal is not blinking".

In the Turkish translation of above sentence, the representation would be as shown in Figure 5.

| A Diver's out water | Diving | Diver's under s | water | |
|-------------------------|---------|-----------------|--------|---------------------|
| B Diver's out water | | | Diving | Diver's under water |
| C Life signal is not bl | linking | To light | n Ma | NA A |
| L | | Te fade | 2020 | 00 O |

Fig. 5. Suggestion for Recurring events on Allen's Temporal Logic

The temporal logic behind Turkish natural language states that even when the event starts with lighting or fading, these two events follow each other and continue while the diver's under the water. This logic cannot be stated in ATL.



Fig. 6. Addition of recurring events to ATL

Figure 6 visualizes all the relations in ATL. The original figure [12], holding only the first 13 relations, does not cover the last row, which indicates the relation type "recurs". This row is added to visualize one of the results of the study on ATL.

One representation missing in TimeML is that of selfreferencing cases. This uncovered case in TimeML, is not unique to Turkish, and the same problem can occur in any language. The speaker can refer to the current speech. For example, the case below is a self-referring case:

"My current talk is about computer science."

In the above sentence, the speaker refers to the current talk, which is the talk itself. This case creates a self-reference which is not covered in ATL, and therefore not in TimeML either. A solution would be to add this reference at the signal level of TimeML, but the signal level does not hold information regarding the relationships betwee events. We suggest a better solution, which is adding these cases into the ATL level as a new relation type, as shown in figure 7.



Fig. 7. Self referring events

The self-reference problem occurs in TimeML and this problem is neither a Reichenbach nor ATL level problem, since Reichenbach doesn't relate to the relation between two or more events and ATL doesn't relate to the time of reference.

Another piece missing from TimeML is the modeling of absence. At first glance, this concept could be created by using the negative of omnipresence. However, this does not translate exactly. The difference between the two is explained by the following examples.

" -- When do you play tennis?

-- Never!"

The above dialogue is answered by the word "never," indicating that the event will never occur. Another example, with a slight semantic difference, would be:

"I never take notes when I attend classes."

This example is slightly different than the first example, because in the second case, the word "never" refers to the negative of an event that is addressable in temporal logic.



Fig. 8 Relation of a negative event

Figure 8 visualizes the relation between the events "Attending classes" and "Take notes," where the event "Take notes" is negative.

This case can be interpreted as "I don't take notes when I attend classes." Thus we can interpret that this event is connected to the "attending of classes" by the "during" operator with a negative perspective. On the other hand, the first example cannot be connected to any other event. Another problem for the first case is the demonstration of the event in a timeline. If an event never exists, the demonstration of the event is also impossible. We suggest the solution of adding a nonexistence to the environment, and connecting the event "never exist" to the nonexistence. This solution is also useful for the events connected by an order operator to nonexistence.

For example, the below sentence contains an event connected to nonexistence.

"Nothing existed before the Big Bang."

In this phrase, the word "nothing" indicates the nonexistence, and the event "Big Bang" occurs after "nothing," so we can conclude that the Big Bang is connected to nonexistence by the "next" operator in temporal logic.

We can demonstrate this solution with a slightly more complicated example:

"Nothing comes from nothing" (Parmanides)..



Fig. 9 Interpretation of "nothing comes from nothing"

Figure 9 visualizes the famous quotation by Parmanides, where one of the "nothings" is in the existence domain and the other is in the non-existence domain.

As a solution, TimeML should contain a non-existence domain of events.

Another problem is the linearity of time. In figure 1, when the linearity of temporal logic is explained, the linearity was considered as a forking of events on time. This concept is called multi-linear time, where a timeline can flow through any of the possible paths.

4. Development of Software

In order to demonstrate that the above solutions are implementable in computer software, a new project is developed in JAVA with a user interface to model the temporal relations indicated above.

| Palette | |
|--|--|
| | e reger () |
| | CORDER Educations 28 |
| Query Panel | rell 😸 (Annaldesse |
| Trst box | (of2 convectional) (brassing) (brassing) (brassing) |
| iecond box | Everil borilitor |
| | |
| | recuribson (sillers |
| roperties | |
| same incessional | |
| the second second second second second | |
| re an | |
| 1 AVECA | D-0-0 |
| 1 (VIC) | 0-0-0 0 0 |
| 1 (V.35) 20 . Adth 30 | 0-0-0 0 0-0-0 |
| a webs I 70 f. Adth 40 f | 0,0-0 0 0 0-0 |
| a webs I 70 5 Width 40 5 | a, o-o o a o-o-o |
| a vezos k zoli width do leight zoli | a_o.o. a a o.o.o |
| a vyzov k zo - vidih | a_o.o a a o-o-o |
| 1 vezes zo i ndth 20 i eight 26 y ype Pelenrei by 2 | 0,0-0 0 0-0-0 |

Fig. 10 Screenshot of software interface with a challenging problem

The temporal framework implemented, makes possible to create all the discussed problems above like self reference, negative events, or recurring events besides the current relation models of AIL or Reichenbach Temporal Logic.

Software is downloadable from www.sadievrenseker.com/temporalframework web site.

Acknowledgement

This study was supported by Scientific Research Projects Coordination Unit of Istanbul University. Project number 16339 and project number YADOP-16728.

REFERENCES

- [1] Şeker, Ş.E., Diri, B., "TimeML and Turkish Temporal Logic", International Conference on Artificial Intelligence,ICAI'10, 2010,
- [2] Hacioglu, K., Chen, Y. and Douglas, B., "Automatic Time Expression Labeling for English and Chinese Text", Computational Linguistics and Intelligent Text Processing, pages 548-559, ISSN 0302-9743, 2005.
- [3] Robust temporal processing of news, Proceedings of the 38th Annual Meeting on Association for Computational Linguistics, Inderjeet Mani, George Wilson Pages: 69 -76, 2000
- [4] Allen Linear (Interval) Temporal Logic –Translation to LTL and Monitor Synthesis– Grigore Ro,sul _ and Saddek Bensalem, CAV'06, LNCS 4144, pp 263-277, 2006
- [5] Reichenbach, H., Elements of Symbolic Logic, New York: Macmillan (1947)
- [6] Frank D. Anger, Debasis Mitra, Rita V. Rodriguez, "Satisfiability in Nonlinear Time: Algorithms and Complexity", Proceedings of the Twelfth International FLAIRS conference, AAAI, 1999

- [7] Ozkirimli A, Türk Dili: Dil ve Anlatım, İstanbul Bilgi University Press, 2001.
- [8] TimeML, TERQAS, (2002), TimeML has been developed in the context of three workshops starting from 2002.
- [9] "TimeML Annotation Guidelines Version 1.2.1" Roser Saur'1, Jessica Littman, Bob Knippen, Robert Gaizauskas, Andrea Setzer, and James Pustejovsky (January 31, 2006)
- [10] Natalia Kotsyba Using Petri nets for temporal information visualization Études Cognitives/Studia Kognitywne, CEEOL (2006)
- [11] Andr'e Bittar, Annotation of Events and Temporal Expressions in French Texts Proceedings of the Third Linguistic Annotation Workshop, ACL-IJCNLP 2009, pages 48–51, Suntec, Singapore, 6-7 August 2009. c 2009 ACL and AFNLP
- [12] Temporal Logics for Real-Time System Specification P. BELLINI, R. MATTOLINI, and P. NESI University of Florence
- [13] Christian Kissig and Laura Rimell, Closing TLink-Relations (Reasoning with Intervals) June 23, 2005

Interactive Computer Program: Packaging DNA into Chromosomes

Xiaoli Yang¹, Yifan Cai¹ and Charles Tseng²

¹Department of Electrical and Computer Engineering ²Department of Biological Sciences Purdue University Calumet Hammond, IN, USA

Abstract - As part of the interactive program for teaching and learning genetics, the module on packaging DNA into chromosomes involves the simultaneous coordination of eyes, mind, and hands for visualization, cognitive feedback, and manipulation, respectively. Computer modeling of various chromatin structures during packaging is based on OpenTK-OpenGL on .Net Platform, which is coupled with an inquiry based content design to enhance the efficiency of teaching and learning. The prototype has been successfully tested in a genetics class at Purdue University Calumet. It should also be applicable to a number of undergraduate biology courses.

Keywords: DNA, Chromosomes, Modeling, Computer Program

1 Introduction

From its central role in real-life forensic investigations to being the basis of major biotechnological applications in medicine, agriculture, and the environment, DNA based genetics is an essential discipline in the life sciences. As fascinating as the subject is, however, teaching and learning genetics has often been fraught with difficulty [1-3]. Confronted with intricate molecular structures, complex packaging schemes, and elaborate mechanisms of action, both teacher and student are frequently at a loss – the teacher in how to convey this material in a clear and understandable way, and the student in how to assimilate all the information usefully. To be sure, the abstract and intangible nature of much of the material is the source of the problem.

Traditional methods of teaching genetics, employing classroom lectures, textbook readings, homework assignments, and laboratory exercises, have not proven to be very effective [4, 5]. Recently, efforts have been made to integrate computer visualization technologies into pedagogy to enhance the learning process [6-8]. Current computer-based tools, however, do not stress cognitive feedback in their designs. The present paper describes an innovative approach to teaching and learning genetics, in which students can visualize a real-time, interactive DNA model, as well as actively control the dynamic process of packaging DNA into a compact metaphase chromosome.

The objectives of the program are to 1) develop, as part of a web-based interactive program, a DNA packaging module suitable for a wide range of college courses and 2) serve as a model for STEM (science, technology, engineering, and mathematics) education via distance learning.

2 Model Development

Models of various structures were developed based on the following system: OpenTK-OpenGL on .Net Platform. The Open Tool Kit (OpenTK) is a free project that allows developers to use OpenGL, OpenGL|ES, OpenCL, and OpenAL APIs from a managed language (e.g. VB.NET). Features include:

- Written in cross-platform C# and usable by all managed languages (F#, Boo, VB.Net, C++/CLI).
- Consistent, strongly typed bindings, suitable for RAD development.
- Usable standing alone or integrated with Windows.Forms, GTK#, and WPF.
- Cross-platform binaries that are portable on .Net and Mono without recompilation.
- Wide platform support: Windows, Linux, and Mac OS X, with iPhone port in process.

2.1 3D model of double helix DNA

DNA is modeled as a double helix. The model is specified by l, the length of the helix, r, the radius of the helix, and w and h, the width and the thickness of one strand of the double helix, respectively (Fig. 1). These parameters generate a group of points, which are used to construct the DNA model. The points are linked together to form a sketch of the double helix. After shading the sketch, a 3D DNA model is created. The double helix model is calculated during runtime based on the equations below:

$$f(x) = t, 0 \le t \le l$$

$$f(y) = r * \sin(\theta + t * \alpha)$$

$$f(z) = r * \cos(\theta + t * \alpha)$$
(1)

Where *t* is the length variable along x-axis, *r* the radius of the helix. α the angle increment, controlling the smoothness of the helix. We chose $\alpha = 6$ from the experiment to make the model smooth. θ determines the initial angle. We chose $\theta = 0$ and 45 from the experiment to generate double helical shapes. $t * \alpha \in [0, 2\pi)$. From the above equations,

f(x), f(y), f(z) determine the Cartesian coordinates x, y and z in 3D space.

By calculating the position of all the points, a helical line can be generated (Fig. 2, left). The quadrupling of the line (Fig. 2, right) is generated by replicating the original line four times.



Fig. 1. Parameters for a helix (perpendicular views)



Fig. 2. Left: helical line; Right: doubling of line

After further duplicating the strand with a different θ value and shading the sketch, two DNA strands of different colors are created (Fig. 3).



Fig. 3. DNA double helix structure with shading

2.2 3D model of nucleotide bases

We use line segments (cuboid) of different colors to represent DNA bases. The points that form the DNA strands (determined above) are used to calculate the points representing the line segments (bases). Assume that p1 (x1, y1, z1) and p2 (x2, y2, z2) are corresponding points on different strands generated by the same *t* value, but different θ values, p3 (x3, y3, z3) and p4 (x4, y4, z4) are the points next to p1 and p2, respectively, *w* is the length of the side, pc1 is the midpoint between p1 and p2, and pc2 is the midpoint between p3 and p4 (Fig. 4).



Fig. 4. Parameters used to calculate the position and shape of bases.

Then all 8 points needed to describe a cuboid can be calculated as follows.

$$\begin{cases} (x1 \pm \frac{w}{2}, y1, z1) \\ (x3 \pm \frac{w}{2}, y3, z3) \\ \left(\frac{x1+x2}{2} \pm \frac{w}{2}, \frac{y1+y2}{2}, \frac{z1+z2}{2}\right) \\ \left(\frac{x3+x4}{2} \pm \frac{w}{2}, \frac{y3+y4}{2}, \frac{z3+z4}{2}\right) \end{cases}$$
(2)

The result is a complete DNA model (Fig. 5).



Fig. 5. Screen snapshot of DNA model from the program

2.3 3D model of histone octomer

The histone octomer is represented by an elongated ball, which is described by the following equations:

$$f(x) = r * \sin \varphi * \cos \theta$$

$$f(y) = \begin{cases} r * \cos(\varphi) - \frac{h}{2}, \varphi < 0\\ r * \cos(\varphi) + \frac{h}{2}, \varphi \ge 0\\ f(z) = r * \sin \varphi * \sin \theta \end{cases}$$
(3)

f(x), f(y), f(z) determine the Cartesian coordinate x, y and z in a 3 dimensional space, where r is the radius of the histone; φ the angle between the diameter and z-axis, $\varphi \in \left[-\frac{\pi}{2}, \frac{\pi}{2}\right]$; θ the angle between the projection of the diameter on the plane and the x-axis. $, \theta \in [0, 2\pi)$ and *h* the height of the elongated ball (Fig. 6).



Fig. 6. Sphere coordinate system

These points generate a sketch of the histone octomer. Shading the sketch with a color completes the model (Fig. 7).



Fig. 7. Left: sketch of histone octomer; Right: shaded histone octomer

2.4 DNA wrapping-formation of core nucleosome

In this step, DNA is simplified as a line, which can be wrapped around the histone octomer. Fig. 8 shows how to calculate the position of the binding points.



Fig. 8. Calculation for DNA-protein binding position

The histone octomer is projected on the *xz*-plane as a circle. Assume that *O* is the center of the circle, *P* is outside the circle, *r* is the radius of the circle, *dx* and *dy* are the differences between *O* and *P* in *x* and *y* components, respectively, *D* is the distance between *P* and *O*, *Pb* is the binding point, α is the angle between line *PO* and the vertical line, and β is the angle between line *PO* and the line *PbO*. Then α and β can be calculated as follows:

$$\alpha = \arcsin(\frac{r}{D})$$

$$\beta = \arctan(\frac{dy}{dx})$$
(4)

The value of θ is

 $\theta = 270^{\circ} - \alpha - \beta$ Finally, *Pb* is represented by (*x*,*y*), where

$$x = r * \sin \theta$$

$$y = r * \cos \theta$$
(6)

(5)

When DNA binds to a histone octomer, it starts to wrap around the octomer. After Pb is determined, the points on the spiral nearest to it can be calculated (Fig. 9).



Fig. 9. Wrapping of DNA around histone octomer. Left: sketch; Right: shaded

2.5 Camera position adjustment

In the interactive module, users can view the model from different angles by dragging the mouse. Our program

implemented this function by adjusting camera position when we developed the model. The camera was located on the surface of sphere with the target at the center of sphere, so that the distance between camera and target never changes. In other words, the size of the target remains the same, so that the model size does not change. Camera position (a point on the surface of sphere) is described by ϕ , θ and r, where r is the radius of sphere. While the value of r never changes, ϕ and θ are variable. Changing ϕ and θ changes the camera position, and they are changed by moving the mouse. Moving the mouse produces component values dx and dy along the xand y-axis, respectively. Therefore, mapping dx and dy to ϕ and θ is an effective way to adjust camera position.

3 The Program Content Design

The design of module focuses on inquiry based methods with cognitive feedback and interactive experiences as important components [9]. In every section, a question is followed by observations and measurements, hands-on experiments, and conclusions. In each of the learning steps, dynamic models of DNA molecules, chromatin fibers, and metaphase chromosomes are presented for interaction through visualization, cognition, and operation. Completion of the program requires comprehension of the entire concept and thus ensures the success of the learning experience. The computer-based content is summarized below:

3.1 DNA and chromosomes in prokaryotes and eukaryotes

Inside the cell, DNA molecules are packaged, with helped of proteins, into thread-like structures called chromosomes. In prokaryotes (such as bacteria), the chromosomal DNA, when open, is often circular. The total length of a bacterial chromosomal DNA (e.g., *E. coli* DNA) may be a thousand times longer than the cell that contains it. Little is known about the packaging of bacterial DNA, although a few major DNA regions anchored by proteins at specific sites in the cell have been noted.

In eukaryotes (such as animals and plants), DNA molecules are linear. Each eukaryotic species has a fixed number of chromosomes. For diploid species (species with 2 sets of chromosomes, one from each parent), chromosomes are paired, so that the total number of chromosomes is always even.

In humans, for example, the father provides a set (also called a genome) of 23 different chromosomes (from sperm), while the mother provides the other set (genome) of 23 different chromosomes (from egg). Thus, each of our somatic (body) cells contains 46 (or 23 pairs) chromosomes.

3.2 How long is our DNA?

The DNA of each human genome is about 3.2 billion (3.2 x10⁹) deoxyribonucleotides long. Each deoxyribonucleotide is 3.4 A (0.34 nm), making the total length of human genomic DNA (0.34 nm)(3.2 X 10⁹) \approx 1 m (meter) per genome. Since

there are 2 genomes per human cell, the total length of DNA per human cell is $2 \times 1 \text{ m} = 2 \text{ m}$.

Assuming that an adult human body contains about 50 trillion (50×10^{12}) cells, the total length of DNA in the human body is $(50 \times 10^{12}) 2 \text{ m} = 100 \times 10^{12} \text{ m} (100 \text{ trillion meters of DNA per human})$. The Sun is $150 \times 10^9 \text{ m} (150 \text{ billion meters})$ from Earth. How many times can you stretch the DNA from the Earth to the Sun? $(100 \times 10^{12})/(150 \times 10^9) = 666 \text{ times}$ (Fig. 10). The distance between the Earth and the Moon is about $3.84 \times 10^8 \text{ m}$. Can you calculate the number of times your DNA can stretch from the Earth to the Moon?



Fig. 10. The total length of human DNA

3.3 Can our genomic DNA fit into a nucleus?

Let's examine the size of the human genome in the cell. In the nucleus (G1 phase) of each cell in the human body are 46 chromosomes. Since a set of 23 chromosomes constitutes a genome, there are 2 genomes per nucleus. Each genome contains approximately 3×10^9 base pairs, so there are 6×10^9 base pairs per nucleus in the G1 phase. Each nucleotide of the base pair measures approximately 3.4×10^{-10} m in length.

Therefore, the total length of DNA in each nucleus is: $(2)(3 \times 10^9 \text{ nucleotides})(3.4 \times 10^{-10} \text{ m/nucleotide}) \approx 2 \text{ m}.$

However, the diameter of an average nucleus is 10×10^{-6} m, making the total length of DNA in the nucleus 200,000 times longer than the diameter of the nucleus:

 $(2 \text{ m})/(10 \text{ x}10^{-6} \text{ m}) = 200,000$. How can such a long strand of DNA fit in such a small nucleus?

Cleary, the DNA must be folded. Let us examine how much space the DNA occupies in the nucleus if it is somehow folded so that it will fit. DNA exists in a double helix, which can be approximated by a cylinder with diameter 20×10^{-10} m. Therefore, the volume the DNA occupies is:

 $\pi(r^2)h = (3.14159)(10 \text{ x}10^{-10} \text{ m})(2 \text{ m}) = 6.4 \text{ x} 10^{-18} \text{ m}^3$

Assume that the diameter of an average nucleus is approximately 10×10^{-6} m, making the volume of the spherical nucleus:

 $(4/3)(\pi)(r^3) = (4/3)(3.14159)(5 \times 10^{-6} \text{ m})^3 = 5.24 \times 10^{-16} \text{ m}^3$

Consequently, the fraction of the nucleus occupied by DNA is:

 $[(6.4 \text{ x } 10^{-18} \text{ m}^3) / (5.24 \text{ x } 10^{-16} \text{ m}^3)] \text{ x } 100 = 1.22\%$

There is clearly enough room in the nucleus for the DNA – and for its activities including transcription, replication, packaging and unpackaging. This leads to the our main topic: How is the DNA packaged into chromosomes?

3.4 Levels of DNA packaging

DNA packaging can be considered at the following levels:

<u>Level 1 - Double helix:</u> The double helical DNA molecule has a width of about 20 A (2 nm) (Fig. 11)

Level 2 - 10 nm chromatin fiber: The DNA of eukaryotic cells is tightly bound to basic (positively charged) proteins known as histones. This nucleoprotein complex is called chromatin. The basic structural unit of chromatin is the nucleosome. A nucleosome consists of a small segment of DNA wrapped around histones. The core nucleosome particle consists of two molecules each of the core histones (H2A, H2B, H3, and H4), forming a histone octomer (Fig. 12), around which is wrapped approximately 146 base pairs of DNA. The core particle is stabilized by a fifth histone called H1 (also called a linker histone). DNA between the core nucleosome particles is called linker DNA. The core nucleosome particle plus the linker DNA is about 200 bases long, as evidenced by digestion with the enzyme micrococcal nuclease.



Fig. 11. Screen snapshot of double helical DNA molecule



(blue: H3; green: H4; Yellow: H2A; red: H2B)

The linear chromatin fiber at this level is about 10 nm in diameter (Fig. 13. a). This represents the state of most chromosomal DNA during interphase and is known as euchromatin. Genes that are actively transcribed (with momentary detachment of histones) are in this less condensed state.

Level 3 – 30 nm chromatin fiber:

Some of the 10 nm chromatin fibers can be packed into 30 nm fibers. To do this, the H1 histones, each attached to a core nucleosome, interact with each other, turning inwards and forming a new spiral structure known as a solenoid or 30 nm chromatin fiber (Fig. 13. b). Each 6-8 nucleosomes constitute one turn of the new spiral. In this state, the chromatin is tightly packed and is referred to as heterochromatin, a state in which DNA is genetically inactive (no transcription or replication).

<u>Higher levels of packaging - looping</u>: The above levels of DNA packaging mainly describe the G1 phase of the cell cycle. If the cell is destined for division, the G1 phase is followed by the S phase, where a DNA molecule replicates semiconservatively to form two identical DNA molecules before being packaged into chromatin fibers and entering the G2 phase. In the early stages of mitosis (or meiosis), the two replicated DNA molecules (in the form of chromatin fibers) continue to condense, and the 30 nm fibers are folded into loops (Fig. 14).



Fig. 13 a. 11 nm chromatin fiber



Fig. 13 b. 30 nm chromatin fiber



Fig. 14. Screen snapshot of looped chromatin fibers

The chromosome reaches its highest condensed state at metaphase. Condensation is the result of the further folding of the 30 nm fibers into different loops until a 700 nm structure is reached (width of a chromatid). Therefore, a metaphase chromosome, which consists of two sister chromatids as thick as 1700 nm, is large enough to be clearly seen under a light microscope (Fig. 15). Two specialized structures, centromere and telomere, can be seen with special staining technique (Fig. 16). Chromatin fibers decondense through unpackaging and stretching as the cell returns to the G1 phase.



Fig. 15. Screen snapshot of DNA in sister chromatids of metaphase chromosome



Fig. 16. Screen snapshot of metaphase chromosome with centromere and telomeres

4 Conclusions

This research uses computer technology to enhance life science education. It is part of the interactive program for genetics education [10-14]. The success of this interactive computer program relies heavily on 1) innovative content design that stimulates cognitive feedback through coordinated hands-on interactions at key points of the learning process, and 2) efficient computer programs that are capable of demonstrating complex concepts and processes.

The computer learning modules involve the simultaneous coordination of eyes, mind, and hands for visualization, cognitive feedback, and manipulation, respectively. In this way, complex concepts are scaffolded, reinforcing the learning process. From a pedagogical standpoint, science is essentially taught with the scientific method, with questions, observations, experiments, and analysis. Science learned in this way is more meaningful – and more memorable.

The DNA packaging module, utilizing these methods, is one of a series of modules for learning genetics that can be adopted for a number of biology courses at both college and high school levels.

5 References

[1] Tibell, L. A. E. and Rundgren, C. J. (2010) "Educational challenges of molecular life science: characteristics and implications for education and research" CBE - Life Sci. Educ. 9: 25-33.

[2] Huang, P. C. (2000) "The integrative nature of biochemistry: challenges of biochemical education in the USA" Biol. Educ. 28:14-17.

[3] Bahar, M., A. H. Johnstone, and M. H. Hansell (1999) "Revisiting learning difficulties in biology" J. Biol. Educ. 33: 84-86.

[4] Brig, J. (1996). "Enhancing teaching through constructive alignment" Higher Education, 32:347-364.

[5] Sheley, S. M. and T. R. Mertens (1990) "A Survey of Introductory College Genetics Courses" J. Heredity 81: 153-156

[6] Sved, J. A. (2010) "Genetics Computer Teaching Simulation Programs: Promise and Problems" http://www.genetics.org/cgi/content/full/gentics.110.116640/ DC1

[7] Lowe, R. (2004). "Interrogation of a dynamic visualization during learning" Learning and Instruction 14:257-274.

[8] Tsui, C. Y. and D. F. Treagust (2001) "Teaching and learning reasoning in genetics with multiple external representations" www.aare.edu.au/01pap/tsu01462.htm

[9] Inquiry Based Learning: www.thirteen.org/edonline/concept2class/inquiry/ [10] Yang X., G. Rong, C. Tseng (2011) "Modeling of DNA Replication" The 2011 International Conference on Modeling, Simulation and Visualization Methods, p.146-149, Las Vegas, July 18-21, 2011.

[11] Wu W., X. Yang, C. Tseng (2011) "Effective Alogorithms for Altering Human Chromosome Shapes" The 2011 International Conference on Modeling, Simulation and Visualization Methods, p. 257-261, Las Vegas.

[12] Yang X., R. Ge, Y.Yang, H. Shen, Y. L and C. Tseng (2009) "Interactive Computer Program for Learning the Genetic Principles of Segregation and Independent Assortment through Meiosis" The 31st Annual International Conference of the IEEE Engineering in Medicine and Biology Society (EMBC 2009), p. 5842-5845, Minneapolis.

[13] Wu W., X. Yang, B. Chen, Z. Zhao, J. Lacny and C. Tseng (2009) "Computer Based Simulation of Chromosome Abnormality" The 2009 World Congress in Computer Science Computer Engineering and Applied Computing (WORLDCOMP 2009) p. 359-363, Las Vegas.

[14] Yang, X., D. Wen, Y. Cui, X. Cao, J. Lacny and C. Tseng (2009) "Computer Based Karyotyping" The Third International Conference on Digital Society (ICDS 2009), p. 310-315, Cancun, Mexico..

ANFIS Inverse Kinematics and Precise Trajectory Tracking of a Dual Arm Robot

Arif Ankarali¹

¹ Mechanical Engineering Department, Selcuk University, Konya, TURKEY

Abstract- The dual arm manipulation of a totally six degrees of freedom revolute robot will be considered in this study. For the inverse kinematics purposes of the robot, the Cartesian trajectory assumed for the end effectors of both arms will be transformed to joint space by Adaptive Neuro Fuzzy Inference System (ANFIS) modeling of the system. The Fuzzy control system designed to drive the joint servomotors will utilize this information as a reference input for the joint angular positions and realize the overall assumed Cartesian space trajectory. The whole system will be modeled by using Simulink with some embedded MATLAB functions and the SimMechanics model of the real mechanical construction. The simulations belonging to the assumed trajectories will be realized and the error will be evaluated for the applied Fuzzy control.

Keywords: Dual arm robot, Inverse kinematics, Adaptive neuro-fuzzy inference system (ANFIS), Fuzzy control

1 Introduction

The precise trajectory tracking control of dual arm manipulators has studied by many researchers. This type of manipulators are used on humanoid robots and additionally for some pick up and place operations in dangerous environments, dual arm robots are utilized. Such problems that researchers deal with are generally nonlinear in nature and new techniques like soft computing are developed and discussed to solve and control such systems [1]. Because, controlling real-world non-linear dynamical systems require the use of new techniques to get the desired system performance if the stability and the chaotic motions are considered. The control problem of the cooperative motion of a two-link dual arm robot during handling and transportation of an object was studied by Hacioglu et al. [2]. They applied a fuzzy logic unit improved sliding mode controller to track the desired trajectory with high accuracy. They also studied on a sliding mode controlled dual arm robotic system and found that the SMC made smaller trajectory tracking errors than the PID controller [3]. For heavy objects, it's necessary to use dual arm robots rather than one arm robots. In such applications, passivity requirement is important. Kosuke Kido and et al. had studied on controlling the coefficient of desired velocity field to prevent an accident because of energy flowing from robot to the environment [4]. Hybrid controlled robotic systems were studied by Tinkir et al. in a similar fashion before [7-8]. Duran and Ankarali had studied on a PUMA type robot control [9].

2 Direct and inverse kinematics

The configuration of the dual arm robot together with the coordinate system assignment is given in Fig. 1. Here, the z directions represent the positive angular directions of motions of the joints. The joint angular motions will be driven by joint servomotors with gears.



Fig.1. D-H coordinate system assignment for dual arm robot

From these coordinate systems one can construct the Denavit-Hartenberg (D-H) parameters of the system as given in Table 1. Here, l_1 is the shoulder width and l_2 and l_3 is the back arm and forearm respectively.

Table 1. D-H kinematic parameters.

| Joint | θ_{i} (deg.) | α_i (deg.) | ai | d_i |
|-------|---------------------|-------------------|-------|-------|
| 1R | 90 | 90 | 0 | l_1 |
| 2R | 180 | 0 | l_2 | 0 |
| 3R | 0 | 0 | l_3 | 0 |
| 1L | 90 | 90 | 0 | l_1 |
| 2L | 0 | 0 | l_2 | 0 |
| 3L | 0 | 0 | l_3 | 0 |

The direct and inverse kinematic equations are obtained by D-H approach by utilizing the homogeneous transformation matrices (HTM) for revolute joints. By using the values given in Table 1 HTMs $^{i-1}A_i$ for each neighboring links can be written and then by multiplying those HTM's direct kinematic equations relating the end effecter coordinate system to the first one may be obtained [5]. The direct kinematic equations for the right arm can be calculated as

$$x = \cos(\theta_1) \left[l_3 \cos(\theta_2 + \theta_3) + l_2 \cos(\theta_2) \right]$$
(1)
$$y = \sin(\theta_1) \left[l_2 \cos(\theta_2 + \theta_2) + l_2 \cos(\theta_2) \right]$$
(2)

$$= \sin(\theta_1) \left[l_3 \cos(\theta_2 + \theta_3) + l_2 \cos(\theta_2) \right]$$
(2)

$$z = \iota_1 + \iota_3 \sin(\theta_2 + \theta_3) + \iota_2 \sin(\theta_2)$$
(3)

$$n_{\chi} = \cos(\theta_2 + \theta_3)\cos(\theta_1) \qquad (4)$$

$$n_y = \cos(v_2 + v_3)\sin(v_1)$$
 (3)

$$n_z = \sin(\theta_2 + \theta_3) \tag{6}$$

$$s_x = -\sin(\theta_2 + \theta_3)\cos(\theta_1) \tag{7}$$

$$s_y = -\sin(\theta_2 + \theta_3)\sin(\theta_1) \tag{8}$$

$$s_z = \cos(\theta_2 + \theta_3) \tag{9}$$

$$a_x = -\sin\theta_1 \ a_y = -\cos(\theta_1), \ a_z = 0$$
 (10)

The inverse kinematic equations by equating the suitable corresponding matrix members are

$$\theta_{1} = atan2(a_{x}, -a_{y})$$
(11)

$$s = \frac{z - a_{z}d_{3} - a_{3}n_{z} - d_{1}}{a_{2}},$$

$$c = x - a_{x}d_{3} - a_{3}n_{x} - a_{1}\cos(\theta_{1}) -d_{2}\sin(\theta_{1})/(a_{2}\cos(\theta_{1}))$$

$$\theta_{2} = atan2(s, c) \text{ and } \theta_{3} = atan2(s_{z}, -n_{z})$$
(12)

2.1 ANFIS inverse kinematics

Equations (1-10) are used to calculate the workspace data matching joint space variables to Cartesian space positions and orientations. To calculate the joint variables corresponding to



Fig.2. Workspace of the right arm for $\theta_1 = 90 - 180 \ deg.$, $\theta_2 = 90 - 180 \ deg.$, $\theta_3 = 0 - 90 \ deg.$



Fig. 3. Reference and ANFIS trajectories



Fig.4. ANFIS generated joint trajectories

Cartesian variables, first of all the data for joint variables are constructed by meshgrid command then data arrays for $(x_i, y_i, z_i, \theta_i)$ designed in matrix form. Later, by using anfis command together with the *data_i* the corresponding Cartesian space variables are matched to joint space variables by adaptive neuro-fuzzy structure. This method gives the advantage of safe solution free from singularities, multiple roots of analytical equations and non-linear nature of the equations of the real physical system. The right arm workspace for educating the ANFIS system is given in Fig. 2. A linear reference Cartesian trajectory is compared with the ANFIS calculated trajectory on Fig. 3. The joint space variable changes calculated by hybrid algorithm to follow this trajectory are given in Fig. 4. For these simulations, l_1 , l_2 and l_3 are taken as 0.15 m, 0.30 m and 0.30 m respectively. The masses of the shoulder (m_1) , back arm (m_2) and the forearm (m_3) are assumed as 1 kg, 3 kg and 2 kg respectively. Although the trajectory looks like fine, it could be improved by increasing the number of membership functions and decreasing the step size. Three membership functions are

introduced for this simulation and the number of fuzzy rules is 64. Number of nodes is 158, total number of parameters is 292 and number of training data pairs is 29791.

2.2 Controls

 θ_i

1

2

has 35.

 p_1

0

0

The fuzzy control structure of the whole system is shown in Fig. 6. DC servomotors with gears are used to drive the joints and the block diagram used in this Simulink diagram can be found in [6]. The calculated torques coming from the arm dynamics are introduced as disturbance torques at summing points. These disturbance effects are directly received from the SimMechanics model of the arm. The masses and inertia terms of the links and also the joints are introduced in the related graphical user interfaces. After appropriate selection of servomotors the parameters like torque constants, resistance and inductance values of each motor are supplied to Simulink diagram. Gears are selected for each joint with a reduction ratio of N_i to decrease the angular velocity and increase the output torque to drive the manipulator. ANFIS generated joint trajectories are approximated and fitted by polynomial functions and given in Fig. 5.

$$\theta_i(t) = p_1 t^5 + p_2 t^4 + p_3 t^3 + p_4 t^2 + p_5 t + p_6$$
(13)

Where i=1, 2,...., 6. The coefficients of the polynomial functions are given in Table 2. The first and 4^{th} polynomials are of third order, the 2^{nd} and 5^{th} polynomials are of 5^{th} order, and the 3^{th} and 6^{th} polynomials are of 6^{th} order.

Table 2. Fitted angular motion polynomials coefficients for right arm joints

 p_4

-0.0047

0.06519

 \mathbf{p}_5

-0.091

-0.264

 p_6

2.5

2.662

 p_3

0.0011

-0.0191

 \mathbf{p}_2

0.00445

0

| 3 | -0.002 | 0.01036 | -0.0061 | -0.0876 | 0.1113 | 1.162 |
|-------|------------|-------------|-------------|-------------|------------|---------------------|
| | | | | | | |
| These | e polyn | omials ar | e introd | uced into | an en | nbedded |
| MAT | LAB fur | nction with | a clock t | ool in whi | ch the sir | nulation |
| time | is given a | as 3 s. Cha | racteristic | s of the se | lected DC | ² motors |
| are g | iven in T | Table 3. T | he gears o | combined y | with the f | first and |
| secon | nd motors | s have a re | eduction r | atio of 81 | and the th | hird one |



Fig. 5. The fitted curves from ANFIS generated data

2.2.1 Fuzzy controller design

The basic theory and configuration of fuzzy systems can be found in Wang [10]. The controller is designed to control the angular velocities of the joints. For this reason, derivatives of the joint trajectories calculated before are introduced to the Fuzzy controller as the reference input. Basically, controller consists of two inputs which are the error $E_i(s)$ and the rate of change of the error $dE_i(s)/dt$, and one output $U_i(s)$ as shown in Fig. 7. K_i is the amplifier gain. Input and output membership functions are selected as triangular ones as shown in Fig. 8.



Fig.6. Simulink and SimMechanics structure of dual arm robot control.



Fig.7. General Structure of the fuzzy controller

There are seven membership functions selected for each input and also for output. Forty nine rules are written in the rules table of the fuzzy tool to build the fuzzy structure of the control system.



Fig.8. Fuzzy controller inputs and output.



Fig. 9. Control surface

The rule table is constructed by defining the relations between the inputs and the output. The control surface generated by the fuzzy rules is given in Fig. 9. The horizontal axis of the membership functions are scaled by taking the velocity range of the servomotors into consideration where the vertical one is selected as unity.

2.2.2 Servomotor modeling



Fig. 10. Block diagram of a dc servomotor [11]

Mathematical modeling of a DC servomotor gives the block diagram given in Fig. 10 [11]. The computed torques coming from the arm dynamics are applied to each of the servomotor as a disturbance torque which is calculated by taking the gear reduction into consideration as τ/n . Here τ is the torque coming from arm dynamics and *n* is the gear ratio. In this model, *L* and *R* represent inductance and the resistance of the motor respectively. K_m is motor torque constant and K_b is back emf constant. J_m and B_m are the moment of inertia and viscous friction coefficient of the motor respectively. Input voltage supplied to the electrical part of the servomotor is represented by V(s) and I(s) is the current. All the characteristic values of the selected motors are tabulated in Table 3.

Table 3. Motor's characteristics used in simulations

| Characteristics | Motors 1-2 | Motor 3 |
|---------------------|--|--|
| Terminal resistance | 0.316 Ω | 0.61 Ω |
| Terminal | 82.0 µH | 119.0 μH |
| Torque constant | 30.2 mNm/A | 25.9 mNm/A |
| Speed constant | 317.0 min ⁻¹ / V | 369.0 min ⁻¹ / V |
| Speed / torque | 3.33 min ⁻¹ mNm ⁻¹ | 8.69 min ⁻¹ mNm ⁻¹ |
| Mech. Time const. | 4.67 ms | 3.05 ms |
| Rotor inertia | 139.0 gcm ² | 33.5 gcm ² |



Fig. 11. Reference linear trajectory versus fuzzy controlled response.



a. Before starting its motion

b. In motion

Fig.12. Simulation machine.

3 Conclusions

ANFIS inverse kinematics for a dual arm robot is realized and fuzzy control is applied for precise trajectory tracking. The system is modeled by Simulink and SimMechanics tools of MATLAB and fuzzy control is also included in the model.

From Fig. 3, it's seen that the reference trajectory is precisely followed by the ANFIS produced trajectory. It's observed that the system is working well with the selected motors and the desired trajectory tracking is performed. Reference linear trajectory introduced into the model and the response of the system is shown in Fig.11. The simulation machine snapshots of the designed model are given in Fig. 12.

The systems with higher degrees of freedom may be studied successfully by the same procedure given here in a similar fashion. The SimMechanics model may be obtained by directly transforming computer aided design models of the real design of the robot to get better responses.

4 References

[1] Oscar Castillo and Patricia Melin, "Soft Computing Models for Intelligent Control of Non-linear Dynamical Systems", Studies in Computational Intelligence, Volume 180/2009, 43-70, 2009

Yuksel Hacioglu, Yunus Ziya Arslan, Nurkan Yagiz,
 "MIMO fuzzy sliding mode controlled dual arm robot in load transportation", Journal of the Franklin Institute, 348, 1886–1902, 2011

[3] Nurkan Yagiz, Yuksel Hacioglu and Yunus Ziya Arslan, "Load transportation by dual arm robot using sliding mode control", Journal of mechanical science and technology, Volume 24, Number 5, 1177-1184, 2010

[4] Kosuke Kido, Akinori Nagano, Zhi-Wei Luo, "Passive Control of A Dual-Arm Cooperative Robot", SICE Annual Conference 2010, August 18-21, 2010

[5] K.S. Fu, R.C. Gonzales and C.S.G. Lee, "Robotics", McGraw Hill Book Company, 1987

[6] Katsuhiko Ogata, "Modern Control Engineering", Prentice Hall International, Inc.,1997

[7] Y. Şahin, M.Tınkır, A. Ankaralı, "Trajectory planning and adaptive neural network based interval type-2 fuzzy logic controller design of 3-DOF robot", 3rd International Conference on Computer Research and Development (ICCRD), 2011

[8] Y. Şahin, M.Tınkır, A. Ankaralı, "Neuro-Fuzzy Trajectory Control of a Scara Robot", The 2nd International Conference on Computer and Automation Engineering (ICCAE), 2010

[9] M. Ali Duran, A. Ankaralı, "Trajectory Planning of a PUMA Type Robot", AMSE 2006, International Conference on Modeling and Simulation, Konya, Turkey, 2006

[10] L. X. Wang, "A Course in Fuzzy Systems and Control", Prentice-Hall, Inc, New Jersey, 1997

[11] Mark W. Spong, M. Vidsayagar, "Robot Dynamics and Control", John Wiley and Sons, 1989

Acknowledgment : The author would like to convey thanks to TUBITAK (The scientific and Technological Research Council of Turkey) for providing the financial means.

Microarray Image Processing for Real Time Scanning with Reduced Dimensional Variables

Deok Hee Nam

Engineering and Computing Science, Wilberforce University, Wilberforce, Ohio, USA

Abstract - The image processing of microarrays is like analyzing the arrays of gene fragments from the sampled images. In addition, microarrays provide the information about the detection of differences in genome sequences so that they can be used to detect and classify genetic characteristics very precisely. The experiments of microarray analyses can produce huge amounts of data because thousands of genes require to be processed in a single experiment. In order to obtain significant information from the various experiments of microarray, it is necessary to develop a compact and simpler technique, which can handle the large number of data. This paper provides a technique to perform the real time scanning of microarray images using embedded variables with applying multivariate analysis to reduce the dimension of the original images in vertically as well as horizontally.

Keywords: dimensional reduction, image processing, microarray, multivariate analysis, data mining

1 Introduction

There are many compounds and drug targets with the advent of genomics and advanced combinatorial chemistry. Simultaneously, various methods are introduced to detect and quantify gene expression levels, including northern blots, differential display, and sequencing of cDNA libraries, S1 nuclease protection, and serial analysis of gene expression. Over the last several years there has been a huge expansion of microarray technology [1] in the fields of biosciences including medical sciences, biotechnology, and pharmaceutical industry.

The most important technology which has been a focus on is how to provide a platform for determining in a single experiment, the gene expression profiles of various kinds of genes in tissue, tumors, cell, or biological fluids. The quick and comprehensive implementation about the bio-scientific technology has been evaluated and predicted on the reduction of its complexities with providing large amounts of highly relevant data successfully.

Among the various biomedical scientific technologies, the gene expression about microarrays is the most important technology to improve the blockage for the discoveries of drugs and disease diagnostics. Using the pattern recognition of microarrays [10], the expression of various genes, can be easily identified. Due to the

increments of genomic and advanced biomedical scientific analyses, the technologies for enhanced data extraction and analysis are demanded more and more. Hence, there are now hundreds of papers related to the topics of biomedical scientific or genomic science fields and the application of microarray technology in biological research.

All microarray images can be recorded by microarray scanning technology [8] and represent the visual information [9]. In the fields of engineering and science, image processing is applied to examine properties of objects or processes encoded in images. Hence, after recording the objects or processes for the studies in the images, the image analysis [6] is applied to extract and quantify the characteristics of the objects and processes about the studies by the statistical analyses [4][5][7][12]. For examples, the pattern recognition [2][3] of characters, fingerprints, or various face images performed by computer systems are image analysis tasks. Especially, detecting, outlining, and measuring bio-scientific objects, like boundaries of blood vessels, structural deformations of the heart, degeneration of the retina, in medical applications, and extracting spots and detecting their boundaries, are very important studies to estimate parameters which quantify gene expression levels in microarray applications.

For the microarray image processing, the fundamental goal is to measure the intensities of the arrayed spots, and based on these intensities to quantify the spots expression levels. In a more sophisticated and complete approach, the array image processing will also assess the reliability of the quantified spot data and generate warnings to the possible problems during the array production and/or hybridization phases.

For the instrument point of view, any scientific instrument must perform consistently over time so that results can be compared from day to day and week to week. At the same time, the instrument needs to show whether its performance is in the real time or not. But sometimes the instrument may need more capabilities to perform the required processes due to the limitation of the instrument. For example, scanning the microarray images using a scanner may not be successfully done whenever it needs because of the processing speed of the scanner. To avoid the delayed scanning process, the scanner which can perform the process in the real time as the human eye is scanning the images of the objects may need to be used. In order to achieve the required performance like a human eye, there are some options to approach the real time processing such as using the smaller or simpler microarray images with the less resolution or the faster instrument which can scan the images much faster. However, when the less resolution microarray images are used, the microarray image quality [10] is the key point to perform the image processing not to lose any meaningful information from the sampled microarray images. Hence, if there is a method which can satisfy the condition such as keeping all the information from the sampled microarray images with relatively less dimension of the microarray images, then the results of scanning the extracted, reduced, or transformed microarray images will be same as the results of scanning the original microarray images. To reduce or extract only the required the microarray images from the sampled microarray images with relatively more dimension, the proposed algorithm with multivariate analyses [7] will be used. Scanning the extracted microarray images may give more realistic time for scanning the microarray images without losing the meaningful information from the sampled microarray images and the time delay for the scanning procedure.

In this research, to develop a new algorithm of microarray image processing using the multivariate analyses for the real time processing. The new proposed method can be implanted to the robust of the microarray image scanners whose scanning time is close to the real scanning time like human eyes.

2 Principal Component Analysis (PCA)

Principal component analysis (PCA) [4] is a technique for dimensionality reduction, by diminishing the number of dimension with grouping relatively highly correlated variables among original variables after extracting principal components. Among the newly identified principal components, the first principal component implies the highest variability in the data and each succeeding component accounts for as much of the remaining variability as possible for the original data. In order to identify the extracted components, the eigenvalue or eigenvector approach of the covariance matrices is often used with singular value decomposition (SVD) of the data matrix. [13]

Consider a random vector $X = [X_1, X_2, X_3, ..., X_n, X_P]$ from a multivariate distribution, and let $X_1, X_2, X_3, ..., X_n$ be a random sample of n observations from the distribution and X_P is the corresponding output for X. Assume that it had been centered to zero mean. Since the scales of the different variables will not be proportionate,

it requires the normalization with a common scale by dividing them by their standard deviations.

The mean value of a variable is defined as the sum of the observed values divided by the number of values. It can be shown as following

$$\overline{X} = \frac{\sum_{i=1}^{n} x_i}{n}$$
(1)

where \overline{X} is a mean of the variable and *n* is the number of observations. The variance is determined by dividing the sum of squared deviations from the mean by n - 1. It is represented as

$$v_x^2 = \frac{\sum_{i=1}^n (x_i - \overline{x})^2}{n-1}$$
 (2)

where v_x^2 is the variance of a variable. The standard deviation is found from the variance. The square root of variance is called the standard deviation. Thus

$$v_{x} = \sqrt{\frac{\sum_{i=1}^{n} (x_{i} - \bar{x})^{2}}{n-1}} \qquad (3)$$

where v_x is the standard deviation. When two variables problem is considered, the covariance can be defined as the sum of the products of the deviations of two variables from their respective means divided by the number of the observation in a variable. Assume if there is the same number of observation for each variable. Thus

$$s_{XY} = \frac{\sum (X - \overline{X})(Y - \overline{Y})}{n - 1} \qquad (4)$$

where X and Y are variables, and \overline{X} and \overline{Y} are means of each variables.

The correlation is a measure of the degree of agreement between two sets of scores from the same individuals. The correlation coefficients between two variables can be expressed as the cosine values between any two variables. Therefore, the correlation coefficients measure the tendency of the points to cluster along a line.

Suppose variables can be grouped by their correlations. All variables within a particular group are highly correlated among themselves but have relatively small correlations with variables in a different group. It is conceivable that each group of variables represents a

single underlying construct, or factor, that is responsible for the observed correlations.

Let R be the correlation matrix with q by q, from the data matrix with n by p, A, can be expressed as following

$$R = \begin{bmatrix} r_{11} & r_{12} & \dots & r_{1q-1} & r_{1q} \\ r_{21} & r_{22} & \dots & r_{2q-1} & r_{2q} \\ \cdot & \cdot & \dots & \cdot & \cdot \\ \cdot & \cdot & \dots & \cdot & \cdot \\ r_{q1} & r_{q2} & \dots & r_{qq-1} & r_{qq} \end{bmatrix}$$
(5)

where

$$q = max \{ n, p-1 \}$$

and

$$r_{ij} = \frac{s_{ij}}{s_i s_j} = \frac{\sum d_{ik} d_{jk}}{\sqrt{\sum d_{ik}^2 \sum d_{jk}^2}}$$
(6)

with unit diagonals. Note that s_{ij} is the covariance of the matrix *A* for any two variables *i* and *j*.

After calculating the correlation, then estimate the eigenvectors with eigenvalues. From the matrix of the eigenvectors, rearrange the eigenvectors for each eigenvalue represented by descending order of the eigenvalues (variances) such as $\lambda_1 \geq \lambda_2 \geq \lambda_3 \geq \dots \geq \lambda_k$ solving the system equation $|R - \lambda I| A = 0$ with $A^T A = I$, where A is the eigenvector matrix. After that, extract the vector which covers the majority part of the matrix R with eliminating the closely related variables from the original data set. (Normally the cumulative value of the eigenvalues, or variances, are greater or equal to 0.9. [14]) Therefore, the principal components of each observation will be scored by

$$Y_1 = A_1^T X, Y_2 = A_2^T X, \dots, Y_p = A_p^T X$$
 (7)

3 Fuzzy C-Means (FCM) Clustering Analysis [2]

Fuzzy C-Means Clustering (FCM) algorithm attempts to partition a finite collection of n elements of the given data into a collection of c fuzzy clusters with respect to some given criterion. Given a finite set of data, for $2 \le c \le$ n, consider a set of n vectors $X = \{x_1, x_2,, x_n\}$ to be clustered into c groups of data. Each of the groups, $x_i \in \mathbb{R}^S$, is a feature vector consisting of s real-valued measurements describing the features of the object represented by x_i . The features could be length, width, color etc. Fuzzy clustering of the objects can be represented by a fuzzy membership matrix called as a fuzzy partition.

The set of all $c \times n$ non-degenerate constrained fuzzy partition matrices is denoted by M_{fcn} and is defined as

$$U = [u_{ij}]_{i=1,2,...,c,\,j=1,2,...,n} = M_{fcn}$$

where u_{ij} expresses the degree to which the element x_j belongs to the *i*th cluster and is a numerical value in [0,1] such that the constraints in

and

$$0 < \sum_{j=1}^{n} u_{ij} < n$$
 for all $i = 1, 2, \dots, c$

 $\sum_{i=1}^{c} u_{ij} = 1 \text{ for all } j = 1, 2, \dots, n$

For the Fuzzy C-Means algorithms, the objective is to find $U = [u_{ik}] \in M_{fcn}$ as the fuzzy c-partition matrix and $V = (v_1, ..., v_c)$ with $v_i \in R^s$ as the cluster center such that

$$J_m(U,V) = \sum_{k=1}^n \sum_{i=1}^c (u_{ik})^m \|x_k - v_i\|^2$$
(8)

is minimized, where u_{ik} is the value of the i^{th} membership function on the k^{th} data point x_k , n is the number of samples in X, and $m \in (1, \infty)$ is a weighting constant. The distances, $d_{ik} = ||x_k - v_i||^2$, are weighted with the membership values u_{ik}^m , where $||x_k - v_i||^2$ is any inner product induced norm on R^s which is called the square of Euclidean distance for i^{th} cluster center and j^{th} data point. The Euclidean distance formula will be

$$d_{ik} = ||x_k - v_i|| = (\{x_k - v_i\}^2)^{1/2}$$
(9)

where d_{ik} will be the distance between i^{th} cluster center and k^{th} data point x_k .

The assumption is that the distance between their corresponding data vectors measures the similarity between objects. Then the necessary conditions to minimize the objective function, $J_m(U,V)$, can be

$$\boldsymbol{\mu}_{ik} = \frac{1}{\sum_{j=1}^{c} \left(\frac{\boldsymbol{d}_{ik}}{\boldsymbol{d}_{jk}}\right)^{\frac{2}{m-1}}}$$
(10)

and

$$v_{i} = \frac{\sum_{k=1}^{n} (u_{ik})^{m} x_{k}}{\sum_{k=1}^{n} (u_{ik})^{m}}$$
(11)

where l < i < c, and l < k < n.

Therefore, the Fuzzy C-Means algorithm is an iterated procedure through those two necessary conditions to minimize the objective function, $J_m(U,V)$.

The following steps summarize the Fuzzy C-Means Clustering algorithm.

- Step 1. Initialize the partition matrix, or membership matrix such that $U^{(0)} \in M_{fcn}$ randomly or prior knowledge.
- Step 2. Calculate the cluster centers, v_i , using the equation (11).
- Step 3. Compute the distance, d_{ik} .
- Step 4. Update the partition matrix $U^{(new)}$ using the equation (10) for u_{ik} .
- Step 5. Until $|| U^{(new)} U^{(old)} || < \varepsilon$ where ε is the termination tolerance $\varepsilon > 0$. If this condition is not satisfied, then go back to Step 2.

4 Proposed Multivariate Analysis

There are various techniques to reduce or diminish the huge data set into an appropriate size of the data without losing any significant meaning from the original huge data. Among the frequently used techniques, multivariate analysis such as Principal Component Analysis or Factor Analysis is refocused by the scientists in these days. To develop the suitable technique for extracting the required microarray image, Fuzzy C-Means Clustering Analysis is used as a postprocessing procedure after the images are processed by Principal Component Analysis. The following algorithm summarizes the proposed algorithm using the Principal Component Analysis followed by Fuzzy C-Means Clustering Analysis.

- Step 1. Read the original data set as a matrix format.
- Step 2. Normalize the original data from Step 1.
- Step 3. Find the correlation matrix of the normalized data from Step 2.
- Step 4. Find eigenvalues and eigenvectors of the correlation matrix from Step 3 using characteristic equation.
- Step 5. Define a matrix that is the eigenvectors from Step4 as the coefficients of principal components using the criterion for extracting components.
- Step 6. Multiply the standardized matrix from Step 2 and the coefficients of principal components from Step 5.

- Step 7. Using the result from Step 6, find the centers of clusters for each clustering technique.
- Step 8. To clustering the components from Step 7, initialize the partition matrix, or membership matrix randomly such that $U^{(0)} \in M_{fcn}$.

Step 9. Calculate the cluster centers, v_i , using the equation

(11).

Step 10. Compute the distance, d_{ik} . Step 11. Update the partition matrix $U^{(new)}$ using the

equation (10) for u_{ik} . If $d_{ik} > 0$, for $1 \le i \le c$, and $1 \le k \le n$, then get the new u_{ik} .

Otherwise if $d_{ik} > 0$, and $u_{ik} = [0, 1]$ with

$$\sum_{i=1}^{c} u_{ik}^{(new)} = 1$$
, then $u_{ik}^{(new)} = 0$

Step 12. Until $|| U^{(new)} - U^{(old)} || < \varepsilon$ where ε is the termination tolerance $\varepsilon > 0$. If this condition is not satisfied, then go back to Step 9.

5 Analysis and Results

To develop the real time scanning for the microarray image, the reduced-dimensional images from the original image with applying the proposed algorithm are analyzed. To reduce the dimension into the less number of pixels, the appropriate number of new components needs to be determined without losing any significant meaning from the original image. In Fig. 1, the relation between the evaluated eigenvalues and the number of newly extracted components are plotted to determine the appropriate number of newly extracted components from the original image.



Fig. 1 The number of extracted components vs. the corresponding eigenvalues

In Fig.2, the images are shown how the original images can be reduced by applying selected proposed algorithms. To analyze the reduced dimensional images comparing with the original image to present a method of simplifying the analysis of large amounts of microarray image data by reducing information without losing validity, the consecutive sampled rows of pixels from each extracted image are compared by evaluating its correlation coefficients [11]. As shown above, when the image (b) is reduced by 11 by 9 using Principal Component Analysis, the correlations between the original image and the extracted image are -0.704, -0.594, -0.914, and 0.036, in TABLE 1, respectively. Since the images are shown that the lighter parts extracted as the darker part and the darker parts are extracted by the lighter parts after applying the procedure. There appears to be some loss of intensity information, as represented by the diagrams shown in the analysis.

From Fig. 2, the extracted image with reducing the dimension of the height and width using the factor analysis with principal components and Fuzzy C-Means Clustering Analysis is shown that the microarray spot can be recognizable even though the image of the bottom-right portion is not clearly recognized as comparing to the original image. From (c) in Fig. 2, the extracted image loses the characteristics of the right-middle portion for the original image with reducing the dimensions in the vertical and horizontal directions. From (d) and (e) in Fig. 2, the images loses the characteristics of the original image as comparing with the image (a). In (f) in Fig. 2, even though the image is shown with a less intensity image as comparing with the image (a), but the characteristic of the original image as comparing with the image (a).

| | AP1 | AP2 | AP3 | AP4 |
|-----|--------|--------|--------|--------|
| bP1 | -0.704 | -0.623 | -0.84 | -0.104 |
| bP2 | -0.853 | -0.594 | -0.925 | -0.035 |
| bP3 | -0.76 | -0.64 | -0.914 | 0.007 |
| bP4 | -0.713 | -0.603 | -0.905 | 0.036 |
| cP1 | -0.662 | -0.707 | -0.886 | -0.064 |
| cP2 | -0.853 | -0.594 | -0.925 | -0.035 |
| cP3 | -0.713 | -0.603 | -0.905 | 0.036 |
| cP4 | -0.704 | -0.623 | -0.84 | -0.104 |
| dP1 | -0.715 | -0.605 | -0.906 | 0.034 |
| dP2 | -0.832 | -0.468 | -0.774 | -0.06 |
| dP3 | -0.877 | -0.39 | -0.713 | -0.079 |
| eP1 | -0.726 | -0.607 | -0.909 | 0.028 |
| eP2 | -0.884 | -0.521 | -0.89 | 0.042 |
| eP3 | -0.721 | -0.673 | -0.905 | -0.029 |
| fP1 | -0.644 | 0.166 | -0.321 | -0.129 |
| fP2 | 0.121 | 0.266 | 0.251 | 0.298 |
| fP3 | 0.026 | 0.607 | 0.503 | 0.275 |
| fP4 | 0.278 | 0.208 | 0.393 | 0.06 |

| m / DI D / | G 1 1 | • | | |
|------------|--------------|---------|---------|--------|
| TABLET | Correlations | between | sampled | nixels |

(a) aP1, aP2, aP3, and aP4 are sampled pixels from the original images.
(b) bP1, bP2, bP3, and bP4 are sampled pixels from the 11 by 9 extracted image extracted by Principal Component Analysis. (c) cP1, cP2, cP3, and cP4 are sampled pixels from the 10 by 9 extracted images processed by Principal Component Analysis followed by Fuzzy C-means Clustering Analysis. (d) dP1, dP2, and dP3 are sampled pixels from the 8 by 9 extracted images processed by Principal Component Analysis followed by Fuzzy C-means Clustering Analysis. (e) eP1, eP2, and eP3 are sampled pixels from the original 6 by 9 extracted images processed

by Principal Component Analysis followed by Fuzzy C-means Clustering Analysis. (f) fP1, fP2, fP3, and fP4 are sampled pixels from the 11 by 9 extracted image extracted by Factor Analysis.

6 Conclusion

The microarray analysis is an important research topic in the biomedical image processing. Moreover, the amount of the processing data is a very important issue with the processing time. To improve the data processing and its execution, the proposed algorithm is developed to help the real time scanning and reduce the amount of the time to process the analysis of the microarray images. Through the various extracted images with reducing the dimensions in vertical and horizontal directions, the characteristics of the original image can be recognized by the dimension-reduced images without losing its validities if the dimension of the image is kept in the reasonable intensity.



(a) Original image of a single microarray spot

(b) 11 × 9 extracted image of a single microarray spot with Principal Component Analysis

(c) 10 × 9 extracted image of a single microarray spot with Principal Component Analysis followed by Fuzzy C-Means Clustering Analysis
(d) 8 × 9 extracted image of a single microarray spot with Principal Component Analysis followed by Fuzzy C-Means Clustering Analysis
(e) 6 × 9 extracted image of a single microarray spot with Principal Component Analysis followed by Fuzzy C-Means Clustering Analysis
(f) 11 by 9 extracted image of a single microarray spot with only Factor Analysis.

Acknowledgement

This project has been supported by United Negro College Fund (UNCF) through Henry C. McBay Fellowship.

7 References

- G. Kamberova and S. Shah, "Microarrays and image analysis: Introduction", DNA array image analysis nuts and bolts, DNA press LLC, Salem, MA, pp. 8 – 15, 2002.
- [2] R. O. Duda, P. E. Hart, and D. G. Stork, Pattern classification, 2nd Ed., John Wiley & Sons, Inc., Danvours, MA, 2001.
- [3] James C. Bezdek, "A Review of Probabilistic, Fuzzy, and Neural Models for Pattern Recognition," Journal

of Intelligent and Fuzzy Systems, Vol. 1, No. 1, pp. 1-25, 1993.

- [4] I.T. Jolliffe, Principal Component Analysis, New York, NY, Springer- Verlag New York Inc., 1986.
- [5] Maurice Kendall, Multivariate Analysis, New York, MacMillan Publishing Co. INC., 1980.
- [6] G. Kamberova, "Introduction to image analysis", DNA array image analysis nuts and bolts, DNA press LLC, Salem, MA, pp. 17 – 50, 2002.
- [7] D. Nam, "Realization of Crop Viruses by Neurofuzzy Systems using Hybrid Mining Algorithm", Proceedings of the 4th International Conference on Cybernetics and Information Technologies, Systems and Applications, Orlando, Florida, pp. 98 – 103, July 12–15, 2007.
- [8] D. Verdnik, S. Handran, and S. Pickett, "Key considerations for accurate microarray scanning and image analysis", DNA array image analysis nuts and bolts, DNA press LLC, Salem, MA, pp. 82 – 98, 2002.
- [9] D. Moody, B. Fadlia, A. Singh, S. Shah, and L. McIntyre, "Quantitative comparison ofimage analysis software," DNA press LLC, Salem, MA, pp. 155 – 166, 2002.

- [10] A. Petrov, S. Shah, S. Draghici, and S. Shams, "Microarray image processing and quality control", DNA array image analysis nuts and bolts, DNA press, LLC, Salem, MA, pp. 99 – 130, 2002.
- [11] Allen L. Edward, An Introduction to Linear Regression and Correlation, San Fransisco: W. H. Freeman and Company, 1976.
- [12] Andrew L. Comrey and Howard B. Lee, A first course in factor analysis, Hilldale, NJ: Lawrence Erlbaum Associates Inc., 1992.
- [13] Hairong Qi, Tsei-Wei Wang, and J. Douglas Birdwell, "Global Principal Component Analysis for Dimensionality Reduction in Distributed Data Mining," University of Tennessee, Knoxville, TN 37996, USA
- [14] N. Cliff, "The Eigenvalues-Greater-Than-One Rule and the Reliability of Components," Psychological Bulletin, Vol. 103, No. 2, pp. 276–279, 1988.

Mitigation of Low Frequency Power Oscillations Generated by a Hydroelectric Generation Station of CFE-México

G. Villa-Carapia¹, O. Mora-Hoppe¹, G. Enriquez Harper¹, F. Sánchez-Tello¹, G. Carreón-Navarro¹, and E. Espinosa-Juárez²

¹Comisión Federal de Electricidad, México,D.F., México ²Electrical Engineering Faculty, Universidad Michoacana de San Nicolás de Hidalgo Morelia, Michoacán, México

Abstract - This paper presents a study to solve the problem of undamped low frequency power oscillations frequency on July 31, 2008, in which the 600 MW hydroelectric generation station "El Caracol" oscillates in respect to the national network. The event was recorded by phasor measurement units making possible the identification of an unstable electromechanical mode and the dominant modes of the oscillation.

Keywords: power oscillation of low frequency, small signal stability, power system stabilizers

1 Introduction

The dynamic stability of the Electric Power Systems (EPS) is determined by the positive contribution of two concurrent forces, a synchronization and a damping. The synchronization torque ensures synchronous operation of all units of the network, while the damping torque can quickly reach a new operating point after a disturbance [1].

The stability studies, an essential tool for planning and operating the EPS, have traditionally focused on the transient behavior of generators in the system. Generally, the damping of the system is considered sufficient to different operating conditions; however, this damping may be adversely affected by changes in the pattern of power flows due to the entry of new energy producers in the electricity industry. Of the upmost importance in Mexico is the effect of the location of each new licensee and the new technology used in type combined cycle power plants and their associated controls.

When an EPS has problems of poor damping, this is manifested through sustained power oscillations, which can grow and cause electrical service interruptions to make protection devices operate to disconnect generating units or loads for high or low frequencies.

Damping problems in power systems are generally associated with the interaction of the control system of generating units, the power demand condition and the electrical network topology. In modern interconnected power systems, the Power System Stabilizer (PSS) is widely utilized to damp low frequency power oscillations [1][2].

The most efficient cost-benefit solution to the problem of oscillations is the implementation and tuning of PSS in the generating units. Some publications also note that the tuning of the Automatic Speed Regulators (ASR) has significantly improved the damping of unstable oscillations of very low frequency [1].

The events reproduction and the solution of the oscillations problems in power systems is difficult because of the high level of modeling required for such studies and the detail needed to represent each generator and its control systems (automatic voltage regulators (AVR), PSS and ASR), which require considerable work in identification, mathematical modeling, validation and testing of control systems of the generating units [3][4].

The particular interest in analyzing the behavior of the "El Caracol" generating station has contributed to the greater participation in the undamped low frequency power oscillations, recorded on July 31, 2008, in the Mexican National Electric Power System. The focus for the solution of the oscillations problem was directed to tuning the PSS of these generating units.

In this paper the methodology for mitigation of power oscillations of low frequency generated by the Hydroelectric Generating Units (HGU) "El Caracol" of the Federal Electricity Commission (CFE) is presented. This is a 600 MW generating station which is part of the Mexican National Electric Power System.

2 Undamped power oscillations which occurred on July 31, 2008

The network associated with the HGU "El Caracol" prevailing on July 31, 2008 is shown in Fig. 1.



Fig. 1 Network topology associated to the HC "El Caracol"

At 14:44 hours, unit 1 increases the generation of 158 MW to 200 MW, summing up a total plant generation of about 600 MW. The change in generation dispatch of hydroelectric plant motivates the presence of an undamped power oscillation of low frequency. The records showed an oscillation frequency of about 1 cycle per second and negative damping values; they displayed an angular difference of 180 degrees to the other generators of the National Interconnected System (NIS). The oscillations event lasted approximately 1.5 minutes and was eliminated naturally by being insolated generation of "El Caracol" by firing the 230 kV transmission line "El Caracol-Mezcala" to operate the directional overcurrent scheme causing a NIS lowering frequency of 59.74 Hz. The behavior described characterizes a typical problem of small signal stability.

The approach to the solution of the problem is directed toward tuning the PSS of the HGU "El Caracol", as these generators are the most active participants in the oscillations presented.

3 Methodology to solve the problem of undamped oscillations

Under the coordination of the Operation Management of the National Energy Control Center (CENACE) of CFE-México, a group of engineers was formed to analyze and to solve the undamped oscillations of power recorded on July 31, 2008. The analysis group included the CFE specialists of the sub-management of CENACE, the Specialized Engineering Unit, the sub-management of Programming, Generation and Transmission, and the Equipment and Materials Test Laboratory. The procedure for solving the oscillations problem is as follows:

- 1) To develop and to validate the mathematical models of AVR, PSS, ASR [3].
- 2) To reproduce the event in the time domain.
- 3) To define PSS operation.
- To obtain the oscillation frequencies and damping of the event.
- 5) To define the new settings AVR and PSS
- 6) To simulate in the time domain and frequency domain for different operating conditions of the network.
- 7) To test in field the performance of PSS tuning.
- 8) To evaluate in the time domain and frequency domain behavior of the network with tuned PSS.

Playback of the oscillation mode of interest represents from the standpoint of analysis, the primary challenge in the study of small signal stability. The difficulty lies in the high level of modeling required for these studies because the detail necessary to represent each generator and control systems, which requires considerable work in identifying mathematical modeling, validating and testing of systems control of generating units. In most cases field tests are necessary, which implies the need for human capital, test equipment and the availability of generating units.

The initial work group's analysis focused on the following: testing the behavior of AVR and PSS of the HGU "El Caracol"; the field review and confirmation of the manufacturer's information on these control systems; development and validation of their mathematical models; and finally the formation of a database that represents properly the dynamic behavior of the NIS [3].

4 Phasor measurements

The event of July 31, 2008 could be recorded by phasor measurement units (PMU) that are installed at various substations of the NIS. The information provided by these measurements is very important because it represents the real dynamic behavior of the NIS; from this you can adjust the mathematical models that represent the dynamic behavior of the system to digitally reproduce the event as closely as possible.

The modal identification analysis of these measurements can determine the degree of geographical spread of the different modes of dominant low frequency oscillation.

Figure 2 shows measurements of the SCADA sequence of events that occurred on July 31, 2008, and Fig. 3 shows the PMU records at the substation Texcoco for the same day.

The records of PMU's installed in different parts of NIS show evidence of oscillations over the NIS on July 31. Figure 4 shows the PMU measurements of plants "Carbón Dos (CBD)" and "Rio Escondido (REC)". Figures 2-4 show that poorly damped oscillations are manifested throughout the NIS, and these were stimulated by the small perturbation of 40 MW increasing the generation of the HGU "El Caracol".



Fig. 2. Active power measurements of the July 31 event at 14:44 hours.



REAL POWER

Ň

SAMPLE NUMBER — REC-U3 — REC-U4 — CBD-U1 — CBD-U2

Fig. 4. Phasor measurements of real power of "Carbón Dos" and "Rio Escondido"

5 Characterization of oscillation modes

The problems of poorly damped power oscillations can cause service interruptions of electrical power, prompting protection devices to disconnect generating units or loads or low frequencies in the system. According to the type of oscillation, the dynamic patterns of behavior can be classified as oscillation modes of inter-area, between plants and between units, each with an oscillation frequency within a typical range depending on the number of generating units involved.

A power swing is characterized by a frequency value and a damping level. For each oscillation frequency, the generating units reveal a specific and particular activity, manifested by the magnitude of deviation of the oscillation's speed and phase angle relative to the rest of the units. The oscillation frequency will be determined by the capacity of each generating unit, its location on the grid and operating condition.

Based on the above indentifying dynamic patterns of behavior of all generating units of EPS is possible in every operating condition. The damping level of each oscillation frequency is a complex function of the interaction of controls voltage and speed for each generating unit, also controls in the grid.

5.1 **Reproduction and event analysis**

With the development and validation of mathematical models of AVR and PSS-generating units of the HGU "El Caracol" and the rest of the dynamic database of the NIS, it was possible to reproduce with certain accuracy electromechanical oscillations recorded. Figure 5 shows PMU records of active power flow in transmission line TEX-LAP and digital simulation results.

Figure 6 shows the digital simulation of the effect to increase the active power of the unit-1 of HGU "El Caracol"(CRL-U1) from 158 to 200 MW.

From Fig. 5 and Fig. 6 it is noted that the event has an oscillation frequency of about 1 Hz and a negative damping ratio of -0.2%. The mode shape for the event on July 31, 2008 features a local oscillation mode, where the units of HGU "El Caracol" oscillate with respect to the rest of NIS.

The mode shape and participation factors that allow the units to find the greatest impact on the recorded oscillation are obtained by the analysis module in the frequency domain DSATools software [5].

Table 1 shows the oscillation modes related to CRL-U1 for the event of increased generation in this unit, and Table 2 shows the participation factors for the mode of 0.947 Hz.



Fig. 5. Comparison of PMU and digital simulation





Fig. 6. Digital simulation of active power increase in CRL-U1

| CRL-U1 [MW] | Frequency [Hz] | Damping ratio [%] |
|------------------|---------------------|----------------------|
| 158 | 0.947 | -0.20 |
| 168 | 0.946 | -0.45 |
| 178 | 0.945 | -0.70 |
| 188 | 0.943 | -0.96 |
| 198 | 0.942 | -1.23 |

Table 1. Oscillation modes related to CRL-U1.

Table 2. Participation factors for the mode of oscillation of 0.9474 Hz.

| Participation Factor Mode 0.947 Hz | | |
|---------------------------------------|--------|--|
| Value | Unit | |
| 1.00 | CRL-U1 | |
| 1.00 | CRL-U2 | |
| 0.99 | CRL-U3 | |

From the obtained results of the event reproduction a negative value can be observed for the damping coefficient, which means that there is a presence of oscillations growing in time. As generation increases the CRL-U1 damping factor seems to be more negative. To correct this situation and improve the damping level of NIS, it was decided to make an adjustment to PSS parameters of the HGU "El Caracol". The particular aim is to dampen the positive oscillation mode of the plant against the rest of the interconnected system.

6 Tuning of PSS

From the point of view of the analysis, the relevant parts of PSS to consider are the phase characteristic and the gain. The phase characteristic is determined by the dynamic equivalent network seen from the AVR control loop generating unit, evaluated in the frequency range of interest, which corresponds to the phase of the function GEP(s) [6].

On the other hand, for the gain of the PSS there are several criteria to determine its final adjustment. One of them establishes a maximum value at the point where the damping mode interest begins to subside, without appreciably affecting the other oscillation modes. Another criterion sets its value directly on the field, where the maximum value is set to one third of the gain value which destabilizes the control of the PSS. Other criteria establish the value obeying minimum damping criteria in a selected set of modes. In this paper the first criterion was preliminarily considered, and the final value is established making adjustments in the field.

The adjustment of phase characteristic to compensate by the PSS is determined by the operating condition of the network; therefore, the variation of this feature will establish a range of settings, which can include inaccuracies in modeling the network and controls, ensuring an increase in damping in the frequency range of interest, and also ensures that the final adjustment follows as well as possible the variations of the characteristic GEP(s) [6].

After calculating the band of variation of the characteristic GEP(s), the best fitting parameters to the phase characteristic of the PSS in this band of variation are found.

In Figure 7 for CRL-U1 the band of variation of GEP(s) for different operating conditions with the complete network and for several contingencies, as well as different adjustments of the phase characteristic of the PSS in the frequency range 0.1-2.0 Hz are shown.

Figure 8 shows for CRL-U1 PMU records of the test of positive reactive step with the PSS connected and digital simulations with different settings for phase and magnitude of the PSS. Graphics in Figure 8 show that a good fit of the PSS at the GEP(s) phase for a frequency range of around 1 Hz is obtained when the phase shift of the PSS stabilization signal is zero. In addition the gain of this signal

is 1.2 pu, but the overall design of these include the setting of common elements in any PSS, compensating block and torsional filter. The focus adjustment of the PSS is centered on a frequency of 0.947 Hz.

In Table 3 the impact of the PSS settings for the CRL-U1 in the interest oscillation mode is presented. It is observed that the proposed adjustments to the PSS provided acceptable positive damping (greater than 5%). The generation increase in CRL-U1 represents the July 31, 2008 condition, where there is an acceptable positive damping factor.



Fig. 7. Impact of the operating conditions in the GEP(s) characteristic



Fig. 8. Variation impact of the PSS gain, considering a phase of 0.5 pu. Test and simulation results.

Table 3. Oscillation modes related to CRL-U1...

| CRL-U1 [MW] | Frequency [Hz] | Damping ratio [%] |
|------------------|---------------------|----------------------|
| 158 | 0.6870 | 11.92 |
| 168 | 0.6857 | 11.65 |
| 178 | 0.6844 | 11.38 |
| 188 | 0.6831 | 11.10 |
| 198 | 0.6817 | 10.82 |

The final results of the PSSs tuning is shown in Fig. 9-Fig. 11, where PMU records of tests at the plant with the proposed adjustments of the PSS and the computer simulations are shown. Figure 9 shows the active power behavior of CRL-U2 before the test of positive reactive step with the PSS disabled.

Figure 10 shows the active power behavior for the CRL-U2 when the positive reactive step test with the PSS connected is applied.



Fig. 9. Positive reactive step with PSS disabled. Test and simulation results

TEST OF POSITIVE REACTIVE STEP WITH PSS $\Delta V = +0.0381$ pu



Fig. 10. Positive reactive step, with PSS disabled. Test and simulation results

Figure 11 shows the behavior of field voltage for the CRL-U2 unit when the test of the reactive step with PSS is applied.

Figure 12 illustrates the dynamic behavior of the HGU "El Caracol" with the PSS proposed adjustments to increase the generation of CRL-U1, which represents the operation condition of the event on July 31, 2008. Figure 12 shows the positive damping provided by PSS tuning of the HGU "El Caracol".



Fig. 11. Field voltage of the CRL-U2 unit for the test of the positive reactive step with PSS. Test and simulation results





Fig. 12. Event simulation with PSS tuning

7 Conclusions

The reproduction of event of July 31, 2008 was successfully performed. The methodological sequence identification, mathematical modeling and the testing of AVR and PSS behavior are the basis of such studies. This highlights the need for a validated dynamic database.

Highly qualified human capital is crucial to develop the identification, mathematical modeling, performance tests and validation of mathematical models to ensure success in solving problems of low frequency oscillations.

In the studied event was identified a mode of oscillation with a frequency of 0.947 Hz and a negative damping ratio of 0.2%, where the HGU "El Caracol " oscillate against the units of the NIS.

The characterization of the oscillation mode shows that units with greater participation in this oscillation mode are those of the HGU "El Caracol".

The results show a positive impact on the damping of the oscillation mode of interest, which is because of appropriate mathematical modeling of controls on generating units, and additionally to the PSS tuning of the HGU "El Caracol ".

8 References

- P. Kundur, "Power system stability and control", EPRI Power System Engineering Series, McGraw Hill, 1994.
- [2] H. G. Far, H. Banakar, P. Li, C. Luo, B. T. Ooi, "Damping interarea oscillations by multiple modal selectivity method", IEEE Trans. Power Delivery, vol. 24, no. 2, pp. 766-775, May 2009.
- [3] G. Villa-Carapia, O. Mora-Hoppe, F. Sánchez-Tello, G. Carreón-Navarro, R. García-Kasusky, A. Guzmán-Terrones, E. Espinosa-Juárez, "Mathematical modeling of automatic voltage regulators and power system stabilizers for a hydroelectric generating unit of CFE-México", MSV'11, Modeling, Simulation & Visualization Methods, WordComp 2011, Las Vegas, Nevada, July 18-21, 2011.
- [4] A. Dysko, W. E. Leithead, J. O'Reilly, "Enhanced power system stability by coordinated PSS design", IEEE Trans. Power Systems, vol. 25, no. 1, pp. 413-422, Feb. 2010.
- [5] Power Systems Technologies, Powertech, "Dynamic security assessment software", DSAToolsTM", 2010.
- [6] P. Kundur, M. Klein, G.J. Rogers, M.S. Zywno, "Application of power system stabilizers for enhancement of overall system stability", IEEE Transactions on Power Systems, vol. 4, no. 2, May 1989.
- [7] R. T. Byerly, E. W. Kimbark, "Stability of large electric power systems", IEEE Press, 1974.

TEST OF POSITIVE REACTIVE STEP WITH PSS $\Delta V = +0.0381$ pu

Zonal Statistics to Identify Hot-regions of Traffic Accidents

Ömer M. Soysal^{1,2,3}, Helmut Schneider^{1,2}, Asim Shrestha^{1,2}, Christy D. Guempel², Pei Li², Harisha Donepudi^{2,3}, Naveen K. Kondoju^{2,3}, Kazim Sekeroglu²

{¹Department of Information Systems and Decision Sciences, ²Highway Safety Research Group, ³Department of Computer Science}, Louisiana State University, Baton Rouge, LA, USA

Abstract - This paper presents how to utilize ArcGIS zonal statistic tool for identification of hot-regions. For this purpose, we used two descriptive statistical indexes Sum and Max over normalized raster data. Our effort was to get a better understanding of the results generated by the zonal statistics tool to identify the traffic accidents hot-spot zones. Our approach makes it possible for comparison of statistical parameters across different zone areas within the same scale. This helps decision makers to pinpoint locations with higher crash index and take preventive measures.

Keywords: Hot-spot, Zonal statistics, traffic accident, raster.

1 INTRODUCTION

Spatial data like weather data, demographics, land management data, and crash data need to be analyzed at different levels for analysis of trends and patterns. In the analysis of crash data, it is crucial to identify hot-spot regions of traffic accidents. For such identification, a comparison of different zones of interest is required based upon statistical parameters.

A difficulty in this type of analysis is that aggregate spatial data (like sum, maximum) exists at different levels of scale on the map, for example, sum or maximum of crime incidents across the nation cannot be directly compared to that of a county or a city. For such type of comparative analysis, we propose to normalize the results obtained from the zonal analysis in ArcGIS. We use these normalized values to identify the hot-spots of traffic accidents over different levels of administrative zones. Furthermore, the use of two different statistics parameters ensures that different hot-spot regions can be identified according to the nature of analysis.

Hot-spot Analysis is a key area for research not just for traffic accidents but a wide range of other topics such as environmental studies, spread of diseases, biological studies, migration studies, etc. Different methodologies have been proposed in the literature to study the patterns of occurrence of spatial data. The methodologies can be cross-domain; meaning that methods applied to find out hot-spots in one area of research can also be applied to another. [1] has used Kernel Density Estimation and K-means clustering to study the spatial patterns of injury related road accidents and to create a classification of road accident hot-spots. There is no universal definition of accidental hot-spots and researchers have used different sophisticated methods to quantify hot-spots. [2] has provided a review of the following three different hot-spot detection techniques:

- Kernel density estimation
- Network analysis
- Census Output Area estimation

An extension to the Kernel density estimation in which the network space is represented with basic linear units of equal network length, known as lixe (linear pixel) and related network topology is proposed by [3]. They argue that the use of lixe facilitates the systematic selection of a set of regularly spaced locations along a network for density estimation and makes the practical application of the network Kernel density estimation feasible by significantly improving the computation efficiency.

[4] has used a process called kernel smoothing for hot-spot analysis. This process creates local estimates of the measure of the spatial intensity using the count of frequency of points within a given distance of each point, relative to symmetric distribution. The author has used a Gaussian kernel of 0.5 kilometer bandwidth for the hot-spot analysis.

Zonal statistics which we employed in this paper follows a raster-based method. Raster-based methods are widely used in environmental and geophysical studies [5], geographic analysis [6], surface modeling [7], planning and design [8], biological studies [9].

In this paper, we have used the sum values of different zones (of the same scale) to identify hot-spots with a greater number of events (crashes) per unit area. The Max values, on the other, hand are used to identify hot-spot zones that have high number of events at sub-zonal units. We have used two statistical parameters instead of just one so that different zones can be identified as hot-spots according to the objective of the analysis. The usage of zonal statistics and the visualization of normalized results obtained from it will help to draw meaningful analysis and interpretations. It serves as a visual aid in knowing where hot-spots occur and helps in finding out the reason behind high number of un-expected crashes in specific locations.

2 METHOD

We first used the zonal statistics tool in ArcGIS to get the sum and max index of each zone for three different zonal levels. The statistical parameters thus obtained were normalized and hot-spots were identified based upon these values.

2.1 Raster Based Approach for Computation of Sum and Max

In this approach a grid is drawn such that it encompasses the area over which hot-spot analysis is to be performed. An area under consideration would contain a finite number of grid cells. Sum is computed as the total number of events (crashes in our case) contained by all the grid cells enclosed within a particular zone. Similarly, Max is computed as the highest number of events enclosed by a single grid cell of a zone. In Figure 1 four different zones marked by boundaries with color red, orange, blue, and green are shown. The Sum and Max values for the four zones are (7, 5), (8, 2), (3, 1), and (1, 1) respectively. If the hot-spot zones are identified according to Sum value then the orange zone would rank at the top, while the red zone is identified as the 'hottest' zone if Max value is chosen as the criteria. Thus, the use of both Sum and Max values ensures that all the hot-spot zones are identified depending upon the objective of the ranking. High values of Sum indicate a high overall total number of events in a zone which may or may not be influenced by local high values. High values of Max on the other hand indicate the presence of a local region within a zone with a high density of events.

2.2 Normalization

The results obtained from zonal statistics are not sufficient for insightful comparisons. Zones with high values of Sum may have those values only because their area is bigger. Also, one zone may have only a single small area with high number of accidents and show a very high Max value. In order to have a better interpretation of crash patterns, we normalize the results by first converting the zonal areas which are in square decimal degrees to square miles, and then getting the Sum for one square mile of area in each of the zones. The total sum per unit square mile area is computed for each zone level. The same process is followed for Max value of each of the three zonal levels.

Since the normalized results of the sum value shows the number of accidents within a one mile square area, comparisons can be made among the different zones within the same scale, i.e. a direct inference can be drawn about whether one parish/census tract/census block group has more accidents per square mile than the other. Sum values suggest which zones have the highest number of total accidents per unit area. Max value gives us an idea about where the peak values (sub-zonal unit areas with highest number of crashes) are located and about their peak value. When interpreting the Max value, it should be remembered that the peak value is within one mile area somewhere in the corresponding zone unit and comparisons are logical only within the same zonal scale level. Figure 4 shows the results obtained by using zonal statistics tool and the normalized results based on our calculations.



Figure 1: Raster based approach for hot-spot analysis

2.3 Hot-spot Identification

Use of only one statistical parameter can be inconclusive for the identification of the hot-spot zones. If Sum is used to rank the hot-spots then zones with high peak values may be left out. Conversely, if Max is used to rank the hot-spot regions then zones with a high total number of events per unit area could be ignored. This point is illustrated in the Figure 2 and explained below.


* All Areas are equal

Figure 2: Hot-spot Identification by using Sum and Max

We can see from Figure 2 that Area 1 is composed of consistently high valued sub-zones, and its rank would be 1 if the sum value is considered. Area 2 on the other hand has one single sub-zone with a very high Max value. Area 2 would be ranked number 1 if Max value is to be considered. In Area 3 there is a high Max value as well as other sub-zonal area units with significant counts. It is therefore necessary not to rely upon a single parameter for the identification of the hot-spots. Depending on the objectives of the analysis, a list of different zones could be identified as hot-spots.

3 RESULTS

We have used a zonal statistics tool in ArcGIS to get the sum and max values of crash data at three different zonal scale levels: 1) Parish level, 2) Census Tract level, and 3) Census Block Group level. Sum and Max of the crash data from different zonal scale levels cannot be compared directly for any insightful interpretations. For each zonal scale level, we normalize the Sum and Max values generated by the zonal statistics to get the total/highest number of accidents within a one mile square area for each zone. The Sum and Max indexes are visualized per unit area of the zone. The map visualization of the normalized results using Sum and Max values are presented in Figure 5 and 6.



Figure 3: Three Different Zonal Levels

| Zonal statistics results | | | | | | | Calculations done by us | | | | | | |
|--------------------------|-----------------|----------------|-----------|-------|----------|------|-------------------------|--------|-----------|------|-------------------|------------------|------------------|
| Atte | ributes of zo | onal statistic | s table | - | | | | | | | | | |
| Ros | wid | COUNTY | ZONE-CODE | COUNT | AREA | MIII | MAX | RAIIGE | STD | SUM | AREA_III_SOR_MILE | SUM_HORM_BY_SORM | MAX_NORM_BY_SORM |
| | 17 East Bato | in Rouge | 17 | 322 | 0.078402 | 1 | 85 | 84 | 10.131253 | 1872 | 311.64963 | 6.006745 | 69.03965 |
| | 9 Caddo | | 9 | 461 | 0.112247 | 1 | 69 | 68 | 4.876241 | 1079 | 446.18161 | 2.418298 | 39.145668 |
| | 26 Jefferson | 1 | 26 | 136 | 0.033114 | 1 | 52 | 51 | 11.017445 | 1048 | 131.62842 | 7.961807 | 99.999992 |
| | 28 Lafayette | | 28 | 219 | 0.053323 | 1 | 58 | 57 | 8.894848 | 1048 | 211.96046 | 4.944318 | 69.265884 |
| | 52 St. Terrm | any | 52 | 512 | 0.124664 | 1 | 25 | 24 | 2.346378 | 945 | 495.54227 | 1.907002 | 12,770432 |
| | 10 Calcasieu | 1 | 10 | 436 | 0.106159 | 1 | 46 | 45 | 3.677274 | 932 | 421.98523 | 2.208608 | 27.593506 |
| | 40 Rapides | | 40 | 590 | 0.143656 | 1 | 32 | 31 | 2.112713 | 886 | 571.03503 | 1.551589 | 14.185136 |
| | 53 Tangipah | oa | 53 | 478 | 0.116386 | 1 | 16 | 15 | 1.477266 | 730 | 462 63519 | 1.577917 | 8.754424 |
| | 32 Livingstor | n | 32 | 405 | 0.098611 | 1 | 39 | 38 | 2.571266 | 692 | 391.98166 | 1.765389 | 25.185184 |
| | 49 St. Landr | Y | 49 | 491 | 0.119551 | 1 | 16 | 15 | 1.592983 | 688 | 475.21732 | 1.447759 | 8.522637 |
| | 8 Bossier | | 8 | 305 | 0.074263 | 1 | 35 | 34 | 3.230206 | 556 | 295.19611 | 1.883494 | 30.012606 |
| | 29 Lafourch | e | 29 | 343 | 0.083515 | 1 | 13 | 12 | 1.357036 | 552 | 331.97461 | 1.862778 | 9.912536 |
| | 35 Natchitoc | hes | 35 | 446 | 0.108594 | 1 | 14 | 13 | 0.794773 | 520 | 431.66379 | 1.204641 | 8.209727 |
| | 55 Terreboni | ne | 55 | 271 | 0.065984 | 1 | 23 | 22 | 2.298316 | 500 | 262 28897 | 1.906295 | 22.196989 |
| | 37 Ouechite | | 37 | 327 | 0.07962 | 1 | 20 | 19 | 1.636483 | 493 | 316.48892 | 1.557717 | 15.996235 |
| | 57 Vermilion | | 57 | 384 | 0.093498 | 1 | 11 | 10 | 0.987467 | 485 | 371.65671 | 1.304968 | 7.491987 |
| | 1 Acadia | | 1 | 394 | 0.095933 | 1 | 8 | 7 | 0.593844 | 460 | 381.33527 | 1.206288 | 5.310425 |
| | 3 Ascensio | n | 3 | 202 | 0.049184 | 1 | 24 | 23 | 2 507263 | 447 | 195.50691 | 2.286364 | 31.073875 |
| | 16 De Soto | | 16 | 370 | 0.090089 | 1 | 5 | 4 | 0.490531 | 417 | 358.10672 | 1.164457 | 3,534303 |
| | 58 Vernon | | 58 | 343 | 0.083515 | 1 | 18 | 17 | 1.260285 | 405 | 331.97461 | 1.219973 | 13.72505 |
| | 6 Beaurega | and | 6 | 356 | 0.086681 | 1 | 11 | 10 | 0.669995 | 399 | 344.55673 | 1.15801 | 8.081244 |
| | 5 Avoyeles | 5 | 5 | 339 | 0.082541 | 1 | 5 | 4 | 0.457731 | 384 | 328.10318 | 1.170364 | 3.857499 |
| | 36 Orleans | | 36 | 124 | 0.030192 | 1 | 32 | 31 | 4,6381 | 380 | 120.01414 | 3.166293 | 67.49379 |
| | 50 St. Martin | Ē. | 50 | 249 | 0.060628 | 1 | 13 | 12 | 1.246236 | 354 | 240.99614 | 1.468903 | 13.654618 |
| | 23 Iberia | | 23 | 198 | 0.04821 | 1 | 20 | 19 | 2.127696 | 348 | 191.6355 | 1.815948 | 26.418024 |
| | 20 Evangelin | e | 20 | 288 | 0.070124 | 1 | 9 | 8 | 0.675149 | 332 | 278.74252 | 1.191063 | 8.173077 |
| | 60 Webster | | 60 | 235 | 0.057219 | 1 | 17 | 16 | 1.602251 | 324 | 227.44617 | 1.424513 | 18.919802 |
| | 45 St. Charle | 15 | 45 | 139 | 0.033844 | 1 | 16 | 15 | 2.396126 | 317 | 134.53198 | 2.356317 | 30.105145 |
| | on here shis as | ** | ~ | | 0.070047 | | 1 | | n 200004 | 240 | 200.2004 | 4 004030 | * 2007* |

Figure 4: Zonal Statistics and Normalized Results



Figure 5: SUM (Census Block Groups Scale, zoomed at a single parish)



Figure 6: MAX (Census Block Groups Scale zoomed at a single parish)

4 CONCLUSION

We have presented a way to improvise the results generated by the zonal statistics in ArcGIS. The normalization step that we propose following the zonal statistics computation is used for hot-spot analysis to find out patterns and/or trends of crashes. Our effort was to get a better understanding of the results generated by the zonal statistics tool to identify the traffic accident hot-spot zones. Our approach makes it possible for comparison of statistical parameters across different zone areas within the same scale. This helps decision makers to pinpoint locations with a higher crash index and take preventive measures. The road condition can be improved, proper signs can be placed and, traffic laws can be imposed more strictly in the areas determined as having high risk by the analysis. This would not only help to prevent loss of millions of dollars' worth of property that are damaged due to crashes, but more importantly save thousands of innocent lives each year.

ACKNOWLEDGEMENT: The authors of this paper would like to thank LADOTD for supporting this work.

5 REFERENCES

- [1] T. K. Anderson, "Kernel density estimation and K-means clustering to profile road accidenthotspots," Accident Analysis & Prevention, 41(3), pp. 359-364, 2009.
- T. Anderson, "Comparison of spatial methods for $\mathbf{6}$ [2] measuring road accident 'hotspots': a case study of London," Journal of Maps, 3 (1), 2007.
- Z. Xie and J. Yan, "Kernel Density Estimation of [3] Environment and Urban Systems, 32 (5), pp. 396-406, 2008.
- [4] L. Schweitzer, "Environmental justice and hazmat transport: A spatial analysis in southern California," Transportation Research Part D: Transport and Environment, 11 (6), pp. 408-421, 2006.
- P. D. Bates and A. P. J. De Roo, "A simple raster-[5] based model for flood inundation simulation," Journal of Hydrology, 236 (1-2), pp. 54-77, 2000.
- [6] Topographic Structure from Digital Elevation Data for Geographic Information System Analysis," Photographic Engineering and Remote Sensing, pp. 1593-1600, 1988.
- P. Y. Julien, B. Saghafian and F. L. Ogden, [7] Resources Association, pp. 523-536, 1995, 31 (3).
- [8] C. Ratti and P. Richens, "Raster Analysis of Urban Form," Environment and Planning B: Planning and Design, pp. 297-309, 2004, Vol 31.
- [9] I. S. Sotheran, R. L. Foster-Smith and J. Davies, "Mapping of marine benthic habitats using image processing techniques within a raster-based geographic information system," Estuarine, Coastal and Shelf Science, pp. 25-31, 1997, 44 (1).

- [10] ESRI, "An Overview of the Zonal tools," [Online]. Available: http://webhelp.esri.com/arcgisdesktop/9.3/index.cfm?Top icName=An_overview_of_the_Zonal_tools. [Accessed 16 May 2012].
- ESRI, "How Zonal Statistics Works (ArcGIS [11] Desktop Help)".

Appendix

6.1 **Zonal Statistics**

Zonal analysis is an important analysis tool in ArcGIS trafficaccidents in a network space," Computers, under its spatial analyst extension. It is useful for several types of GIS-related analysis or studies such as environmental monitoring, demographic studies, land management, traffic data analysis and so on. Zonal analysis is the creation of an output raster (or statistics table) in which the desired function is computed on the cell values from the input value raster that intersect or fall within each zone of a specified input zone dataset. The zonal tools in which the zones are defined by a single input value raster either calculate statistics or quantify the characteristics of the geometry of the input zones.

Zonal tools are categorized as zone shapes (Zonal Area, S. K. Jenson and J. O. Domingue, "Extracting Zonal Centroid, Zonal Perimeter), zone attributes (Zonal Max, Zonal Min, Zonal Sum), and determine specified zones (Zonal Fill). The Zonal Statistics tool records the specified statistic of the values of all cells in the value Dataset that belong to the same zone as the output cell in each output cell "Raster-Based Hydrological Modeling of Spatially-varied [10]. An input zone for the zonal statistics tool can be a vector Surface Runoff," Journal of the American Water file or an integer raster file. Similarly, an input value raster includes any raster file that contains values that can be analyzed visually and statistically. The input zone dataset is only used to define the size, shape, and location of each zone, while the input value raster contains the values to be used in the evaluations within the zones.

> Maximum: For this statistic, the zone input must be an integer. The data type of the output will be the same as the value input.



Figure 7: Maximum Calculation (Zonal Statistics) [11]

In our work, Maximum value gives an idea about where the peak values are located and what the peak value is. During interpretation of the maximum value, it should be remembered that the peak value is within a one mile area somewhere in the corresponding zone unit.

Sum: For this statistic, the zone input must be an integer. The data type of the output raster is a floating point. This is

because the value for the Sum tends to be quite large, and may not be possible to represent with an integer value.

For example, if a zone has 2500 rows and columns of cell in size, and the value of each cell is 1000, the sum for that zone would be 2500*2500*1000 = 6.25 billion [11]. If an integer output is required and the range is within ± 2.147 billion, then the INT function can be applied [11].



Figure 8: Sum Computation (Zonal Statistics) [11]

In our work Sum value shows the number of accidents within a one mile square area.

A Graph Grammar Model of Financial Statements with Heterogeneous Parts

Takeo Yaku¹, Koichi Anada², Koushi Anzai³, Shinji Koka¹, Miyuki Shimizu¹, and Yuki Shindo¹

¹Department Computer Science and System Analysis, Nihon University, Tokyo, Japan ²Waseda Research Institute for Science and Engineering, Waseda University, Tokyo, Japan ³Department Economics, Kanto Gakuen University, Gunma, Japan

Abstract - Spreadsheets are frequently used in information processing. In those information processing, automatic generation and automatic verification of spreadsheets are required. However, it is difficult to verify the structure of spread sheets with flexible calculation ranges and with heterogeneous structures. Accordingly, the modeling of spreadsheets is important. We deal with the modeling of spreadsheets such as financial statements. In this paper, We formalize the structure of financial statements with heterogeneous parts using a context sensitive graph grammar. We propose 62 rewriting rules that provide the spatial order of items of heterogeneous financial statements. Furthermore we also show derivation of financial statements with the grammar.

Keywords: modeling of spreadsheets, context sensitive graph grammar, financial statements, syntax-directed recognition of logical structures

1 Introduction

The formalization of business documents as in Figure 1 has become an important subject with the progress of ecommerce and e-government (see, e.g. [3]). In order to formalize financial statements, we have to specify the spatial order of items, and specify calculation methods of categorized items. A context sensitive graph grammar was proposed in [5] that specify financial statements only with homogeneous part.



Figure 1. A form of financial statements with heterogeneous parts.

In this paper, we formalize the structure of financial statements with heterogeneous parts using a *context sensitive graph grammar* (CSGG) (see, e.g. [1]). We first propose rewriting rules that provide the spatial order of items in the financial statements. Next, we also show a derivation process that provides syntax directed recognition process of logical structures.

2 Graph representation of financial statements

We represent financial statements by octgrids [4]. Figures 2 and 3 show a financial statement and its corresponding octgrid. Each node in Figure 3 corresponds to each cell of the financial statement in Figure 2.

| | | 2011 | 2012 | change |
|-----------|----------------------------|------|------|--------|
| Section 1 | Domestic Operations | | | |
| | Overseas Operations | | | |
| | Sub total | | | |
| Section 2 | Domestic Operations | | | |
| | Overseas Operations | | | |
| | Sub total | | | |
| Total | | | | |





Figure 3. The octgrid for Figure 2.

3 Graph grammar for heterogeneous parts

Financial statements have the two dimensional grid structures. Accordingly they could not be formalized by context free graph grammars (see, e.g. [2]). Thus, we construct a CSGG that formalize financial statements with heterogeneous parts (cf. [4]).

3.1 **Rewriting Rules**

The CSGG for financial statements is a system $GG_F = (N_F, T_F, M_F, P_F, S_F)$, where $N_F = \{S, \rightarrow, \downarrow, tr\}$, $T_F = \{DO, OO, St, DO', OO', St'', ye, ch, ch'', Em, Se, To, To'', pe\}$, $M_F = \{nwe, swe, ewe, wwe\}$ and $S_F = S$.

The terminal label DO stands for "Domestic Operations", OO for "Overseas Operations", St for "Sub total", ch for "change", Em for an empty cell, Se for "Section", To for "Total", and pe for "perimeter". Fig. 4 shows a part of the production rules of GG_F . Each production rule means the rewriting of the left-hand side graph by the right-hand side graph, where vertices are accompanied by their vertex number.



Figure 4. Production rules in GG_F .

3.2 Derivation

Following Fig. 5 shows a derivation in GG_F . The derivation is consist of following Phases 1~20.

- Phase 1: Generate the frame and determine the number of columns.
- Phase 2: Change labels of nodes in the 1st and the last line.
- Phase 3-4: Generate the 2nd line and change labels of nodes in the 2nd line.
- Phase 5-6: Generate the 3rd line and change labels of nodes in the 2nd line into terminal ones.
- Phase 7-8: Generate the 4th line and change labels of nodes in the 3rd line into terminal ones.

Phase 9-19: Repeat to generate lines and change labels. Phase 20: Change labels of nodes in the last line into "perimeter"



Figure 5. A derivation process for Figure 2 by GG_F .

4 Conclusion

In this paper, we constructed a CSGG, GG_F , which generates graphs for heterogeneous financial statements. GG_F has 62 rewriting rules. Furthermore, we showed process. However, it is not verified whether all of heterogeneous financial statements are generated by those 62 rules or not. As future works, we discuss attribute rules that formalize scopes of spread sheet calculation defined by GG_F .

We would like to thank Professors G. Akagi and K. Tsuchida for their valuable suggestions.

5 References

[1] Y. Adachi, S. Kobayashi, K. Tsuchida, and T. Yaku, An NCE Context-Sensitive Graph Grammar for Visual Design Languages. *Proc. VL*, 1999, pp. 228-235.

[2] T. Arita, K. Tsuchida, and T. Yaku, Syntactic Characterization of Two Dimensional Grid Graphs. *IEICE TRANS*. Vol E89-D/2, pp. 771-778.

[3] M. Burnett, A. Sheretov, and G. Rothermel, Scaling Up a "What You See Is What You Test" Methodology to Spreadsheet Grids. *IEEE Symp. on Visual Languages*, pp.30-37, 1999.

[4] M. Shimizu, G. Akagi, K. Tsuchida, and T. Yaku, A Formalization of Business Documents by an Attribute Graph Grammar. *IPSJ SIG Tech. Rep.* Vol.2006 No.29, 2006-MPS-058, pp.39-42.

[5] Y. Shindo, K. Anada, K. Anzai, S. Koka, and T. Yaku, A Graph Grammar Model for Syntaxes of Financial Statements, *Proc. VL/HCC* 2011, pp. 265-266.

SESSION

MODELING, SIMULATION AND VISUALIZATION METHODS: NOVEL ALGORITHMS AND APPLICATIONS

Chair(s)

Prof. Hamid R. Arabnia

Analysis, Design, and Simulation of a Novel Current Sensing Circuit

Louiza Sellami Electrical and Computer Engineering Department US Naval Academy Annapolis, Maryland, USA sellami@usna.edu

Abstract— This paper investigates the behavior of a novel active current sensing circuit based on the circuit proposed by Maciej Kokot. The circuit is a modification of an op-amp low side current sensor whereby the output sinks the current to ground through a resistor. Supporting analysis, design equations, as well as Multisim simulation results are included.

Keywords: Current Sensors, Op-Amp Circuits, Design, Simulation, Multisim

1 Introduction

A current sensor is a device that detects and converts current to an easily measured output voltage, which is proportional to the current through the measured path. There are a wide variety of sensors, and each sensor is suitable for a specific current range and environmental condition. No one sensor is optimum for all applications.

Current measurement, current sensing, or the monitoring of current flow into and out of electronic circuits is a fundamental requirement in a wide range of electronic applications. Typical applications that benefit from current sensing include battery life indicators and chargers, current and voltage regulators, DC/DC converters, ground fault detectors, handheld communications devices, medical diagnostic equipment, motor speed controls and overload protection, battery chargers, and battery-operated circuits for which one must know the ratio of current flow into and out of a rechargeable battery. As more applications become portable, the demand increases for dedicated current monitors that accomplish their task in a small package and with low bias current [1]-[4].

A particular current sensing circuit of interest is the one proposed by Kokot in [5] whereby by combining the Robert W. Newcomb Electrical and Computer Engineering Department The University of Maryland College Park, Maryland, USA newcomb@umd.edu

advantages of low-side and high-side current sensors, it is possible to come up with a simple current sensor with a return current and zero input voltage. The circuit uses two op-amps configured so that the top op-amp sinks the current to the output of the second op-amp, while the latter returns an equal current. We point out that [5] does not provide any analysis or design method for implementing the circuit or to support the claims above.

In this paper we present a novel a current sensing circuit, develop a design method through analysis, as well as simulate the circuit. Specifically, we derive the conditions on the resistors to guarantee that the return current matches the input current, on the one hand, and, on the other hand, that the sensed current at the output is also the same as the input current. Further, we examine the conditions on the input signal to prevent the op-amps from saturating. Also, the non-ideal case of the effect of small voltage perturbation at the non-inverting terminal of the top op-amp on the sensed current is examined.

2 Basic Current Sensors

There are two types of current sensing circuits: passive and active. Current sensing resistor is of the passive type, and is the most commonly used. It can be considered a current-to-voltage converter, where inserting a resistor into the current path, the current is converted to voltage in a linear way (Ohm's law). The main advantages of current sensing resistors include low cost, high measurement accuracy, measurable current range from very low to medium, and the capability to measure DC or AC currents. The disadvantages include the introduction of an additional resistance into the measured circuit path, which may increase source output resistance and result in undesirable loading effect, power loss due to dissipation [3]. Therefore, current sensing resistors are rarely used beyond the low and medium current sensing applications.

Two techniques for active current sensing applications are used: low-side current sensing and highside current sensing. Each technique has its own advantages and disadvantages, discussed in more detail in the following sections. In both methods current is converted to a proportional voltage that is readily measurable.



Figure 1: Low-side Current-to-voltage Converter.

2.1 A Simple Low-side Current Sensing

As shown in Figure 1, a simple low-side current sensor is the typical operational amplifier current-tovoltage converter in which the op-amp sinks the incoming current through the feedback resistor. In general, low-side current sensing connects the sensing or probing resistor between the load and (virtual) ground. Normally, the sensed voltage signal is so small that it needs to be amplified by subsequent op amp circuits (e.g., non-inverting amplifier) to get the measurable output voltage. Advantages include low input common mode voltage, ground referenced input and output, simplicity and low cost. Disadvantages include ground path disturbance, load is lifted from system ground since it adds undesirable resistance to the ground path, and high load current caused by accidental short goes undetected. Considering the advantages, low-side current sensing are used where short circuit detection is not required, and ground disturbances can be tolerated [2].

2.2 A Simple High-side Current Sensing

High-side current sensing is typically selected in applications where ground disturbance cannot be tolerated, and short circuit detection is required, such as motor monitoring and control, overcurrent protection and supervising circuits, automotive safety systems, and battery current monitoring. Figure 2 shows a high-side current sensing whereby the sensing or probing resistor is inserted in the current path. The circuit uses a current-tovoltage differential amplifier to measure the resulting voltage drop. Typically, the sensed voltage signal is amplified by subsequent op amp circuits to get the measurable output voltage. Advantages include eliminating ground disturbance, load connects system ground directly, detects the high load current caused by accidental shorts. The drawback is that it must be able to handle very high and dynamic common mode input voltages, in addition to complexity and higher costs [3].



Figure 2: High-side Differential Amplifier Current to Voltage Converter.

3 A Novel Current Sensing Circuit

Comparing the two sensing circuits of Figures 1 and 2, we note that in the first circuit the input voltage drop is on the order of microvolts to millivolts, depending on the quality of the IC, but the measured current flows into the sensing node with no return current to the circuit. With this method one can only measure currents flowing to ground. In the second circuit, however, the same current flowing in flows out of the circuit, but a significant voltage drop occurs at the input across resistance R_p. A circuit that combines the benefits of a return current and zero input voltage is shown in Figure 3 [5]. The circuit in Figure 1 in the sense that the top op-amp sinks the input (measureable) current. However, the bottom op-amp output sources back an equal outgoing current.

In Section 3, we analyze the behavior of this circuit, namely derive the conditions on the resistors to guarantee that the return current matches the input current, on the one hand, and, on the other hand, that the sensed current at the output is also the same as the input current. Further, we examine the conditions on the input signal to prevent the op-amps from saturating. Except for the third case, we assume ideal op-amps, meaning currents into the op-amp terminals are zero, and the voltage difference between the op-amp terminals is also zero. Note that this is not always the case because the circuit uses both positive and negative feedback.



Figure 3: A Novel Current Sensing Circuit.

3.1 Conditions to force the return and the input currents to be the same

By virtue of the ideal characteristics of the op-amps, $v_3 = v_2 = v_4$, $i_i = i_3$, $i_0 = i_4$, and $i_1 = i_2$. By Ohm's Law, we have

$$\dot{i}_i = \dot{i}_3 = \frac{v_3 - v_{02}}{R_2}.$$
 (1)

$$i_0 = i_4 = \frac{v_{01} - v_4}{R_4} \cdot$$
(2)

$$\dot{i}_1 = \dot{i}_2 = \frac{v_{01} - v_2}{R_1} = \frac{v_2 - v_{02}}{R_2}.$$
(3)

Using (3) to solve for v_2 and substituting its expression in (1) and (2) gives:

$$i_i = \frac{R_2(v_{01} - v_{02})}{R_3(R_1 + R_2)}$$
 and $i_0 = \frac{R_1(v_{01} - v_{02})}{R_4(R_1 + R_2)}$.
(4)

In order for the two currents in (4) to be equal, it is required that $\frac{R_1}{R_4} = \frac{R_2}{R_3}$.

3.2 Conditions to force the sensed and the input currents to be the same

Defining $v_{01} = v_{01} - v_{02}$ and on using (4) and Ohm's Law, we get:

$$v_0 = \frac{R_3(R_1 + R_2)}{R_2} i_i = (R_1 + R_2) i_1 \cdot$$
(5)

From (5), to achieve i1=ii = io, we require $\frac{R_3}{R_2}=1$.

Combining this condition and the previous result yields the general condition that all four resistors must be equal in order for the sensed current to be the same as the input current which is equal to the return current.

3.3 Conditions to avoid op-amp saturation

Since the output of the op-amps cannot exceed the rail voltages V_{s+} and V_{s-} , and also for testing purposes a voltage source v_i with internal resistance R would be connected at node v_3 while v_4 would be grounded, it follows that v_2 and v_3 become virtual grounds, and $v_{01} = -v_{02}$. Also since

$$i_{i} = \frac{v_{i}}{R} = i_{1} = \frac{v_{0}}{(R_{1} + R_{2})} \Longrightarrow \frac{R}{(R_{1} + R_{2})} V_{s}^{-} \prec v_{i} \prec \frac{R}{(R_{1} + R_{2})} V_{s}^{+}$$
(6)

3.4 Effect of perturbation in v3 on the input current

We keep the same set up as in section 3.3 above. For analysis purposes, we assume the op-amps ideal, except for a disturbance which causes v3 to deviate from v2 (=v4=0) by an amount δ , so that $v3' = v3 + \delta$. This in turn will cause the input current to deviate by a certain amount, which is determined by substituting v3' in (1). The resulting current is then given by:

$$\dot{i}_i = \dot{i}_i + \frac{\partial}{R_3}.$$
(7)

Note that this perturbation has no effect on the sensed and return current which are still given by (3) and (4). As seen in (7), in order to minimize the deviation of the input current, R3 has to be large, but still in the $k\Omega$ range. This way the input, return and sensed currents are kept the same.

4 Simulation Results

Multisim simulations both in DC (though we do not include the results here) and AC (transient analysis in Multisim) were conducted on the circuit of Figure 4, with the indicated values for the resistances and input signal, the choice of which was dictated by the conditions derived in section III. The op-amps were a pair of 741's (though other ICs may be used) biased symmetrically with $\pm 15V$. The graphs of Figure 5 show the plots of the input and sensed currents which are equal. Figure 6 shows the plots of the input and returned currents, which are equal but out of phase since they flow in opposite directions. Though not shown here, currents through R, R1, R2, and R3 are all equal, thus verifying the analysis and design equations.

5 Discussion

Though the simulation results show that the circuit performs according to the theory, physical limitations of the devices will undoubtedly affect the results, especially if the values of the resistances are not chosen properly. For instance, values of R1, R2, R3, and R4 in the Ohm range lead to the saturation of the op-amps, as the differential terminal input voltage deviates from zero. Also the frequency response of the devices limits the overall frequency response for the circuit.

Though the effect of the frequency of the input signal on the discrepancy between the input, sensed and output currents was not analyzed here, a Multisim simulation was done, and the results are shown in Figures 7 and 8. It is noted that even beyond the cutoff frequency of the op-amps, the magnitudes of the input and sensed currents remain the same. That is not the case of the phase which starts to deviate around 40 kHz, which incidentally is also above the cutoff. Similar behavior is observed with the return current also, as illustrated in Figure 8.



Figure 4: Simulated Circuit with input voltage applied to node 3 via the source resistance, and node 4 connected to ground.



Figure 5: Plots of the input and sensed currents.



Figure 6: Plots of the input and return current which here is the same as I(R4).



Figure 7: Frequency response comparing the magnitude and phase of the input current to the sensed current



Figure 8: Magnitude plot of the input current compared to the return current.

6 References

[1] Robert F. Coughlin and Frederick F. Driscoll, Operational Amplifiers and Linear Integrated Circuits, 5th ed., Prentice Hall, 1998, pp.125–139.

[2] Neil Zhao, Wenshuai Liao, and Henri Sino "High-Side Current Sensing with Dynamic Range: Three Solutions," Analog Dialogue, 2010, pp. 1-5.

[3] B. Yuan and X. Lai, "On-chip CMOS Current sensing for DC-DC Bucj converter," Electronics Letters, Vol. 45, No. 2, 2009, pp. 102-103.

[4] B. Giavanni et al., "MAGFET based current sensingfor power integrated circuits," Microelectronics Reliability, Vol 43, No. 4, 2003, pp. 577-584.

[5] Maciej Kokot, "Measuring small currents without adding resistive insertion loss,"EDN Europe.

Visualization, Analysis and Error prediction of Supersonic Shock Angles over diamond shape Aerofoil Using Hydraulic analogy

Akshay Garg Aerospace Engineer, Mahindra Satyam, Bangalore, INDIA

Abstract - The analogy between the shallow water flow with the free surface and two-dimensional gas flow has been found to be useful for qualitative as well as quantitative study of high speed flows. This technique is valued highly because of the fact that many practical problems in supersonic flows involving shock and expansion waves, which require a sophisticated and expensive wind tunnel and instrumentations for their analysis, may be studied in an inexpensive manner with simple water channel facility. We have considered two values at Mach numbers, 1.5 and 2 and shown the results. Hence the theoretical and practical phenomenon matches and desirable results achieved.

Keywords: Visualization, Supersonic Shock, MSV'12, Hydraulic Analogy, Shallow Water

1 Introduction

The flow visualization plays a dominant role by way of paving a path for the development of techniques for solving important but complex problems such as separated flows, jets, and so on. In fact, in many flow situations the information given by the flow visualization cannot be obtained by any other technique or instrument. However, these visualizations lack authentic mathematical procedure. They simply serve as a tool to understand certain complex flow fields in a simple manner. To overcome this shortcoming of visualization, certain methods were developed to study the fluid flow problems without actually going into the complexities involved. They are termed *analogue methods*.

In analogue methods, fluid flow problems are solved by setting up another physical system, such as an electric filed, for which the basic governing equations are of the same form with corresponding boundary conditions as those of the fluid flow. The solution of the original problem may then be obtained experimentally from measurement on the analogous system. Some of the well known analogy methods for fluid flow problems are the Hele-Shaw analogy, electrolytic tank, and surface waves in a ripple tank [4]

In this project, we discussed the various aspect of visualization of shocks around the body. We experiment on

wedge and semi-wedge airfoil under two conditions i.e. M = 2 and M = 1.5

2 Hydraulic Analogy

The analogy between shallow water flow with a free surface and two-dimensional gas flow has been found to be useful for qualitative as well as quantitative study of high speed flows. In principle, any regulated stream of shallow water can be used for this analogy. In laboratories it is usually done in a water flow table. This technique is valued highly because of the fact that many practical problems in supersonic flows involving shocks and expansion waves that require a sophisticated and expensive wind tunnel and instrumentation for their analysis may be studied in an inexpensive manner with a simple water flow channel facility. [1]

The analogy may be used to understand many challenging problems of practical interest, such as transonic unsteady flows past wings, transient phenomena in highspeed flow and shock wave interaction, and so on. Simulation of such flows in a wind tunnel is highly complicated and costly and hence primary information explaining the basic nature of such flows will be of immense use in solving such problems. In such situations, the analogy channel proves it handy and helps us by providing preliminary information and a qualitative picture of the flow field.

The essential feature of this analogy is that the Froude number of shallow water flow with a free surface is equivalent to a gas stream with Mach number equal to that of Froude number.

2.1 Theory of the analogy

The basic governing equations of flow of an incompressible fluid with a free surface, in which the depth of flow is small compared to its surface wave length, forms the shallow water theory. The similarity between the governing equations of motion of a two-dimensional isentropic flow of a perfect gas and two-dimensional shallow water flow forms the analogy.

2.2 Shallow Water Flow

Consider a continuous flow (without hydraulic jump) of a perfect liquid in the absence of any external forces such as electromagnetic forces, viscous forces, surface tension, and so on, except gravitational force. For this flow, neglecting the variation of pressure (if any) on the free surface and treating the flow as a two-dimensional with constant energy, we can write the basic governing equation as follows-[7] The continuity equation is-

$$\frac{\partial h}{\partial t} + \frac{\partial (hu)}{\partial x} + \frac{\partial (hv)}{\partial y} = 0. \qquad 2.2(a)$$

The momentum equations along the flow (x) and transverse (y) directions, respectively, are-

$$\frac{\partial u}{\partial t} + u \frac{\partial u}{\partial x} + v \frac{\partial u}{\partial y} = -\frac{g}{2} \cdot \frac{1}{h} \cdot \frac{\partial (h^2)}{\partial x}$$
$$\frac{\partial v}{\partial t} + u \frac{\partial v}{\partial x} + v \frac{\partial v}{\partial y} = -\frac{g}{2} \cdot \frac{1}{h} \cdot \frac{\partial (h^2)}{\partial y}$$
$$2.2(b)$$

Where u and v are the velocity components along x- and ydirections, respectively, g is the acceleration due to gravity and h is the depth of the fluid stream.

Now, let us define a fictitious pressure $p_{\rm f}\,$ and a fictitious density $\rho_{\rm f}$ as

$$p_{f} = \int_{o}^{h} p \, dz$$
$$= \int_{o}^{h} \rho g z \, dz$$
$$p_{f} = \rho h,$$

Where ρ is the density of water and z is the coordinate in the depth direction. Introduction of these expressions for p_f and ρ_f in to Eq.2.2(a) yields –

$$\frac{\partial(\frac{\rho f}{\rho})}{\partial t} + \frac{\partial(\frac{\rho f}{\rho}u)}{\partial x} + \frac{\partial(\frac{\rho f}{\rho}v)}{\partial y} = 0.$$

This simplifies to -

$$\frac{\partial(\rho f)}{\partial t} + \frac{\partial(\rho f u)}{\partial x} + \frac{\partial(\rho f v)}{\partial y} = 0. \qquad 2.2(c)$$

Because ρ is a constant.

Similarly, introducing p_f and ρ_f , Eqs. 2.2(b) can be reduced to

$$\frac{\partial u}{\partial t} + u \frac{\partial u}{\partial x} + v \frac{\partial u}{\partial y} = -\frac{1}{\rho f} \cdot \frac{\partial \langle Pf \rangle}{\partial x}$$
$$\frac{\partial v}{\partial t} + u \frac{\partial v}{\partial x} + v \frac{\partial v}{\partial y} = -\frac{1}{\rho f} \cdot \frac{\partial \langle Pf \rangle}{\partial y}$$
$$2.2(d)$$

2.3 Gas flow

Consider a two-dimensional isentropic flow of a perfect gas in the absence of any external forces such as

electromagnetic force, gravitational force and so on. The basic equations for this flow are the following.

The continuity equation is

$$\frac{\partial(\rho g)}{\partial t} + \frac{\partial(\rho g u)}{\partial x} + \frac{\partial(\rho g v)}{\partial y} = 0. \qquad 2.3 \text{ (a)}$$

And the momentum equations are

$$\frac{\partial u}{\partial t} + u \frac{\partial u}{\partial x} + v \frac{\partial u}{\partial y} = -\frac{1}{\rho g} \cdot \frac{\partial (P)}{\partial x}$$
$$\frac{\partial v}{\partial t} + u \frac{\partial v}{\partial x} + v \frac{\partial v}{\partial y} = -\frac{1}{\rho g} \cdot \frac{\partial (P)}{\partial y} \qquad 2.3 \text{ (b)}$$

Where ρ_g is the density of gas.

Comparing the Eqs. 2.3(a) and 2.3(b), it can be seen that the shallow water flow equations are similar to the gas flow equations. In terms of ρ_w and ρ_f , the fictitious pressure p_f can be expressed as[2]

$$P_{\rm f} = \rho_{\rm w} \, {\rm g} \, {\rm h}^2 / {\rm 2} = {\rm g} \, \rho_{\rm w} (\rho_{\rm f})^2 / 2(\rho_{\rm w})^2 = {\rm g}(\rho_{\rm f})^2 \rho_{\rm w} / {\rm 2}(\rho_{\rm w})^2$$

 $P_{\rm f}/(\rho_{\rm f})^2 = {\rm g}/2\rho_{\rm w} = {\rm constant.}$

This expression corresponds to the polytrophic process equation for gases with polytrophic index n=2. This imposes a further condition to be satisfied for complete analogy between these two streams of gas and water. Here we should note that the value of index n can never be equal to two for isentropic flow of gases. This index becomes γ , namely the isentropic index, which is the ratio of specific heats C_p and C_v . Also, we know that γ varies from 1 to 1.67 only. Therefore, n=2, required for hydraulic analogy is bound to introduce some error in the quantitative results obtained with hydraulic analogy. In spite of this shortcoming, the hydraulic analogy is found to yield results with reasonable accuracy. This aspect as well as the ease with which the results for gas flows at high speeds can be obtained with this technique makes it attractive, in spite of the above mentioned drawback.

It can be seen from the definition of fictitious pressure and density that the gas pressure is proportional to the square of water stream depth h and the gas density is directly proportional to the water stream depth. That is, we can express that

$$P/P_0 = h^2/(h_0)^2$$
 2.3(c)

$$P/P_0 = h/h_0$$
 2.3(d)

By the perfect gas state equation, we have

$$P = \rho RT.$$

Using this Eq.2.3(c), we obtain

 $T/T_0 = h/h_0$ 2.3(e)

Where p_0 , ρ_0 , T_0 , h_0 are the stagnation pressure, stagnation density, stagnation temperature of the gas and the stagnation depth of water stream, respectively.

In gas flow, the velocity of propagation of a small pressure disturbance is the speed of sound 'a'. The corresponding parameter in hydraulic flow is a_w , the velocity of propagation of a surface disturbance with a large wavelength compared to water depth. The wave speeds 'a' and a_w may be expressed as

$$a = (\gamma P / \rho)^{1/2}$$
$$a_w = (2g\rho_w h^2 / 2\rho_w h)^{1/2} = (gh)^{1/2}$$

We know from our basic studies on fluid flows that it is possible to identify a group of dimensionless parameters between any two dynamically similar flows. Therefore, from the above discussion on gas and water flows, it is possible to identify the Mach number and Froude number as the nondimensional similarity parameters, respectively, for gas and water streams. That is, for the gas and water streams to be dynamically similar, the Mach number of the gas stream and the Froude number of the water stream must be equal. The Mach number M and the Froude number F may be expressed as

$$M = V_g / a \qquad 2.3(f)$$

$$Fr = V_w / (gh)^{1/2}$$
 2.3(g)

Where V_g and V_w are the velocities of gas and water streams, respectively.

- When M<1 the flow is termed *subsonic*.
- When M>1 the flow is termed *supersonic*.
- When Fr<1 the flow is termed *subcritical*.
- When Fr>1 the flow is termed *supercritical*

Furthermore, there exists a degree of similarity between the flow patterns in the water and gas streams. In supersonic flows, $\mu_g = \sin^{-1}(1/M)$ is the smallest angle at which disturbance prevails in the flow with a given Mach number M, where μ_g is the Mach angle. The analogous quantity in shallow water is $\Psi_w = \sin^{-1}(1/Fr)$, where Ψ_w is called Wave angle.[3] At this stage, it must be noted that in the above analysis of similarity between shallow water and gas streams, nowhere was the assumption of zero vorticity made. The expression for vorticity in a two-dimensional gas flow as well as in a water flow can be written as-

$$\operatorname{CurlV=K}\left(\frac{\partial v}{\partial x} - \frac{\partial u}{\partial y}\right) \qquad 2.3(h)$$

Where V is the resultant velocity, u and v are the velocity components along the x- and y- directions, respectively, and k is the directional vector in the z- axis direction. Hence, in a principle it is possible to compare rotational flow in the two media. However, in practice the real fluid effects complicate the comparison. From the above discussion it is evident that with $\gamma=2$ the governing equations for the gas flow are identically the same as those for shallow water stream. The analogy holds for flows with Mach number smaller and greater than 1.

Table 2.1: Analogy comparison between the Gas dynamics and hydraulic dynamics[5]

GAS DYNAMICS
DYNAMICSHYDRAULIClocal speed of sound $a = \sqrt{KRT}$ \sim local speed of wave propagation $C = \sqrt{\frac{gh}{n+1}}$ local speed of sound $a = \sqrt{KRT}$ \sim local speed of wave propagation $C = \sqrt{\frac{gh}{n+1}}$ local Mach number $M = \frac{u}{a}$ = local Mach number $M = \frac{u}{C}$ local density ratio ρ' = local water-depth ratio h'^{n+1} local pressure ratio ρ' = local water-depth ratio h'^{n+2} local temperature ratio T'= local water-depth ratio h'flow similarity number $\left(\frac{agl_a}{l_a}\right)$ = flow similarity number $\left(\frac{\sqrt{\frac{gh_a}{n+1}}}{l_w}\right)$ function of gas specific heat ratio $\left(\frac{1}{K-1}\right)$ = function of channel cross-section exponent

2.4 Basic Assumptions

- The fluid is frictionless so that the conversion of mechanical energy into heat is excluded both in gas and in water.
- The flow is one-dimensional and is in a duct (for gas) or channel (for water) of uniform cross section. This implies that the transverse components of fluid velocity are negligible compared to the axial velocity.
- The viscosity and the surface tension of the water are neglected. In the analysis of waves, they are considered to be long waves so that the change on elevation of water is considered to be small in

comparison to h_0 , the water depth at equilibrium position.[5]

• The vertical acceleration of water is negligible compared to the acceleration due to gravity. Under this assumption the static pressure at a point in a field flow is assumed to depend linearly on the vertical distance under the free surface at that position. In other words,

$$\mathbf{p}=\mathbf{\rho}\mathbf{g}~(\mathbf{h}-\mathbf{y}).$$

It is further assumed that the velocity is uniform and constant over any cross section perpendicular to the flow direction. Justification for this assumption for the case one-dimensional unsteady flow

 In one-dimensional unsteady flow of gas dynamics, all parameters (pressure, temperature, velocity) are assumed to be uniform and constant across any section perpendicular to the direction of flow.

Table 2.2: Analogy Implication[5]

Analogous equations and variables for two-dimensional gas and shallow hydraulic flows

| Gas flow | Hydraulic flow | Implication |
|---|---|---|
| $\frac{T}{T_0} = 1 + \frac{\gamma - 1}{2} M a^2$ $\frac{\partial}{\partial x} (\rho u) + \frac{\partial}{\partial y} (\rho v) = 0$ $\frac{\rho}{p_0} = \frac{\rho}{\rho_0} \frac{T}{T_0}$ | $\frac{h}{h_0} = 1 + \frac{1}{2}Fr^2$ $\frac{\partial}{\partial x}(hu) + \frac{\partial}{\partial y}(hv) = 0$ | $\frac{\frac{T}{T_0} \equiv \frac{h}{h_0}; \ \gamma = 2}{Ma \equiv Fr; \frac{\rho}{\rho_0} \equiv \frac{h}{h_0}}$ $\frac{\frac{p}{p_0} = (\frac{h}{h_0})^2}{\frac{p}{p_0} = (\frac{h}{h_0})^2}$ |

3 Application of Hydraulic Analogy

3.1 Aerodynamic Forces on Airfoils

Predictions of aerodynamic force coefficients, namely C_L and C_D , of wings in steady two-dimensional flow are of great importance in the field of aerodynamics. Extensive exact and approximate theories have been developed for this study. The purpose of the present experiment is to show the validity of the analogy to steady supersonic flow past airfoils.

The airfoil chosen is half-wedge airfoil. The aerodynamic coefficients can be determined accurately, using the shock expansion theory. The necessary condition for applying this theory is that the sharp edge model kept at an angle of attack in a supersonic stream must have the shock wave attached to the airfoil at the leading edge. The theory employs a step-by-step application of shock relations for shock waves and expansion wave

3.2 Hydrodynamic Forces on Airfoils

When an airfoil is placed in a high velocity shallow water stream of fixed Froude number with free surface, the resulting water depth distribution around the model will be analogous to pressure distribution around the same model in supersonic gas stream Mach number equal to Froude number. The conditions under which the analogy will be valid are :

- The hydraulic model should be geometrically similar to the airfoil.
- The Froude number of the water should be equal to the Mach number of the gas stream.
- The chord plane of the model must be perpendicular to the water channel bottom floor.

From hypercritical water flow, through the measurement of depth distribution, the hydrostatic pressure around the model can be determined. Hence, the hydrodynamic lift, drag, and moment acting on the model can be computed. [4]

3.3 Measurements with a Wedge Airfoil

A wedge airfoil of wedge angle 7^0 has been used for the study. The experiments were conducted in a water flow channel with a test section cross section of 600mm x 1800mm.

Figure 3.1: Wedge shape Airfoil



The wedge airfoil was tested at a Froude number of 2.13 at various angles of attack. Using the properties of analogy, it can be shown that [6].

 $C_{L}=2(h_{4}/h_{1})^{2}\cos \alpha -(h_{2}/h_{1})^{2} \left\{ \cos (\phi -\alpha)/\cos \phi \right\} -(h_{3}/h_{1})^{2} \left\{ \cos (\phi +\alpha)/\cos \phi \right\}$

 $C_{D} = 2(h_{4}/h_{1})^{2} \sin\alpha - (h_{2}/h_{1})^{2} \left\{ \sin(\phi - \alpha)/\cos\phi \right\} - (h_{3}/h_{1})^{2} \left\{ \sin(\phi + \alpha)/\cos\phi \right\}$

 $C_{\rm M} = (h_4/h_1)^2 - (h_2/h_1)^2 (1/4\cos^2\phi) - (h_3/h_1)^2 (3/4\cos\phi)$

Where C_L , C_D , and C_M are the lift drag, and pitching moments respectively, and α and φ , respectively, are the angle of attack and the semi wedge angle. The subscripts 1, 2, 3, and 4 to the water stream depth h refer to the zones of the flow field around the airfoil. By measuring the depths at various locations around the wedge the hydrodynamic coefficients can be obtained.

4 Hydraulic Jumps (Shocks)

It is known that in shooting water under certain conditions, the velocity may decrease over short distances and the water depth suddenly increase. An unsteady motion of this type is known as a hydraulic jump. Hydraulic jumps occur only in shooting water that is in water streams with flow velocity greater than the wave propagation velocity. To have a better understanding of hydraulic jump, let us examine flow through a large sluice gate. Let us imagine the forward water ahead of the gate at rest, and that from behind there arrives the front of a water wave which arose from the opening of the sluice gate. If the wave were very small, it would move forward with the basic velocity $(gh_1)^{1/2}$. Because, however it has finite height (h_2-h_1) , it moves to a first approximation, with the velocity

$$U_1 = [g(h_1+h_2)/2]^{1/2} (h_2/h_1)^{1/2}$$
 4(a)

Where subscripts 1 and 2 refer to states ahead of and behind the water wave. This velocity is much larger than $(gh_1)^{1/2}$ or $(gh_2)^{1/2}$. In the present flow system, water may be considered as moving with the velocity u_1 , with respect to the wave. That is, apparently the wave is made to remain at rest in peace. Now the water ahead of the wave flows with velocity u_1 , and it is greater than $(gh_1)^{1/2}$. From the above argument it can be seen that the hydraulic jump will remain stationary only in shooting water, such as flow through sluice gate. If the wave existed in streaming water, it would, because of its propagation velocity (which in this case is larger than flow velocity), travel upstream. There would be the usual outflow from upper to lower level without shock. Hence the term shock will be used interchangeably with hydraulic jump, and naturally has nothing to do with the compressibility of water. A shock in which the wave front is normal to the flow direction is called a right hydraulic jump. It naturally has the property that the propagation velocity of the shock wave relative to the wave water is equal and opposite to the water velocity ahead of the jump. [3]

The second kind of hydraulic jump that is along a line oblique to the flow direction is called a slant hydraulic jump. Let the water flow from left to right out of an open sluice gate. The water depth decreases and the velocity increases. The water flows from constant upper water level into a basin with constant lower water level. Because the difference in head is greater than one third of the upper water depth, the water after escaping from the sluice receives a larger velocity than the basic wave propagation velocity, so that it shoots. It is thus possible for it to accelerate so rapidly that water surface of the flow becomes lower than the lower water level. There is a portion of the flow for which there is considerable pressure rise over a short distance. In this flow, however, the jump does not take place normal to the velocity, but instead along a line oblique to the flow direction, and we have a slant jump.

The slant jump, like the right jump, occurs only in shooting water. In order to be able to give a simple numerical treatment of the slant hydraulic jump, we make the assumption that the motion is entirely unsteady; that is, the water jumps suddenly along the jump line from the lower water level to the level after the jump. The simplest case of such a jump is obtained if a parallel flow is deflected by an angle. A similar deflection in a gas stream with supersonic velocity will result in oblique shock. Flows past such shocks are treated exhaustively in standard books on gas dynamics. Here, however, for the shock of the shooting water, the analogy with a compressible gas flows for γ =2 no longer strictly hold.

At this stage, it is essential to realize that most of the assumptions made in the establishment of the analogy cease to be valid at discontinuous flows such as a hydraulic jump. For instance, the assumption that the vertical velocity or acceleration is negligibly small is not valid when a jump occurs, because the vertical component of velocity is considerable at jump. But equations of continuity and momentum are still valid on either side of discontinuity.

In supercritical shallow water flows, discontinuities in the form of a hydraulic jump occur. This discontinuity results in a sudden rise in water depth and change in velocity. At a jump, surface tension and viscosity act as equilibrating forces and they cannot be neglected. Energy losses occur at the hydraulic jumps that are dissipated through heat losses in turbulence at strong jumps and undulations in weak jumps.

In supersonic gas flows shocks occur. In such a case viscosity and heat conduction are no longer negligible. A part of them are used to increase the internal energy and the rest affect the dynamic properties of the flow. The flow process across a shock is no longer isentropic but adiabatic.

Because the flows with shocks or jumps are not isentropic it might be expected that the quantitative correspondence breaks down. The entropy change, of course, is third order in shock strength so that the weak shocks are nearly isentropic. In weak jumps also energy losses are negligibly small. These changes are of the third order with reference to the difference in depth across the jump. Hence, the quantitative similarity should hold good to a good approximation.

4.1 General Equations for Attached Oblique Shocks

For an oblique shock with shock angle β and the flow turning angle θ in a gas flow the Mach number after the shock may be expressed as (Rathakrishnan, 1995)

$$(\mathbf{M}_2)^2 \sin^2(\beta - \theta) = \{1 + (\gamma - 1/2)(\mathbf{M}_1)^2 \sin^2 \beta\} / \gamma(\mathbf{M}_1)^2 \sin^2 \beta - (\gamma - 1/2)$$

Where the subscripts 1 and 2, respectively, refer to conditions upstream and downstream of the shock wave. The density ratio across the shock can be expressed as

$$\rho_2 / \rho_1 = (\gamma + 1)(M_1)^2 \sin^2 \beta / (\gamma - 1)(M_1)^2 \sin \beta + 2$$

$$\rho_2 / \rho_1 = \tan\beta / \tan(\beta - \theta)$$

The pressure ratio across the shock is given by

$$P_2/P_1 = 1 + (2\gamma/\gamma + 1)\{(M_1)^2 \sin^2 \beta - 1\}$$

In terms of density ratio, the pressure ratio may also be expressed as

$$P_2/P_1 = \{(\gamma+1/\gamma-1)(\rho_2/\rho_1)-1\}/\{(\gamma+1/\gamma-1)-\rho_2/\rho_1\}$$

4.2 General Equations for Slant (Oblique) Attached Hydraulic Jumps

Let subscripts 1 and 2 refer to conditions immediately before and after the jump, respectively. The Froude numbers $Fr_1(= V_{1w} / (gh_1)^{1/2})$ and $Fr_2 (= V_{2w} / (gh_2)^{1/2})$ and the depth ratio h_2/h_1 can be expressed as[8]

$$\operatorname{Fr}_{1} = (1/8\sin^{2}\beta) [\{\tan\theta(1-2\tan^{2}\beta)-3\tan\beta/\tan\theta-\tan\beta\}-1]^{1/2}]$$

$$Fr_2 = (h_1/h_2)\{1/\sin(\beta - \theta)\}[1/2(1+h_2/h_1)]$$

 $h_2/h_2 = \frac{1}{2} \left[\left\{ 1 + 8(Fr_1)^2 \sin^2 \beta_w \right\}^{1/2} - 1 \right]$

From the general equations for the oblique shocks in gas and jump in liquid, it is obvious that these equations are similar.

5 EXPERIMENTAL SETUP

Applications are illustrated here by

- Studying the flow fields around hypercritical airfoils in the steady and unsteady transonic regimes of flow
- Examining the validity of the analogy to steady supersonic flows

For investigating the flow field around airfoils qualitatively, two profiles, namely a shock less lifting airfoil and a quasi-elliptical non-lifting airfoil section have been chosen. The experiments have been carried out in a water bed facility.

For demonstrating the quantitative aspects of the analogy, experiments have been done with a semi-wedge airfoil. From depth measurements the aerodynamic coefficients are computed and the flow pattern around the model wave recorded photographically. For this a water flow channel facility was used.

6 EXPERIMENTAL STUDY

The experimental runs were made with a free stream Froude number keeping semi-wedge airfoil model at 0 angle of attack. The Froude number required was obtained by a trial-and-error procedure, as follows. The volume flow in the water channel was measured by flow meter for the positions of bay-pass valve opening. Knowing the free steam depth and channel width, the Froude number was calculated.

The position of the valve opening was adjusted to different flows. The adjustments were continued until an opening of the valve resulted in the free steam Froude number (1.5 to 2). The corresponding depth of the steam was 5 mm. This satisfied the condition for the flow to be a shallow water stream and thus was an appropriate depth for the analogy.

The flow field around the semi-wedge model was photographed with a flash camera. Using the measured values of the water stream depths, h_1 , h_2 , h_3 and h_4 in the

different zones of the flows fields around the airfoil. The different downstream properties such as pressure, density lift, drag and moment coefficients were calculated.

The imaged is processed in the Roborealms for line generation around the body. The processed imaged is processed in Matlab for generating the angle around the body. The calculated value was put in C to calculate the properties.

The validity of the analogy can be checked by calculating the C_{l} , C_{d} , C_{m} independently, with shock-expansion theory and comparing the results of the two methods.

It is seen from these flow patterns that all the features of a supersonic gas floe field around a semi-wedge ,such as the shock waves and the expansion fan, which can be viewed by employing a Schlieren technique in the gas flow, are exhibited by the analogous water flow field. There are many waves caused by the disturbance. Therefore, the experiment has to exercise extreme caution in interpreting the analogous gas dynamics around the model.

7 RESULT AND DISCUSSION

The experimental value of wave angle found is in total correlative with the value found with help of formula of gas dynamics. Hence analogy stands with minor errors.

7a: RoboRealms Results

Figure 7.1: Semi-wedge 7 angle at F=2



a) Real Image

b) Robo Realms Image



7b. Matlab Results

Table 7.1: Matlab Results of Semi- Wedge Angle @ 7

| MODEL HALF ANGLE | $\mathbf{F} = 2$ | F = 1.5 |
|---------------------|------------------|---------|
| Wave Angle | 39.2332 | 34.5647 |

7c. Turbo C Results

 Table 7.2: C Results of Semi- Wedge Angle @ 7

| Coefficients | Angle at 7 at F | Angle at 7 at F | |
|--------------|-----------------|-----------------|--|
| | = 1.5 | = 2 | |
| CL | -0.011 | -0.024 | |
| CD | 0.013 | 0.017 | |
| См | 0.012 | 0.021 | |

8 Error Conclusion

- The wave angle formation was disturbed by interference effect of transverse wave.
- The direction of jet is not constant through the water bed.
- The height of water is measured through rough idea by means of scale.
- Due to the voltage fluctuation there is no continues flow of water over the bed.
- Pressure of flow is not uniform over the bed.
- The leveling of ground is not proper.
- Intensity of light is not proper over the water bed.
- The angle of attack is not properly zero in some cases.
- The changes in the vertical motion of water level over bed introduce some error in height

measurement around the model. Since this analogy holds for non-viscous fluids, it was found that there was some viscosity due to dust and addition of kerosene.

9 Conclusion

The analogy between the shallow water flow with the free surface and two-dimensional gas flow has been found to be useful for qualitative as well as quantitative study of high speed flows.

The essential feature of this analogy is that the Froude Number of shallow water flow with the free surface is equivalent to a gas stream with Mach number equal to that respective Froude number. This analogy is valid for one-dimensional and two-dimensional flow only. The water analogy is well established for a gas with $\gamma = 2$.

It was seen from the results that all the features of supersonic gas flow field around the semi-wedge such as shock wave and expansion fan can be exhibited by the analogous water flow field.

It is important to realize here that, in addition to expansion and compression waves, there are many waves caused by the disturbance which affects the formation of wave angle around the body.

Therefore, experimenter has to excises extreme caution in interpreting the analogous gas dynamics waves around the model.

The error arises as a consequence of neglecting vertical acceleration are proportional to h and $(Fr)^2$. The error also depends on the shape and size of the model and varies inversely with the size of the model. Hence the theoretical and practical phenomenon matches and desirable results achieved.

10 Reference

[1] E. PREISVCERK, "Application of the Methods of Gas Dynamics to Water Flows With Free *Surfaces*," NACA TM-934 and 935.

[2] Arnold Sommerfeld: Partial Differential Equations in Physics; Lectures on Theoretical Physics, Vol. VI, 1964

[3]http://www.ifm.uni-hamburg.de/ers-

sar/Sdata/oceanic/shipwakes/intro; Ship Wakes

[4] Munson, B.,R.; Young, D., F.; H. Okiishi, T., H.: *"Fundamental of Fluid Mechanics"*. Third edition. Oiwa State, University, Ames, Iowa, USA, 1998.

[5] Kelvin Waves Theory, Creeper

[6] Lavicka, D.: *The study analogy between compressible flow and free surface flow (experiments and diagnostics)*, EGIM laboratory, Marseille, France, 2005.

[7] J.J. Stoker: Water Waves; The Mathematical Theory with Applications; 1992

[8] Rani, S. L., and Wooldridge, M. S., "Quantitative Flow Visualization Using the Hydraulic Analogy," Experiments in Fluids, Vol. 28, No. 2,2000, pp. 165–169.

Onion: a Visual Formal Method for Workflow Design in Cloud Computing

Jinho On, Sujeong Woo, Moonkun Lee Department of Computer Engineering Chonbuk National University Jeonju, Jeonbuk, Republic of Korea, 561-756 {jjinghott, wpig04}@gmail.com, moonkun@jbnu.ac.kr

Abstract—This paper presents a new approach for workflow design in cloud computing.

Generally the workflow design in cloud computing is specified by BPEL/BPMN, which is transformed into the existing formal methods for analysis and verification of the design. However the main paradigms of the existing methods reveal some limitations to the design due to their structural characteristics: process algebras mainly focus on in-the-large (ITL) views, and, reversely, state machines mainly focus on in-the-small (ITS) views. Therefore it is necessary to transform ITL views to ITS views, or vice versa, based on some equivalence relations. How can we avoid this notion of the equivalence and transformations to achieve some maximum integration views of the design?

A new visual formal method, called *Onion*, is presented in this paper to integrate these two different ITL and ITS views in a single view. In Onion, processes and their transitive actions are graphically represented in one single entity, just like those of a real onion. Further the temporal properties of actions in the processes are specified in a geotemporal space. Once these are done, the requirements for the design are graphically specified on the processes and actions using a visual logic. Finally, the design is analyzed and verified through simulation in order to see whether it satisfies the requirements and restriction.

The comparative study shows that the Onion approach is very effective and efficient for BPEL/BPMN and overcome some limitations of BPMN/BPMN.

keywords: Cloud Computing, BPEL, BPMN, Onion, Geotemporal Space, Visual Logic, Analysis, Verification.

I. INTRODUCTION

This paper presents a new method for workflow design in cloud computing.

Generally the workflow design in cloud computing is specified BPEL[1] and BPMN[2], which is the visual representation of BPEL. Once the design is completed, it is analyzed and verified by the existing formal methods[3][4][5].

However the existing methods reveal some limitations to the cloud computing due to their structural characteristics, that is, of process algebras and state machines[6][7].

The cloud computing requires strong visualization in process modeling and workflow design, especially, that of process mobility, interactions, reconfiguration, requirements, verification, etc. in the target geo-temporal space[8].

However, the methods in process algebras seem to fail to show the detailed views of the designs since they focus mainly on in-the-large (ITL) views. Reversely, the methods in state machines seem to fail to show the abstract views of the designs since they mainly focus on in-the-small (ITS) views.

Due to these reasons, the equivalence between process algebra and state machine has been a critical issue to show both ITL and ITS views in one design.

The main objective of the research in this paper is this: How to overcome these limitations, satisfying all the requirements of the workflow design in cloud computing?

The method presented in this paper to achieve this objective is the integration of these two different representations in a new visual language, namely, *Onion* [9], with the basic properties of processes, timed actions (communications and movements), as well as hierarchical structure, in a geo-temporal space.

The main distinctive characteristic of Onion is the representation of processes, just as in process algebra, and their transitive actions, just as in state machine, in one singular entity. That is, each process in the design is represented as a process node in Onion, and the interactions and movements of the process are represented as the layered circular leaves of the node, just like those of a real onion, in a geographical space.

This representation shows both ITL and ITS views of the design, and satisfies the requirements of the design, that is, the visualization of process mobility, interactions, reconfiguration, etc.

Once it is done, the temporal properties of interactions and movements of the processes are specified on a geotemporal space, which is the expansion of the geographical space to the temporal dimension. Each timed action is specified in a block with the temporal attributes, such as, *ready time, timeout, execution time, deadline,* etc.

Once the temporal properties are specified, the requirements of the processes and actions are graphically specified with a visual logic on geo-temporal blocks of processes and their actions.

Finally the static and dynamic verification of the requirements are performed and displayed on the geo-temporal space.



Figure 1. Cloud Computing / Workflow Design

The overall approach in the paper is shown in Fig. 1:

- Firstly, an initial workflow design for a target service from cloud computing is specified in BPEL/BPMN and translated into processes and actions in Onion on a geographical space.
- Secondly, the extended features, i.e., mobility, interactions, and control of the processes in the design are added to the translated design on the geographical space.
- Thirdly, the temporal properties of the actions, i.e., communications and movements, are specified on a geo-temporal space in Onion.
- Fourthly, the requirements of the processes and their actions are specified visually on the geo-temporal space with a visual logic.
- Finally, the requirements are visually verified, both statically and dynamically, on the geo-temporal space.

All these processes are performed visually with the support of the Onion System [9].

The paper is organized as follows. Sections II and III present the basic theory of Onion. Section IV presents an iCloud example. In Section V, the approach will be comparatively analyzed with other approaches. Finally conclusions will be made.

II. ONION LANGUAGE

A. Onion Textual Language (OTL)

OTL is a universal process algebra in text for Onion Visual Language (OVL). The basic syntax and semantics of OTL are collected from CSP[10], CCS[11], π -calculus[12] and MA[13]. The extended features are the notions of mobility and process control.

The syntax of OTL is defined as follows:

$$P ::= \alpha' \cdot P'$$

$$| P' | Q' | P' + Q' | P' \cdot Q'$$

$$| P' \wedge (n) | P' \wedge (b)$$

$$| P \setminus \{V\} | P / \{V\}$$

$$| P \triangleright Q ;$$

$$\alpha' := \alpha | (b \rightarrow \alpha) ;$$

$$P' := P | (b \rightarrow P) ;$$

$$\alpha := 0 | \varepsilon | \sigma | \delta | \theta | \tau | \lambda ;$$

$$\sigma := c!x | c?x | c!?x | c?!x ;$$

$$\delta := inQ | outQ | getQ | putQ ;$$

$$\theta := exits | exita | killsQ | killaQ ;$$

Note that b, c and V are a Boolean condition, a channel name, and a list of variables to hide, respectively.

- (1) P: A process in a sequence of actions.
- 2 α': A conditional action, which is an action α with b.
 (i) 0: No action.
 - (ii) \mathcal{E} : An empty action. If $\neg b$, then $(b \rightarrow \alpha) = \mathcal{E}$.
 - (iii) σ : Communication actions.
 - (a) c!x, c?x: Asynchronous send/receive.
 - (b) c!!x, c!!x: Synchronous send/receive.
 - (iv) Movement actions (δ):
 - (a) inQ, outQ: Active or autonomous-in/out.
 - (b) get Q, put Q: Passive or heteronomous-in/out.
 - (v) Process termination actions (θ):
 - (a) exits, exita : self-termination.
 - (b) kills Q, killa Q: Synch/asynch kill.
 - (vi) Thread control actions:
 - (a) P' | Q': Parallel.
 - (b) P' + Q': Nondeterministic choice.
 - (c) $P' \bullet Q'$: Sequential.
 - (d) $P'^{(n)}$: Recursion. $P'^{(n)} = P' \cdot (P'^{(n-1)})$, $n \ge 1$. $P'^{(1)} = P'$, $P'^{(0)} = 0$.
 - (vii) Other actions:
 - (a) $P \triangleright Q$: Exception handling action:.
 - (b) $P \setminus \{V\}$, $P \setminus \{V\}$: hide and reveal.

Note that the timing properties of an action is represented by $\alpha_{[r,to,e,d,p]}$, where *r*, *to*, *e*, *d* and *p* represent *ready time*, *timeout*, *execution time*, *deadline*, and, optionally, *period*, of the action, respectively.

B. Onion Visual Language (OVL)

OVL is a visual representation for OTL.

In OVL, a process is represented by a node, and its actions by a sequence of nested round layers, just like those of a real onion. For example, $P := \alpha \cdot \beta \cdot \gamma$, that is, a process (*P*) consisting of three sequential actions (α , β , and γ), is represented as shown in Fig. 2.



Figure 2. Active Process with Actions

Note that actions are denoted in two different types of circles: straight and dotted. These imply *active* (α) and *inactive* (β , γ) actions. A process with an outermost active action is defined as an *active process*, and a process with the outermost inactive as an *inactive*. Note that all the conditional and nondeterministic actions are inactive.

For communication, a port is defined and represented by a dot on a circular action leaf. A communication between two processes is represented by a straight arrow between two ports. The synchronicity/asynchronicity of communication is distinguished by *open/close* mark on the dot. Note that only the arrow remains dotted unless both actions of two processes for a communication are active. Fig. 3 shows $\alpha ::= x!?m$, $\beta ::= x?!m$ and $\gamma ::= x?m$.



Figure 3. Process with Ports.

An *autonomous* movement is represented by a curved arrow with closed round/arrow marks at the ends, and a *heteronomous* movement by a curved arrow with open round/arrow marks. Once a movement occurs, the process will be relocated in another place, called *future process*, as the target of the movement. For example, Fig. 4 shows both autonomous and heteronomous movement actions. Note that the targets of the arrows are future processes.



Figure 4. Movement Actions with Future Processes.

A process can be nested in another process by visual containment.

C. Transformation of BPEL Key Statements Onion

Table 1 shows the relations between BPEL key statements to Onion, extended from those to LOTOS in [14]. This allows the basic BPEL specification of workflow design to be translated to Onion. Once it is done, the further specification of mobility and reconfiguration of processes in the design is possible in Onion.

| | Sample BPEL | Sample LOTOS | OTL | OVL | |
|---|---|--|------------------------|---------------|--|
| 1 | < act1 > <assign> <copy> <from expression="5"></from> <to var="x"></to> <copy> </copy></copy></assign> < act2 > | act1; exit(5) □ accept x:Nat inact2 | | | |
| 2 | <receive variable="m"> </receive> | g?m:Nat; | g?m | | |
| 3 | <reply variable="m"> </reply> | g!m:Nat; | g!m | (g) g) | |
| 4 | <sequence> < act1> < act2> </sequence> | act1;act2 | act1 □ act2 | act act | |
| 5 | <while condition="bpws:getVariableData(x)>=0"> <act1></act1> </while | | P^(x>=0); P::=act1; | P ((x)=0) | |

TABLE I. TRANSLATION BPEL/MPMN TO OTL/OVL VIA LOTOS

| 6 | <invoke <br="" invar="mS">outvar="mR"> </invoke> | gSImS:Nat; gR?mR:Nat; | g!mS□ g?mR | A C C C C C C C C C C C C C C C C C C C |
|---|---|---|--|--|
| 7 | <pre><pick> <onmessage variable="m1"> < act1> <onmessage> <onmessage> <onmessage> < act2> </onmessage> </onmessage></onmessage></onmessage></pick></pre> | (g1?m1:Nat;act1) [] (g2?m2:Nat;act2) | g1?m1+g2?m2 | 9 (91) (91) (92) (91) (92) (91) (92) (92) (92) (92) (92) (92) (92) (92 |
| 8 | <flow> <act1> <source <br="" linkname="link1"/>condition="cond1"/> </act1> <act2> <target linkname="link1"></target> </act2> </flow> | act1; ([cond1]->link1 !1; [] [not(cond1)] ->link1 !0;) (link1 ?x:Bool; ([x=1]->act2 [] [x=0]->i)) | P Q P::=(cond1->link1) Q::=link1 | P conditionint Q link1 |
| 9 | <switch> <case condition="bpws:getVariableData(x)>=0"> <.act1.><!--.act1.--> <otherwise> <.act2.><!--.act2.--> </otherwise> </case </switch> | | (x>=0->act1) □ (x<0->act2) | P (b=0-sat1) · (x<0-sat2) |

III. GEO-TEMPORAL SPACE AND VISUAL LOGIC

A. Onion Geo-Temporal Space (GTS)

The *Geo-Temporal Space* (GTS) for Onion is the representation of geographical space over temporal expansion in Onion. *Geographical space* (GS) is the one-dimensional vertical representation of OVL. *Temporal Space* (TS) is the one-dimensional horizontal representation of execution of processes and actions in OVL. Consequently, GTS is the two dimensional representation of OVL.

A process P in GS is represented by a vertical line between two virtual geographical points, i.e., t (at the top) and b (at the bottom), and is denoted by P(t,b). Any nested process should reside in the space of its parent process in the same pattern. Examples are shown in Fig 5. By default, P(t,b) is represented by P.



The TS over GS for a process *P* is defined as a *geotemporal block* (GTB), namely, *process GTS*, and represented by a horizontal space over the GS. There are two types of GTB: *discrete* and *continuous*. The discrete space is GTB at a discrete time, t_1 , i.e., $P[t_1]$. The continuous space is GTB at the continuous time between t_1 and t_2 , i.e., $P[t_1, t_2]$, as shown in Fig. 6.

A timed action of a process *P* is also represented by a GTB namely, *action GTS*. Its timing requirements are based on the definition of the timed action: *r*, *to*, *e*, *d* and *p*. The example of the actions for $P := \alpha \cdot \beta \cdot \gamma$ is shown in Fig. 7. The GTBs of the actions are represented with dotted lines since they are defined to be inactive. Note that the types of actions are denoted by the squared yellow marks in legend.



Figure 7. Action with Timing Requirements.

A timed interaction between P and Q is represented by a direct edge between action GTBs in process GTBs, as shown in Fig. 8.



Figure 8. Timed Interaction (Synch Communication).

B. Visual Logic (VL) on Geo-Temporal Space

Visual Logic (VL) is a language to specify visually the requirements of processes and actions in GTS. VL defines the requirements as some inclusion or precedence relations

between/among process and action GTBs in GTS. The syntax and semantics of VL are described in Table II. Note that A and B in the table are process GTBs in geographical requirements, or action GTBs in temporal requirements.

TABLE II. SYNTAX OF VISUAL LOGIC



Fig. 9 shows some examples in VL on a GTS. For example, R_1 to R_6 define the following requirements:

- R_1 : Process A must reside in Process B all the time.
- R_2 : Action b_1 in Process *B* must precede Action c_1 in Process *C*.
- R_3 : Actions in TS_2 of C and actions in TS_3 of D can execute concurrently.
- R_4 : The total execution time of actions in TS_4 of D should be less than 100 time units.
- R_5 : C cannot move into A within the TS_1 interval of A.
- R_{6} : All the actions of *D* have to be terminated in 600 time units.

Note that requirements can be expressed recursively in conjunctive and disjunctive forms.



IV. AN EMS EXAMPLE

A. 911 Service in EMS

911 service is the main service in *Emergency Medical* Service (EMS), where, in case of an medically urgent situation occurred to *Patient* in *House*, *EMS Center* is informed and sends a message to 911 to transport the patient to *Hospital*, and, upon the message, an ambulance from the 911 go to the house and take the patient to the hospital in time.



Figure 10. Workflow model example in BPMN

This service can be described in BPMN as shown in Fig. 11, where the movements of ambulances and patients are represented in the form of message-passing due to the limitation of expressivity for movements in BPMN. Consequently there are considerable difficulties in specifying the movements of patients, ambulances and doctors in the EMS example. Further there is no way to display the overall reconfigurations of the example in timely progression.

B. OTL and OVL

In the Onion method, a BPMN specification for the EMS can be translated into OTL, and the movements and actions of the processes in the EMS are specified with related data further in OTL and OVL as shown in Fig. 11 and 12.



Figure 11. OVL for iCould Example



Figure 12. OTL for EMS Example

In OVL, the key players of the EMS can be visually configured as shown in Fig. 11, and their interactions and movements are visually recognizable.

- The pictorial description is as follows:
- $a_1 \sim a_6$: A critical event is detected from a patient's sensor and sent to the EMS center. The center informs both a 911 center and a hospital, where the patient has been under treatment, of the patient's situation.
- $a_7 \sim a_8$: An ambulance from the 911 goes to the patient's house.
- *a*₉: The patient is carried out of the house and gotten on the ambulance.
- $a_{10} \sim a_{11}$: The ambulance goes to the hospital.
- $a_{12} \sim a_{13}$: The patient is gotten off the ambulance and carried into the emergency room of the hospital.
- a₁₄~a₁₅: A medical doctor in the hospital is informed of the patient with a message.
- $a_{16} \sim a_{17}$: The doctor goes to the emergency room.

Note that, in OVL, the movements, communications, and reconfiguration of processes can be represented in hierarchically organized structure.

As stated, the most important feature of OVL is the integrated representation of ITL and ITS views of the service.

C. Onion GTS

Once the EMS is specified in OVL, it is necessary to define temporal properties of actions of processes in the EMS. Fig. 13 shows GTBs of each processes and actions in them, as well as the interactions, from a_1 to a_{17} , among them. Note that the temporal properties of each action are scaled over its temporal dimension. Its temporal scale is presented at the top of the figure.

D. VL

The requirements for the example are specified visually in Fig. 14. These imply the following conditions to be satisfied:

- *R*₁: The doctor has to arrive at the emergency room before the patient is carried into the room.
- *R*₂: The patient has to be carried to the hospital in the time units of 400 after he is gotten on the ambulance.

• *R*₃: The sensor has to be attached to the patient's body all the time.



Figure 13. GTS and GTBs for iCould Example



Figure 14. VG Requirements for iCould Example

E. Analysis and Verification

In the Onion method, the above requirements can be verified, dynamically. The dynamic verification implies the verification based on run-time analysis. The results of the verifications can be visually displayed on the verified GTS.

For example, Fig. 15 shows the results of the dynamic verification for the example. The similar results are generated. The difference is that it shows the run-time results as it is being executed. The figure shows the results of the execution at the time units of 700, as shown at the top of the figure. All the executed GTBs are represented as the straight lines.



Figure 15. Results of Dynamic Verification for iCloud Example

V. COMPARISON WITH OTHER METHODS

Generally workflow design in cloud computing with BPEL/BPMN consists of the following steps: 1) specification, 2) analysis, and 3) verification. Once the design is specified in BPMN, the specified design is transformed to some formal methods for analysis and verification. The properties of the design to be analyzed and verified are basically dependent on the capability and characteristics of the formal methods. As stated, process algebras mainly focus on ITL views, and, reversely, state machines mainly focus on in-the-small (ITS) views. Due to these characteristics, it is necessary to transform processes algebras to state machines, or vice versa, to see both ITL and ITS views by means of some equivalence relations.

Onion attempted to handle this problem by integrating both views in one view in Onion. Further Onion tried to visualize all the contents of specification, requirements, analysis and verification of the design, as well as all the temporal properties of the contents, using OVL, GS, GTS, GTB, Visual Logic, etc.

VI. CONCLUSIONS AND FUTURE RESEARCH

This paper presented the Onion method for workflow design in cloud computing. The method has the following innovative features:

 Integration and visualization of characteristics of process algebras and state machines,

- Visual specification of processes and their actions/interactions on a geographical space,
- Notion of GTBs for processes and their actions on GTS, as well as interactions among processes,
- Notion of visual logic and graphical requirements on GTBs in GTS, and
- Notion of visual analysis and verification.

The paper also demonstrated that the method would be well suited for the design by translating BPMN into Onion and going through analysis and verification. It also showed that some limitations of BPEL/BPMN could be overcome by using Onion.

The future research includes the development of the complete set of the Onion tools, its application to the real industrial examples, etc.

ACKNOWLEDGMENT

This research was supported by Basic Science Research Program through the National Research Foundation of Korea (NRF) funded by the Ministry of Education, Science and Technology (2010-0023787), and by research funds of Chonbuk National University (2011) and by second stage of Brain Korea 21 Project in 2012.

REFERENCES

- P. Sarang, M. Juric, and B. Mathew, Business Process Execution Language for Web Services BPEL and BPEL4WS 2nd Edition, 2nd ed. Packt Publishing, 2006, p. 372.
- [2] A. Grosskopf, G. Decker, and M. Weske, The Process: Business Process Modeling using BPMN, 1st ed. Meghan-Kiffer Press, 2009, p. 182.
- [3] W. Li, Z. Yang, P. Zhang, and Z. Wang, "Model Checking WS-BPEL with Universal Modal Sequence Diagrams," in Computer and Information Science (ICIS), 2011 IEEE/ACIS 10th International Conference on, 2011, pp. 328–333.
- [4] H. Cao, S. Ying, and D. Du, "Towards Model-based Verification of BPEL with Model Checking," in Computer and Information Technology, 2006. CIT '06. The Sixth IEEE International Conference on, 2006, p. 190.
- [5] F. Van Breugel, "Models and Verification of BPEL," Monograph on Testing and Analysis of Web ..., 2006.
- [6] L. Baresi, D. Bianculli, C. Ghezzi, S. Guinea, and P. Spoletini, "Validation of web service compositions," Software, IET, vol. 1, no. 6, pp. 219–232, 2007.
- [7] M. Ter Beek and A. Bucchiarone, "Formal methods for service composition," Annals of Mathematics, 2007.
- [8] P. Wong and J. Gibbons, "Formalisations and applications of BPMN," Science of Computer Programming, 2011.
- [9] J. On, "Onion: A Graphical Language for Process Algebra," in Computer Software and Applications Conference (COMPSAC), 2011 IEEE 35th Annual, 2011, pp. 708–711.
- [10] C. Hoare, "Communicating sequential processes," Communications of the ACM, 1978.
- [11] R. Milner, A calculus of communicating systems (Lecture notes in computer science). Springer-Verlag, 1980, p. 171.
- [12] D. Sangiorgi and D. Walker, The Pi-Calculus: A Theory of Mobile Processes. Cambridge University Press, 2003, p. 596.
- [13] L. Cardelli, "Mobile ambients," Theoretical computer science, 2000.
- [14] K. Turner, Using Formal Description Techniques: An Introduction to Estelle, Lotos, and SDL. 1993.

Full digital design and fabrication of building components by laser forming

G. Kiani¹, R. Binti Ibrahim², and C. Co-author²

¹Institute of forestry & Forest products, University of Putra Malaysia, Serdang, Selangor, Malaysia ²Faculty of Design & Architecture, University of Putra Malaysia, Serdang, Selangor, Malaysia

Abstract - The current study was performed to evaluate the criteria in conducting full digital process of freely designed building components. Current construction Industry is not capable of coping properly with its design concepts and stylish forms of digital era. Conspicuous move from previous mold base and highly craft dependent works to a comprehensive precise numeric base seems to be the solution. First of all a thorough literature reviewing was administered to understand the current functional industrial system for the purpose. Secondly, elaborated the selection of appropriate material which can be fitted into determined Digital production system. Thirdly, simulation and modeling were employed in order to illustrate the level of workability for the suggested process. The study is limited to materials which can be sintered or melted, and also solidify quickly. Based on the structural behavior prediction of the building components; stability of the model was judged. The aim in the current study is to facilitate the industrial design and manufacturing by using precise computational control.

Keywords: Digital process, computational design, numeric control production, Construction Industry

1 Introduction

In current system of mass production and construction industry the dominant approach is production as panel format [1] by means of casting. Consequently, our design and production will be controlled by molds limitations.

Application of RM in architecture yet to be limited in modeling and prototyping [2] but the other industries like aerospace and automobile [3,4] borrows this system to produce actual and final production which benefits them by higher flexibility, ease of procedure, less waste and lapse besides structural functionality. [5]

In 1997 Penga argued about possibility of employing automotive system for selectively produce sand and cement 3D volume. But his works was not continued since Behrokh Khoshnevis in 2000 [6] contributed a layer base, structural wall production, mostly similar to 3D printing, but named as contour crafting.

The latest application of Rapid Manufacturing in construction industry was adopting the system for producing molds of complex design. [7]

But lately Loughborough University embark on applying concrete to adapted 3D printing system [8] in order to verify the possibility of exercising this machinery in prefab building s Selective laser sintering is [9] a process of slicing the 3D CAD data of your design into layers then a computer-directed heat laser fuses together layers of matrix powder into 3D model, solidified layer and again melted another layer on top. This will generate layers on top of each other in thickness as small as 0.15mm, dependent upon orientation, providing the highest level of functionality combined with speed that is yet to be obtained.

SLS offered the perfect solution for complex design and geometry without any difference compared with simple one. Taken directly from scanned 3D CAD data, elaborated and high- tech details can be easily copied and previously impossible' forms created.

This system, which it seems very optimistic to be applicable in actual production phase is the aim in this study. [10]



Figure 1: Proposed concept of production process for digital fabrication of building components

The current effort tries to elaborate the tremendous advantages of utilization of full digitally fabrication of Rapid Manufacturing over the semi craft based system in-use. This fashion will be followed by defining the required criteria to convert this idea to a fact. Hence, an example; which is miniature model production, is presented for illustrating how the process expected to perform. It will show how this new routine can bring ease and acceleration of producing curvilinear or Non-Euclidean forms in prefab way. Latterly, based on the procurement the workability of process can be evaluated.

The promising process of design and production will be as following diagram (figure 1).

2 FEA Analyzing of the curvilinear

volume

2.1 Design process

In a complete building unit production which is going to be produced in a very nonlinear shape via RM system, as a result the egg shape one story building was resolved and the design details were accomplished subsequently, the 3D format of design in 3D AutoCAD was prepared and fragmented in appropriate segments(Fig 2). The joints were determined to be tongue and groove, as one of the most famous ones for prefabricated volumes [11]. Each ring was fragmented into 4 sections including the roof, walls and floors.

2.2 Simulation

Predicting the structural necessities and determining material composition; which should be compatible with the system and also meet the structural needs, are constitutive of this stage. Expected loads and boundary conditions also needed to be defined. Latterly, these integrated data will be illustrated in a 3D format for the production system while feed the system with the selected material and the production phase will be commenced.

Ansys12.1 was utilized for structural behavior simulation. Ansys is known as reliable software to predict the behavior of very dynamic structures and composite materials while applying the expected loads [12]. For performing this phase, define the following principals seems crucial:

- The volume considered to be fixed to the ground and the joints are tongue and grooves.
- Dead load: Dead Load weight of each panel is calculating from:

 $V^*\rho^*g = W$

Total Dead Loads of component built by ABS= 46135.8 kg while density is: 1010 kg/m3 and gravity of 10

- Earthquake load was calculated [8] and the magnitude was: 7243.32 kg/m2
- In accordance with ASCE, (2003) wind loads was determined.

Just after the forenamed specifying, it is probable to start on structural evaluation and analysis with Ansys. The steps which Ansys had in order to dispense the results included the following items:



Figure 2: Egg Shaped Building Design selected for This Study

2.2.1 Structural behavior simulation phase

The element type of Shell & Linear layer 99 was selected for this simulation [13]. This element consists of eight nodes with six degrees of freedom for each. Linear layer 99 benefits from fewer time for elements with several layers and also it has the capability of attributing specific properties for each different layers.



Figure 3: Stress distribution



Figure 4: Total deformation

The data were input to the software and the problem was solved and the diagram and related results were plotted. As it is shown in the figure 3, the highest stress is implementing over the roof where the area left. The minimum generated stress was 2.8 Pa and the maximum stress was equal to 1.6 MPa while the design allowable for composite was calculated as 25 MPa therefore the generated stress still is in safe domain and considered as tolerable.

Maximum deformation happens over the roof of the volume while the area left unincorporated (figure 4). The predicted amount for this deformation predicted to be 0.001 m which is not critical.

TABLE 3FEA RESULTS OF KCRP.

| Elongation entered | Expected Avg. Stress (MPa) | Found Stress Interval over volume (MPa) | |
|--------------------|-------------------------------|--|--|
| 0.4-0.5 | 15-19 | 1.6-7.9 | |

2.3 Modeling Process Phase utilizing Rapid Manufacturing Technology

The modeling was conducted with prototyping material in order to verify the design and production process without material verification. Therefore, ABS was employed which is regularly used in Rapid manufacturing modeling machine.

2.4 Similar Features of RM Systems

Generally, the prototypes crafted with the existing and growing RP processes greatly resemble each other [14]. The researchers need to electronically section the solid or surface CAD model into layers of prearranged thickness. Shapes of the parts are collectively characterized by means of these sections. Later, the information about each section will be transmitted electronically to the RP machine layer by layer.





Figure 5: (a) the Sliced Model in STL format (b) Successful assembly of components based on

Following steps are essentially employed by the RP process for fabricating the prototypes:

- Creating a CAD model of the design.
- Converting the CAD model to STL (Stereolithography) file format (fig 5a).
- Slicing the STL file into 2-D cross-sectional layers.
- Growing the prototype.
- Post-processing or finishing (fig 5b).

3 Discussion

The whole production phase took 21 hours to produce a model in scale of 1:20 which means to produce a 1:1 model it will take almost 420 hours with the same speed. In this procedure (fig 6) which utilize full digital design and production system, every intricate detail should be designed and predicted prior to production.

The major move in this process is merging material production with finalized product production. In this method the digitally design data is translating directly for the production system, which will decreases a lot of complex stages and material wastage. The main time consuming effort happened before embarking on production phase and its assembly due, which is in contrary with the general prefabricated production process in which the most effort should be made after preparation of design and material determination, during assembly phase. As the joints in this study are designed during the production process, therefore the assembly phase will be easy and fast. It was revealed from the modeling and simulation phases that the whole process can be conducted easily with the least human effort and highest precision if the appropriate material can be adopted for the construction and the scaled up system be available for robust size construction components.



Figure 6: The selective laser sintering of building components

Figure (7) shows that there is a significant process shortening while comparing the full digitally process of this study with known prefabrication design and production.



Figure 7: Comparison between general prefabrication process (right) and suggested process (left)(adapted: Pastor, J. M., et al., (2001) [7]

It should be notified that the main efficiency of the new system is: enabling production of very complex and nonlinear design in the same ease as simple design production, which is in contrary with the previous system that imposes a lot of complex, time consuming and craft made steps.

It can be concluded from the modeling phase that although there is wide range of benefit from direct digitally fabrication of building components to borrow RM system from other industry into construction industry, there is a need to apply some changes into the system to make it compatible with robust size of building components, while increasing the speed.

This study contributed:

• Adopting new futuristic style for prefab construction is not encountering any limit anymore.

• Complexity in parts production and assembly are not an issue any more.

- Tooling for prefab construction will be eliminated.
- Integrating material and design as structural member.

• Predicting the whole design and production detailing prior to production brings us the least waste for material, energy and time.

Further relevant study can be conducted about adoptable SLS system and size for construction industry, also physical and mechanical properties of applicable constructional material to the system.

4 Conclusions

It was demonstrate that by applying appropriate material to the right system and merging structural calculation, material optimization and the design specification, new direct digitally fabrication of complex building components can be founded.

In order to evaluate the production process idea emerged from Selective laser sintering production procedure concept of parts, in this research project a simulation process was defined and conducted to evaluate the variation of the layer by layer processing the parts of the intended structure. The simulation is confined to the chosen design and appropriate material for the research purpose.

ACKNOWLEDGMENT

The authors would like to thank the University of Putra Malaysia Forest Product Laboratory for facilities, and space provided during the course of the study and also to Satavand Shiraz & Chavigan Company which provide the financial supports.

5 References

[1] Cynthia E. Johnson, Rachel Kennedy. (1900-1960). House in a Box:Prefabricated Housing in the Jackson Purchase Cultural Landscape Region

[2] Sriraman, V. D. (2002). Selecting the appropriate rapid prototyping system for an engineering technology program. Journal of Engineering Technology.

[3] Reeves, P. (2007). Rapid Manufacturing – A Business Case For Developing Reusable Multimedia For Engineering Education. International Symposium for Engineering Education (pp. 75-81). Dublin City, Ireland: Econolyst Ltd.

[4] Gideon N. Levya, Ralf Schindela and J.P. Kruthb. (2007). Rapid Manufacturing And Rapid Tooling With Layer Manufacturing (Lm) Technologies, State Of The Art And Future Perspectives . CIRP Annals - Manufacturing Technology, 589-609.

[5] D.L. Bourell, H. M. (1992). Selective laser sintering of metals and ceramics. International Journal of Powder Metallurgy, 369-381C. J. Kaufman, Rocky Mountain Research Lab., Boulder, CO, private communication, May 1995.

[6] Behrokh Khoshnevis, J. Z. (2010). Contour Crafting Process Plan Optimization. Journal of Industrial and Systems Engineering , 33-46.

[7] Kendall, D. (2008). Building the future with FRP composites. Reinforced Plastics , 26-29, 31-33.

[8] Ogden, R. (2008, december). Buildoffsite. Retrieved november 20, 2009, from buildoffsite: www.buildoffsite.com

[9] D.L. Bourell, H. M. (1992). Selective laser sintering of metals and ceramics. International Journal of Powder Metallurgy , 369-381C. J. Kaufman, Rocky Mountain Research Lab., Boulder, CO, private communication, May 1995.

[10] R.A. Buswella,R.C. Soara, A.G.F. Gibbb and A. Thorpe. (2007). Freeform Construction: Mega-scale Rapid Manufacturing for construction. Automation in construction , 224-231

[11] Schelecht, R. o. (1969). Patent No. 3,436,881. USA.

[12] Gengdon, L. G. (2005). Reliability Analysis for Complex Large-scale Structures Based on ANSYS Software. Buildng Structure.

[13] Ansys 12.1 help

[14] Dickenson, N. H. (2003). Analysis of rapid manufacturing- using layer manufacturing processes of production. Mechanical

Optimization Design of the U-shaped Metal Bellows

Ling Yang¹, Mingjin Yang^{1*}, Feng Liu¹, and Guocai Yang²

¹College of Engineering and Technology, Southwest University, Beibei, Chongqing, China ²School of Computer and Information Science, Southwest University, Beibei, Chongqing, China

Abstract - Please Bellows is the flexible element of an expansion joint used to absorb dimensional changes, and the design of a bellows is complex. In this paper, the optimization method was employed in the design of a U-shaped metal bellows. A 4-variable 14-constraint optimization model was built according to related standards and practical engineering experience. Weight factors of mass and flexibility in the objective function were set according to importance of the objectives. Special strategies were adopted to solve the multi-objective optimization problem with mixed discrete and continuous variables. Case study shows that the optimization model presented and program developed can meet requirements of design of the U-shaped metal bellows with high quality and efficiency.

Keywords: optimization design; multi-objective; modeling; programming; metal bellows

1 Introduction

Bellows is the flexible element of an expansion joint used to absorb dimensional changes, such as those caused by thermal expansion or contraction of a pipeline, duct or vessel. The design of a bellows is complex in that it involves the evaluation of pressure capacity, stress due to deflection, fatigue, and instability, etc. In many cases, the design for a particular application involves a compromise of conflicting requirements[1,2]. Optimization techniques are widely used to achieve design goals of the complex design problems with high quality and efficiency[3].

Bellows may be either U-shaped or Ω -shaped (toroidal) in cross section. The U-shaped bellows is superior for great deflection but has a lower pressure capacity for the same material thickness. Conversely, the Ω -shaped bellows is limited to small deflection but has a higher pressure capacity. The use of external reinforcement of the U-shaped can provide a combination of great deflection and high internal pressure capacity, and the pressure capacity can also be increased by the use of multi-ply construction or by increasing the material thickness of the bellows. The Ushaped bellows is mostly manufactured and used in industries, and the Ω -shaped bellows has limited application in case of high pressure and small deflection requirements[4]. The design and optimization process of an unreinforced U-shaped bellows was detailed in this study.

2 Problem description

The structure of unreinforced U-shaped bellows is shown in Fig. 1. The bellows consists of one or more convolutions, and the convolution is the smallest flexible unit of a bellows. $L_{\rm b}$, $L_{\rm c}$, and $L_{\rm t}$ are bellows convoluted length, tangent collar length, and tangent length, respectively. An acceptable design of a bellows must meet requirements of pressure capacity, fatigue life expectancy, instability, spring rate per convolution, etc., and it is complicated by the numerous variables involved such as diameter, pitch, height, number of plies, and material thickness. The basic design equations for the unreinforced bellows are as follows.



Figure 1 Structure of unreinforced U-shaped bellows.

Bellows circumferential membrane stress due to pressure S_2 is:

$$S_2 = \frac{PD_{\rm m}K_{\rm r}q}{2A_{\rm c}} \tag{1}$$

where, *P* is pressure; D_m is the mean diameter of bellows convolutions, and $D_m = D_b + w + nt$; *n* is number of bellows material plies of nominal thickness, *t*; K_r is circumferential stress factor; *q* is convolution pitch; A_c is cross sectional metal area of one convolution.

Bellows meridional membrane stress due to pressure S_3 is:

^{*} Correspondence author
$$S_3 = \frac{Pw}{2nt_p} \tag{2}$$

where, w is convolution height; t_p is bellows material thickness for one ply, corrected for thinning during forming, and $t_p = t \sqrt{D_b / D_m}$.

Bellows meridional bending stress due to pressure S_4 is:

$$S_4 = \frac{P}{2n} \left(\frac{w}{t_p}\right)^2 C_p \tag{3}$$

where, C_p is factor used to relate U-shaped bellows convolution segment behavior to a simple strip beam.

Bellows meridional membrane stress due to deflection S_5 is:

$$S_5 = \frac{E_b t_p^2 e}{2w^3 C_f} \tag{4}$$

where, $E_{\rm b}$ is modulus of elasticity at design temperature; *e* is total equivalent axial movement per convolution; $C_{\rm f}$ is factor used to relate U-shaped bellows convolution segment behavior to a simple strip beam.

Bellows meridional bending stress due to deflection S_6 is:

$$S_6 = \frac{5E_{\rm b}t_{\rm p}e}{3w^2C_{\rm d}} \tag{5}$$

where, $C_{\rm d}$ is factor used to relate U-shaped bellows convolution segment behavior to a simple strip beam. Fatigue life $N_{\rm c}$, number of cycles to failure, is:

$$N_{\rm c} = \left(\frac{c}{S_{\rm t} - b}\right)^a \tag{6}$$

where, a, b, and c are material and manufacturing constants; S_t is pressure, and $S_t=0.7(S_3+S_4)+(S_5+S_6)$.

Limiting internal design pressure based on column instability for single bellows (both ends rigidly supported) P_{sc} is:

$$P_{\rm sc} = \frac{0.34\pi C_{\theta} f_{\rm iu}}{N^2 q} \tag{7}$$

where, C_{θ} is column instability pressure reduction factor based on initial angular rotation; f_{iu} is bellows theoretical initial axial elastic spring rate per convolution; N is number of convolutions in one bellows.

Limiting design pressure based on inplane instability and local plasticity at temperature below the creep range P_{si} is:

$$P_{\rm si} = \frac{1.3A_{\rm c}S_{\rm y}}{K_{\rm r}D_{\rm m}q\sqrt{\alpha}} \tag{8}$$

where, S_y is yield strength at room temperature of the actual bellows material after completion of bellows forming and any applicable heat treatment; α is inplane instability stress interaction factor.

Bellows theoretical axial elastic spring rate per convolution f_{iu} is:

$$f_{\rm iu} = \frac{1.7 D_{\rm m} E_{\rm b} t_{\rm p}^{3} n}{w^{3} C_{\rm f}}$$
(9)

3 Optimization model description

Modeling is of prime importance for optimization of bellows design. The mathematical model usually consists of variables, constraints, and objectives.

For a design, all variables form a design domain, denoted by a vector $X = [x_1, x_2, \dots, x_n]^T$. *n* is number of domain dimensions, and *X* can be considered as a point in the n-dimensional design domain. Variables are often subjected to some constraints, denoted as $g_i(X) \le 0$. *i* is number of the constraints, and $i = 1, 2, \dots, n_c$. Objectives can be expressed by the objective function f(X) based on certain criteria. Then, the optimal design is a group of variables having some values that make f(X) minimum.

3.1 Variables

The variables of an optimization model must be the independent variables and be determined through optimization. Number of variables need to be minimized on the premise of good optimization performance for purpose of the decrease of model complexity.

The structure parameters of an unreinforced U-shaped bellows are as follows: mean diameter of bellows convolutions D_m , inside diameter of bellows convolutions D_b , convolution height w, convolution pitch, mean radius of bellows convolutions r_m , bellows nominal material thickness of one ply t, number of bellows material plies n, and number of convolutions in one bellows N, etc.

Since bellows are compatible with pipeline, duct or vessel linked, inside diameter of bellows convolutions D_b is always defined according to the diameter of pipeline, duct or vessel. There is determined relationship between D_b and D_m , and $D_m = D_b + w + nt$. Thus, D_b and D_m are not independent variables. In most cases, convolution pitch q and mean radius of bellows convolutions r_m are preset by manufacturer for simplicity of the molds for bellows production. Therefore, the independent variables are convolution height w, bellows nominal material thickness of one ply t, number of bellows material plies n, and number of convolutions in one bellows N, then the variable vector is:

$$\boldsymbol{X} = [x_1, x_2, x_3, x_4]^{\mathrm{T}} = [w, t, n, N]^{\mathrm{T}}$$
(10)

3.2 Constraints

The value of optimization objective depends on variables, and there are some constraints acting on these variables. They can be either boundary constraint or behavior constraint.

3.2.1 Boundary constraint

The constraint used to define the scope of variable values or the relations between or among variables.

The ratio w/q has much influence on the performance of a bellows, and it should be confined in the range of 0.6 to 1.6. Then:

$$g_1 = 0.6 - w/q \tag{11}$$

$$g_2 = w/q - 1.6 \tag{12}$$

There are some defined specifications for bellows nominal material thickness of one ply t, and they can be any value of the following: 0.1, 0.2, 0.3, 0.4, 0.5, 0.8, 1.0, 1.2, 1.5, 2.0, 2.5, and 3.0mm. Then:

$$g_3 = 0.1 - t \tag{13}$$

$$g_4 = t - 3 \tag{14}$$

Number of bellows material plies n should be confined in the range of 1 to 5, and number of convolutions in one bellows N should be confined in the range of 1 to 15. Then:

$$g_5 = 1 - n \tag{15}$$

$$g_6 = n - 5 \tag{16}$$

$$g_7 = 1 - N \tag{17}$$

$$g_8 = N - 15$$
 (18)

3.2.2 Behavior constraint

The constraint derived from requirements of design or performance of a bellows and expressed as function of variables.

The stresses should be evaluated for pressure capacity as follows:

$$g_9 = S_2 - C_{\rm wb} S_{\rm ab} \tag{19}$$

$$g_{10} = S_3 + S_4 - C_m S_{ab} \quad \text{(below the creep range)} \tag{20}$$

$$g_{11} = S_3 + S_4 / 1.25 - S_{ab}$$
 (in the creep range) (21)

where, $C_{\rm wb}$ is longitudinal weld joint efficiency factor; $S_{\rm ab}$ is allowable material stress at design temperature; $C_{\rm m}$ is material strength factor at temperature below the creep range.

Fatigue life should be evaluated as:

$$g_{12} = N_{\rm c} / N_{\rm ab} - n_{\rm f} \tag{22}$$

where, N_{ab} is allowable number of cycles to failure; n_f is safety factor of fatigue, and $n_f = 10$.

Instabilities should be evaluated as:

$$g_{13} = P - P_{\rm sc} \tag{23}$$

$$g_{14} = P - P_{\rm si} \tag{24}$$

where, P is design pressure.

3.3 Objective function

Bellows functions as the flexible element of an expansion joint. Total equivalent axial movement per convolution e, an index of flexibility evaluation for a bellows, can be selected as an objective for the optimization, and the bigger flexibility, the better. From Eqs. (4) and (5), there is:

$$e = \frac{(S_{5} + S_{6})w^{2}}{E_{b}t_{p} \left[5/(3C_{d}) + t_{p}/(2wC_{f}) \right]}$$
(25)

Bellows material is high quality austenitic stainless steel with strict requirements. On conditions of performance warranty of bellows, decrease of the material use can result in good benefit for the bellows production. Then, the total mass of a bellows m can be another objective for the optimization, and it is expressed as:

 $m=2\pi\rho NntD_{\rm m} \left(w-2r_{\rm m}+\pi r_{\rm m}\right)+2\pi\rho ntL\left(D_{\rm b}+nt\right)$ (26) where, *L* is belows tangent length including collar, $L = L_{\rm t} + L_{\rm c}$.

By means of evaluation function method, objective function can be expressed as:

$$\min f(X) = \frac{m^{\omega_1}}{e^{\omega_2}}$$
(27)

where, ω_1 and ω_2 are weight factors of mass and flexibility, respectively, and set $\omega_1 = 1$, $\omega_2 = 1.1$ according the importance of objectives.

4 Solving strategies

Three are 4 independent variables and 14 constraints for the optimization problem. Variable of convolution height w is a continuous variable, and variables of bellows nominal material thickness of one ply t, number of bellows material plies n, and number of convolutions in one bellows N are discrete variables. It is a nonlinear optimization problem with mixed discrete and continuous variables, and some strategies are employed to solve the problem[5,6].

4.1 Strategies

1) Initial discrete complex: For the n-dimensional domain, the initial discrete point $X^{(0)}$ has boundary conditions as follows:

$$x_{i\min} \le x_i^{(0)} \le x_{i\max} \tag{28}$$

where, *i* is the *i* th component of $X^{(0)}$, and *i*=1, 2, ..., *n*; $x_{i\min}$ and $x_{i\max}$ are the floor bound and ceiling bound of $x_i^{(0)}$, respectively.

Vertices of the initial discrete complex are defined as:

$$\begin{cases} x_i^{(1)} = x_i^{(0)} \\ x_i^{(j+1)} = x_i^{(0)} \quad (i \neq j) \\ x_j^{(j+1)} = x_{j\min} \\ x_i^{(n+j+1)} = x_i^{(0)} \quad (i \neq j) \\ x_j^{(n+j+1)} = x_{j\max} \end{cases}$$
(29)

totally, there are k vertices defined, and k = 2n+1.

2) One dimensional discrete search: Suppose $X^{(H)}$ is a vertex having the maximum value f(X) among all the initial discrete complex vertices, and $X^{(C)}$ is the geometric center of all vertices except $X^{(H)}$, and $x_i^{(C)} = \frac{1}{k-1} \sum_{l=1}^k X_i^{(l)}$, $l = 1, 2, \dots, k; l \neq H$. The search direction **S** is from $X^{(H)}$ to $X^{(C)}$, and $S_i = x_i^{(C)} - x_i^{(H)}$.

Then, the new vertex after discrete search $X^{(R)}$ is:

$$\begin{cases} x_i^{(R)} = x_i^{(H)} + \alpha S_i & (i = 1, 2, \dots, n) \\ x_i^{(R)} = \left[x_i^{(R)} \right] (i = 1, 2, \dots, p; p \le n) \end{cases}$$
(30)

where, α is increment factor; $\begin{bmatrix} x_i^{(R)} \end{bmatrix}$ is the discrete value that is nearest to $x_i^{(R)}$; *P* is number of discrete variables.

3) Constraints on vertices and search: As the vertices defined above are not subjected to constraints, some may locate in the areas out of feasible domain D. Define effective objective function F(X) as:

$$F(X) = \begin{cases} f(X), X \in D\\ M + C \sum_{i \in I} g_i(X), X \notin D \end{cases}$$
(31)

where, M is a constant much greater than f(X); C is a constant; $l = \{i | g_i(X) > 0\}$.

Starting from a vertex $X \notin D$ along search direction S, seek the minimum of $C \sum_{i \in I} g_i(X)$. If $C \sum_{i \in I} g_i(X) \le 0$, then the new vertex is located in domain D, and effective objective function of the vertex is F(X) = f(X).

4) Stop criterion: Define d_i the difference of the maximum and minimum discrete values of the *i* th variable, namely:

$$d_i = a_i - b_i \tag{32}$$

where, $a_i = \max \{x_i^{(j)}\}$; $b_i = \min \{x_i^{(j)}\}$; $i = 1, 2, \dots, n$, $j=1, 2, \cdots, \cdots, k$.

Let Δ_i and ε_i be the increments of continuous and discrete variables, and preset integer $E \in [n/2, n]$, and n is number of variables. If R is number of variables that satisfy $d_i \leq \Delta_i$ for discrete variable or $d_i \leq \varepsilon_i$ for continuous variable, the stop criterion of the iteration process is:

$$R - E \ge 0 \tag{33}$$

If stop criterion cannot be met after search of k-1times, move all vertices except $X^{(L)}$ to $X^{(L)}$ some distance, often 1/3 of the distance between each vertex and $X^{(L)}$.

5) Verification: The vertex $X^{(L)}$ that meets stop criterion is an optimal vertex in the domain of neighborhood of $X^{(L)}$, which cannot guarantee it is the optimal vertex in the whole domain. Therefore, let $X^{(L)}$ be another initial discrete point $X^{(0)}$. Repeat the optimization process mentioned above, and get a new vertex $X^{(*)}$ that meets stop criterion until $X^{(*)} = X^{(L)}$. Then, the final optimal vertex $X^{(*)}$ can be obtained.

4.2 **Program design**

According to strategies presented, a program is developed for solving the optimization model. The flowchart of the program is shown in figure 2.



Case study and discussions

5.1 Case study

5

1) Input: Costumer input and select of preset conditions for the optimization are as follows:

Bellows type: unreinforced U-shaped bellows;

Support: both ends rigid support;

Material: austenitic stainless steel 0Cr18Ni10Ti;

Temperature during operation: 400 °C;

Pressure: 1.8 MPa;

Applied axial movement in compression or extension:

20 mm; Fatigue life: 1500 cycles; Safe factor of fatigue life: 10; Inside diameter: 300 mm; Convolution pitch: 32 mm.

2) Output: Optimization results are as follows: convolution height 24.1 mm; bellows nominal material thickness of one ply 8 mm; number of bellows material plies 3; number of convolutions in one bellows 8.

The performance data of the bellows from optimization are as follows: bellows circumferential membrane stress due to pressure 61.49 Mpa; bellows meridional membrane stress due to pressure 9.43 Mpa; bellows meridional bending stress due to pressure 163.57 Mpa; bellows meridional membrane stress due to deflection 19.28 Mpa; bellows meridional bending stress due to deflection 983.45 Mpa; limiting internal design pressure based on column instability for single bellows 2.1 Mpa; Limiting design pressure based on inplane instability and local plasticity 1.8 Mpa; bellows theoretical axial elastic spring rate per convolution 6535.89 N/mm; total equivalent axial movement per convolution 5.14 mm.

5.2 Discussions

Intensive understanding of design criteria, forming methods and process flows is the basis of bellows design and optimization.

While optimization modeling and solving, the researcher should have practical engineering experience of bellows design and analysis, and know how to deal with common problems related to bellows design and use.

The weight factors of objectives for the objective function have much influence on the output of the optimization results, and they should be fully evaluated and verified.

6 Conclusions

Optimization method was employed in the design of an unreinforced U-shaped metal bellows. A 4-variable 14-

constraint optimization model was built based on Standards of the Expansion Joint Manufactures Association, Inc., 9th edition, and General Specification for Metal Bellows Expansion Joint, GB/T 12777-2008. Weight factors of mass and flexibility in the objective function were set according to the importance of objectives. Some special strategies were adopted to solve the optimization problem with mixed discrete and continuous variables. Case study shows that the optimization model presented and program developed can meet the requirement of design of the unreinforced Ushaped metal bellows with good quality and efficiency.

Acknowledgment

The study was supported by Fundamental Research Funds for the Central Universities of China XDJK2009B028, 2009C004. Correspondence author: Mingjin Yang.

References

[1] Expansion Joint Manufactures Association, Inc. "Standards of the Expansion Joint Manufactures Association, Inc., the 9th Edition". 2008. 1-4.43.

[2] Standardization Administration of the People's Republic of China. "General Specification for Metal Bellows Expansion Joints". Beijing: Standards Press of China. 2008. 1-36.

[3] Ye Yuanlie. "Mechanical Optimization Theory & Design". Beijing: China Metrology Publishing House. 2001. 1-88.

[4] Li Yongsheng, Li Jianguo. "Practical Technology of Bellows Expansion Joint —Design, Manufacturing and Installing". Beijing: Chemical Industry Press. 2000. 1-189.

[5] Zhang Weigang, Zhong Zhihua. "Advanced Design Methods". Beijing: China Machine Press. 2005. 1-137.

[6] Sun Jingmin. "Optimal Design of Machine, the 4th Edition". Beijing: China Machine Press. 2010. 1-250.

Numerical investigation of AC electrokinetically induced fluid flow in a circular microchannel

Vai Kuong Sin

Department of Electromechanical Engineering, University of Macau Av. Padre Tomás Pereira S.J., Taipa, Macao SAR, China

Abstract - AC-driven electrothermal flow is used to enhance the chaotic mixing through vortex flow in micro-fluidic systems by nearly an order of magnitude. Diffusion of heat and momentum can be increased because of the electrothermal force acting on the fluid through the applied AC electric field. Numerical simulation of vortex-assisted mixing in circular microchannel flow is investigated by solving a set of governing partial differential equations with computational fluid dynamics (CFD) software

Keywords: circular microchannel, electrothermal flow, electrokinetic effects.

1 Introduction

With fast development of biotechnology and bioengineering, micro or nanoscales fluidic mixing becomes a trend in numerous applications in engineering and technology[1]. Microfluidic devices such as micro-pump, micro-valve and micro-mixer have recently raised intensive research. AC electrokinetics, however, has received relatively little attention in micro-fluidical research. In micro scale, the flow is laminar because of the low Reynolds number and mixing is mainly due to molecular diffusion which is not efficient.

Numerical simulation of AC electrokinetically induced fluid flow can give a result which can be interpreted from different point of view. Electrokinetic effects have been studied both experimentally and numerically[2][3]. In previous study[3], the microchannel consists of a pair of coplanar electrodes at the bottom of a two-dimensional channel. Our effort in this study is to utilize computational fluid dynamics software to solve the Navier-Stokes equations in conjunction with the electrostatic equation and the energy equation in a circular channel

2 Mathematical formulations

The governing Navier-Stokes equations for incompressible AC electro-kinetically induced fluid flow can be written as[4]

$$\nabla \cdot \vec{V} = 0 \tag{1}$$

$$\frac{\partial V}{\partial t} + (\vec{V} \cdot \nabla)\vec{V} = -\frac{1}{\rho}\nabla P + \mu_k \Delta \vec{V} + \vec{F}$$
(2)

Where $\vec{V}(x, y, z) = ui + vj + wk$ is velocity with components u, v and w along the x, y and z-axis, P is pressure, ρ is density, μ_k is kinematic viscosity and \vec{F} is electrothermal force acting on the fluid. Noting that electro-thermal force in the incompressible fluid is given by[5]

$$\vec{F} = -0.5 \left[\left(\frac{\nabla \sigma}{\sigma} + \frac{\nabla \varepsilon}{\varepsilon} \right) \cdot E \frac{\varepsilon E}{1 + (\omega \tau)^2} + 0.5 |E|^2 \nabla \varepsilon \right]$$
(3)

where *E* is electric field, σ is the conductivity, $\varepsilon = \varepsilon_r \varepsilon_0$ equals the fluid's permittivity, ω represents the electric field's angular frequency, and $\tau = \varepsilon/\sigma$ is charge relaxation time of the fluid medium, and the incremental temperature-dependent changes are:

$$\nabla \varepsilon = \frac{\partial \varepsilon}{\partial T} \nabla T$$
 and $\nabla \sigma = \frac{\partial \sigma}{\partial T} \nabla T$ (4)

where T is temperature field. The first term on the right hand side of equation (3) is the Coulomb force, and is dominant at low frequencies. The second term is the dielectric force, and is dominant at high frequencies.

Equation (3) can be simplified by combining with (4). Gradients in temperature imply gradients in electrical permittivity and conductivity in the fluid. For water, the dependence of electrical conductivity in temperature is $(1/\sigma)(\partial\sigma/\partial T) = +2\%$ and the dependence of electrical permittivity on temperature is $(1/\epsilon)(\partial\epsilon/\partial T) = -0.4\% \text{ K}^{-1}[6]$. It is assumed that those values to be true for dilute solutions which reduces to

$$\vec{F} = \left(\frac{0.012\varepsilon\nabla T \cdot E}{1 + (\omega\tau)^2}\right)E + 0.001\varepsilon |E|^2 \nabla T .$$
 (5)

The electro-thermal force depends on the value of electric field and temperature field, and thus, we consider the electrostatics and thermal equations in the form

$$-\nabla \cdot (\varepsilon_0 \varepsilon_r \nabla V) = 0 \tag{6}$$

$$\nabla \cdot (-k\nabla T) = Q - \rho C_p \vec{V} \tag{7}$$

where k is thermal conductivity, C_p is specific heat at constant pressure and Q is heat source and is given as

$$Q = \sigma(E_{ems})^2 \tag{8}$$

where E_{ems} is the root mean square value of the AC electric field.

3 Numerical modelling

Electrokinetic effects are studied numerically by solving Equations (1) through (8) in a circular channel with radius R=50 μ m and length L=600 μ m. Let's consider a micro-fluidic circular channel consisting of a pair of electrodes kept at the lower half of the surface as shown in Figure 1. There is a gap of 25 μ m between those two electrodes. The electrodes are energized at applied voltage of V= 6 V_{rms} and frequency of f = 200 kHz.



Figure 1. Coordinate system of electrothermal flow in circular channel

A numerical solution of the full Navier-Stokes, electrostatic and energy equations is sought to solve this electrothermal flow problem. To simulate the fluidic characteristics of the electrothermal flow, commercially available Computational Fluid Dynamics software called Comsol Multiphysics is adapted to perform all simulation. Finite elements have been refined continuously to resolve the large velocity and temperature gradient near the wall. The boundary condition and fluid physical properties in the micro-channel electrothermal flow are given as follow:

(1) Incompressible Navier-Stokes: No slip and no-flux boundary conditions are imposed on all solid surfaces. The boundary conditions of tube are defined as

$$V(x, R, R) = 0 \tag{9}$$

To specify the velocity distribution at the inlet is not an easy work since experimental data are not available at this time. To simplify the analysis, we consider the following assumption. In extremely small scale, Reynolds number becomes low and which results in laminar flow, so it is assumed that the inlet velocity distribution is laminar for the numerical simulation and is equal to

$$u = u_{ave} \left(1 - \frac{r^2}{R^2} \right)$$
 (10)

where $u_{ave} = 67 \ \mu m/s$ is the average inlet velocity inside the channel. The fluid is 0.05M KCL aqueous solution and is treated as Newtonian fluid with the physical properties defined as[7]: dynamic viscosity $\mu = 6.920215314 \ x \ 10^{-4} \ Pa^{+}s$.

(2) Electrostatics: The applied voltages of two electrodes are $V = \pm 0.5 V_{rms}$, and the other boundaries are zero charge, V= 0. The physical properties of the fluid are defined from the references[8][9]: relative permittivity $\varepsilon_r = 77.5$, $\varepsilon_0 = 8.854188 \times 10^{-12}$ F/m and electric conductivity $\sigma = 0.66$ S/m.

(3) Heat Transfer: The selected boundary condition in this model is isothermal (i.e., $-n \cdot (-k\nabla T) = 0$) at the two electrodes and the other boundaries are provided a constant heat flux $q_0 = -k^*(T-T_0)/1$ [mm]. Then,

$$T_0 = 310.15 \text{ K}$$
 and $-n \cdot (-k\nabla T) = -k^*(T-T_0)/1[mm]$

The physical properties of the fluid are defined from the references[10][11]: thermal conductivity k = 0.61 W/ (m*K), density $\rho = 1003.725$ kg/m³.

4 Results and discussion

Results obtained from circular microchannel and twodimensional model are in good agreement at location far away from the electrodes with maximum error less than 1%. Large discrepancy occurs for values near the electrodes. This is because of the fact that the geometry is simplified to two dimensions. When the electrodes begin to supply the voltage, vortex appears obviously at the top of the right electrodes as shown in Fig. 2. Vortex affects the factors on enhanced mixing. Once the flow becomes steady, the distribution of the flow field is smooth and the flow roughly along the wall geometric appearance. Pressure near electrodes is greater than that near the upper surface. This leads flow rolling up and forms a trailing vortex. Figure 2 also shows velocity and temperature distribution at the XOY plane. The thermal color scale was illustrated to the brightest slice at the high temperature and to dark areas at the low temperature $(T_{min}=314.5K)$. Electrothermal forces distort the parabolic flow around the electrode gap to produce a circulating flow.



Figure 2. Velocity and temperature field at the plane z=0 in circular microchannel with electric field of $V_{rms} = 6$ volt



Figure 3. Secondary flow in circular channel at (a) x=0.0002m, (b) x=0.00027m, (c) x=0.0003m, (d) x=0.00033m, (e) x=0.0004m

The secondary flow of the channel in different circular cross sections is shown in Figure 3. Two opposed vortices appear in the regions which is above the gap of electrodes and in the vicinity of the wall boundary where viscous stress are dominant. Because of non-uniform electric field which induces non-uniform zeta potential, and thus, the vortex is located close to the electrode. If vortices are intentionally controlled, then it could be used for enhanced mixing.



Figure 4. Recirculation speed versus voltage at 1 kHz.



Figure 5. Maximum temperature versus voltage at 1 kHz.

Maximum recirculation velocity and temperature distribution at various supply voltages at 1kHz for circular microchannel are shown in Figures 4 and 5, respectively. The electrothermal effect arises through non-uniform Joule heating caused by ac electric fields. The fluid properties are a function of temperature and hence they vary spatially. When applying a non-uniform electric field to a fluid, thermal energy will be transferred into the fluid, leading to a nonuniform temperature field. Once the supply voltage rose, the thermal energy transferring into the fluid would also increase. Moreover, the electrothermal force is produced by the interaction of the electric field and the gradients of the conductivity and permittivity, giving rise to fluid motion. As a result, an increase of the supply voltage also tends to increase the velocity of the fluid motion.

5 Conclusions and future work

Electric field is generated by electrodes in a microfluidic channel, which creates vortex above the gap of electrodes and near the wall. To simulate the model by COMSOL Multiphysics, numerical findings show that velocity and temperature profiles obtained from circular microchannel are under-estimated when compared with those obtained from two-dimension model. But they are still in good agreement, especially for the region far away from electrodes. Geometric effect is required to take into account in the region near electrodes through micro-fluidic channel because of the secondary flow near the electrodes. Maximum velocities of the circular microchannel take place somewhere near the surface on the right electrode. Based on these results, it is suggested that microstirring enhanced binding simulations can be performed in the future.

6 Acknowledgments

The research was partially supported by Research Committee of University of Macau through the grants RG079/09-10S//11T/SVK/FST and MYRG040(Y1-L1)-FST12-VKS

7 References

[1] J. P. Shelby and D. T. Chiu, Lab on a Chip 4, 168 (2004)

[2] Ramos, A., et al., Journal of Physics D-Applied Physics 31, 2338 (1998)

[3] D. Wang, M. Sigurdson and C. Meinhart, Exp. Fluids 38, 1 (2005)

[4] F. M. White, Viscous Fluid Flow, McGraw-Hill book Company, New York (1974)

[5] D. Wang, Investigation of AC Electrokinetics using Micro-PIV, (2004)

[6] D. R. Lide, CRC Handbook of Chemistry and Physics, CRC Press, New York (2000)

[7] S. Mao and Z Duan, International Journal of Thermophysics 5, 1510 (2009)

[8] D. Q. Craig, Dielectric Analysis of Pharmaceutical Systems, Taylor & Francis (1996)

[9] Y. Lin, G.J. Gerfen, D. L. Rousseau, and S.R. Yeh, Analytical Chemistry 75, 5381 (2003)

[10] MLV Ramires and CA Nieto de Castro, International journal of thermophysics 21, 671 (2000)

[11] I.M. Abdulagatov and V.I. Dvoryanchikow, Fluid phase equilibria 143, 213 (1998)

The RTO-RTDB Real-Time Data Model

Z. Ellouze¹, N. Louati², and R. Bouaziz² ¹CES-ENIS, Sfax University, BP 1173, 3038 Sfax, Tunisia ²MIRACL-ISIMS, Sfax University, BP 1088, 3018 Sfax, Tunisia

Abstract—This paper presents a data model of a realtime object-oriented database. In addition to the logical consistency constraints, a real-time database must support both temporal aspects of data and timing constraints of transactions upon the data. The RTO-RTDB data model support these requirements by specifying objects that contains: time-constrained data; time-constrained methods; and concurrency control mechanisms.

Keywords: Data Model; Real-Time; Object Model; Database;

1. Introduction

A Real-Time (RT) database is a database in which both the data and the transactions upon the data may have timing constraints. The mostly used data model for RT databases is the relational model [1]. However, due to the nature of many RT applications that must handle complex real world objects with short deadlines, many researchers believe that the object-oriented model is more natural and powerful than the relational model [2]. A RT database have all requirements of traditional databases, such as the management of accesses to structured, shared and permanent data, and they require management of time-constrained data and time-constrained transactions [3]. RTO-RTDB (Real-Time Object for Real-Time DataBases) is a RT object-oriented database model that incorporates these concepts [4].

The RTO-RTDB model is based upon an earlier model, called *real-time object-oriented data model* [5], on which it incorporates time-constrained data, time-constrained methods and concurrency control mechanisms. As shown in Figure 1, each RTO-RTDB object has four components: (i) a set of RT attributes, (ii) a set of RT methods, (iii) a mailbox, and (iv) a local controller.

This paper represents an extension of the RTO-RTDB data model. We have drawn from our experience in the design of RT databases [6] to identify the basic research issues involved in the design of RT object-oriented databases. The next section provides background information on RT databases. The third section surveys related work, and summarizes the features that are unique to RT databases. The fourth and the fifth sections introduce the RTO-RTDB object model. They present the native RTO-RTDB components introduced in [4] and advanced ones proposed in this paper, respectively. The last section concludes the paper and describes issues for future investigations.

Fig. 1: Overview of the RTO-RTDB object model.

2. Real-Time Databases

A RT database is a time-constrained database designed to handle not only transactions with timing constraints, but also data with timing constraints. The timing constraints on data are defined as how well the content of the database models the actual state of the real world while the timing constraints on transactions are expressed in the form of deadlines which indicate a certain time in the future by which the transactions must be completed. In general, there are two forms of data timing constraints: absolute and relative. The absolute timing constraint requires that a data item's age must be within a certain interval of the current time. The relative timing constraint represents the required correlation among data used together [1]. To meet the operational deadlines from event to system response, a RT database may apply different timing constraints on transactions such as absolute timing constraints (e.g. execution time, earliest start time, latest allowed finish time), periodic timing constraints (e.g. frequency of transaction initiation), and relative timing constraints (e.g. precedance constraints, distance constraints, and freshness constraints). Given the added dimension of time on data and transactions, two of the interesting areas of study in RT databases are that of transactions scheduling and concurrency control policies. Not only the schedule must meet timing constraints of transactions, it must also maintain data temporal consistency.

A RT database manipulates RT data and executes RT transactions [1]. RT data are divided into two types: sensor data and derived data. Sensor data are periodically collected from the physical environment through sensors, whereas derived data, they are computed from sensor data. When a sensor data item changes, all derived data items that are based on it need to be recomputed. RT transactions are



classified into two categories: update transactions which are used to update values of RT data in order to reflect the state of the real world and user transactions which represents user requests [7]. Update transactions are executed periodically to update sensor data, or sporadically to update derived data. User transactions arrive aperiodically.

3. Related work

In recent years, several works on RT databases have been proposed to deal with data modeling issues. Only a few of these works address RT object-oriented data modeling.

DiPippo and Ma [8] describe the RTSORAC model which qualifies the basic parameters for modeling a RT database object-oriented data model. It is composed of three components: objects, relationships, and transactions. Objects represent database entities. Relationships represent associations among the database objects. Transactions are executable entities that access the objects and relationships in the database. RTSORAC model is too complex to support hard RT applications. It support features such as semantic concurrency control which make bounding execution time of many application difficult. Nevertheless, its support for transaction timing constraints and temporal consistency do make it appropriate for soft RT applications.

In [9], Lee et al. describe a simple RT object-oriented data model with atomic objects and a class manager. Atomic objects are basic entities that ensure atomicity of transactions. Class manager is the major vehicle that lessens the complexity involved in transaction management. In [9], the authors have introduced new features to support active pursuit of timely and deterministic processing of transactions that are why it has better support for RT transactions scheduling than RTSORAC model [8]. But it has still weaknesses for modeling RT databases. It can not capture different aspects that are required for RT databases. It does not account for data characteristics, Quality of Service (QoS), and quality of data.

The G-CPN [10] model allows the modeling of syntactical and semantical features related to the objects in a RT database. The semantical features can be mapped to Petri net constructions. G-CPN model has introduced extensions in an object colored Petri net to model a RT database. But certain weaknesses exist. First, the model does not to capture the basic features of RT databases, as well as the expected activities against the database. Second, it does not describe schedulability aspect of transactions. Third, it uses Petri net notations which depend on designer's knowledge.

The RODAIN [11] model is a RT object-oriented database architecture for intelligent networks. It supports two kind of attributes: regular and RT. A RT object has predefined attributes for isolation level and access type. It is referenced by RT transactions that have an explicit deadlines. In [11], Taina and Raatikainen have tailored a data model for the needs of telecommunications applications which make it not general enough to describe any RT database applications.

The BeeHive [12] model extends traditional objectoriented data models by incorporating semantic information regarding RT, fault-tolerance, and QoS requirements. It has some similarity in terms of the structure of objects to the RTSORAC object model [8]. One of the main differences is that while RTSORAC model holds only RT and approximation requirements, BeeHive model supports a rich set of types of requirements and their trade-offs. One limitation of the BeeHive model is that transaction execution times should be determined offline for admission control. This approach is only applicable to a limited set of RT data services in which transactions and their arrival data access patterns are known in advances [13].

In [5], Idoudi et al. describe a RT object-oriented data model where objects contain a set of RT attributes and a set of RT methods. RT attributes are divided into two types: sensor attribute and derived attribute. RT methods are classified into three classes: periodic methods, sporadic methods, and aperiodic methods. In [5], the data model offers solutions to manage data and transactions characteristics. But, it does not describe the schedulability aspect.

4. Native RTO-RTDB

Data objects are classified into either non RT or RT data [7]. A non RT data is a classical data that does not become outdated due to the passage of time, whereas a RT data has a validity interval beyond which it become useless. The RTO-RTDB data model defines two types of attributes: classical attributes and RT attributes. A classical attribute is used to store a non RT data, while a RT attribute stores a RT data. RT data are classified into either sensor or derived data [7]. Thereby, RTO-RTDB data model defines two kinds of RT attributes: sensor attribute and derived attribute. A sensor attribute is used to store a sensor data which is periodically updated in order to reflect the real world state of the environment. A derived attribute is used to store a sensor attribute value, used in its computation, changes.

RTO-RTDB model specifies three types of RT methods: periodic methods, sporadic methods, and aperiodic methods. Periodic methods update periodically sensor data. They are write-only methods that obtain the state of the real world and write the sensed data to the database. Sporadic methods calculate sporadically derived data. The access mode of the sporadic method to derived data is always write. Aperiodic methods do not write any RT data, but they can read/write non RT data and only read RT data. In RTO-RTDB data model, each method execution is considered as a RT transaction

The RTO-RTDB data model originally published in [4] and outlined in the previous section, was missing important features in terms of object orientation and RT database

properties. On the object-oriented side, native RTO-RTDB lacks the notion of object identity, encapsulation, nesting, inheritance which are the main properties of object-oriented databases. On the RT side, native RTO-RTDB extends the object model with RT attributes and RT methods, but still lacks features for expressing RT data features, expressing RT transactions features, handling RT transactions concurrency control, scheduling RT transactions, expressing RT data versions. To address these issues, a solution is proposed hereafter with an enhanced RTO-RTDB data model.

5. Advanced RTO-RTDB

In this section, we first define our model of a RT objectoriented database. Then, we describe our enhanced RTO-RTDB object model and how it extends the object model to suit RT databases features.

Our RT object-oriented database is a collection of RT classes and instances of these classes. Both RT classes and instances are referred to as RT objects. A RT class defines a set of RT attributes for its instances and RT procedures through which instances can be handled. The RT procedures associated with a RT class are referred to as RT methods, and a RT method may invoke other methods on other objects in the RT database. In this model, we allow inheritance of attributes and methods between classes.

The RTO-RTDB data model includes features that support the requirements of a RT database into an extended objectoriented data model. It has a main component that models the properties of a RT object-oriented database which is RT object. RT objects represent RT database entities. A RT object consists of six components, < N, A, M, C, LC, B >, where N is a unique name, A is a set of attributes, M is a set of methods, C is a set of constraints, LC is a local controller and B is a set of behaviors. Figure 2 illustrates an example of an Aircraft RT object for storing information about an air traffic control system in a database.

5.1 Attributes

The second component of a RT object, A, is a set attributes, where each attribute is characterized by $\langle N, T, M, V, Va, [Ti], [Vd], [Mde], [Uo], [Nv] \rangle$.

- N (Name): is the name of the attribute.
- T (Type): is the type of the attribute which can be integer, real, string, etc.
- M (Multiplicity): indicates how kinds of values or objects a RT attribute can obtain.
- V (Visibility): represents the visibility of the attribute: public, protected or private.
- Va (Value): is used to store the real world attribute value captured by the last update correspondent method. This field is used by the system to determine the logical consistency constraints of the attribute value.
- Ti (Timestamp): is used to store the instant at which the attribute value was last updated. The timestamp

determines the temporal consistency constraint of the attribute value. For example, in the Aircraft RT object, there is a property for storing the speed, called *speed*, to which a sensor periodically provides readings. This update is expected every 30 seconds, thus the *speed* property is considered temporally inconsistent if the update does not occur within that time. The timestamp value of the *speed* property must be utilized by the RT database system to determine that the update operation did not happen as expected. There are many ways to define timestamps [6]. In our work, we consider that the timestamp is the time when the value is written.



Fig. 2: Aircraft RT object.

- Vd (Validity duration): it indicates the amount of time during which the attribute value is considered valid. This field permits to determine, in association with the timestamp, the absolute consistency constraint of the RT attribute. A RT data is considered absolutely fresh with respect to time as long as the age of the data value is within a given interval [7]. For instance, the *speed* value is considered valid if the current time is earlier than the timestamp of *speed* followed by the length of the absolute validity interval of *speed*, i.e. {*speed.Ti* + *speed.Vd* > *currentTime*}.
- Mde (Maxiumum data error): is used to memorize the absolute maximum data error tolerated on the attribute

value [6]. This value is the upper bound of the deviation between the current attribute value in the RT database and the reported one. Recently, the demand for RT services has increased in most RT database based applications where it is desirable to execute transactions within their deadlines. They also have to use fresh data in order to reflect the real world state. However, it seems to be difficult for the transactions to both meet their RT constraints and to keep the database consistent. To support this kind of applications, the data error concept is introduced in [14] to indicate that data stored in the RT database may have some deviation from its value in the physical world.

- Uo (Update operation): is used to update the value and timestamp fields of a RT attribute. For example, in the Aircraft RT object, there is a RT method, called *updateSpeed()*, which periodically updates the *speed* RT attribute.
- Nv (Number of versions): is used to preserve RT attribute version history. The multi-version attributes permits to maintain for every attribute multiple versions for a data item. This minimizes data access conflicts between RT transactions and reduces the deadline miss ratio [15]. In order to respect the RT database size, the number of versions of each RT attribute is limited. It does not have to exceed a maximum data versions number [15].

Note that the fields N, T, M, V, and Va characterize classical attributes as well as RT attributes, whereas the fields Ti, Vd, Mde, Uo, and Nv characterize only RT attributes. Here an example of two attributes in the Aircraft RT object: the first is a classical attribute and the second is a RT attribute.

 $\{N = destination, T = string, M = 1, V = private, Va = Paris\}$

 $\{N = speed, T = real, M = 1, V = private, Va = 600, Ti = 01/05/2012 \ 10 : 05 : 23, Vd = 30s, Mde = 10, Uo = updateSpeed(), Nv = 5\}$

5.2 Methods

The third component of a RT object, M, is a set of methods. Each method is characterized by $\langle N, [V], [Arg], [Exc], Op, [Mc], Mst, Mrd, Mct, [Rt], [Per], Pri, Cp >$.

- N (Name): denotes the name of the RT method.
- V (Visibility): indicates whether a RT method is public, protected, or private.
- Arg (Arguments): is a set of arguments for the RT method, where each argument has the same structure as an attribute, and is used to pass information in the method.
- Exc (**Exc**eptions): is a set of exceptions that may be raised by the RT method to signal that the method has terminated abnormally.

- Op (**Op**erations): is a set of operations which represent the impementation of the method.
- Mc (Method constraints): is a set of methods constraints. A method constraint is of the form < N, OpSet, Pred, Er > where N is the name of the method constraint, OpSet represents a subset of the operations in Op, Pred is a boolean expression which is specified over OpSet to express execution constraints, timing constraint and precedence constraints, and Er is a enforcement rule. The enforcement rules are used to specify the actions to take if the predicate (i.e. Pred) evaluates to false. A complete definition of an enforcement rule is described in the next subsection on constraints. Here is an example of a method constraint predicate in the Aircraft RT object:

Pred: getLane().Mct < currentTime + 5s

A deadline of currentTime +5s has been specified for the completion of the *getLane()* method. Note the use of the *Mct* property which represents the completion time of the executable method.

- Mst (Method start time): indicates the execution start time of the RT method.
- Mrd (Method relative deadline): specifies the deadline of a method execution.
- Mct (Method completion time): indicates the time at which the method finishes its execution.
- Rt (Return type): specifies which kinds of value or object a RT method can return.
- Per (**Per**iod): indicates the frequency of the RT method initiation.
- Pri (**Pri**ority): Priority specifies the priority order of a RT method.
- Cp (Concurrency policy): specifies the concurrency policy of a RT method. A concurrency policy may be reader, writer or parallel [4]. A reader RT method implies that multiple calls from concurrent methods may occur simultaneously and will be executed simultaneously if there is no writer methods using one or more data that the reader method needs. A writer RT method implies that multiple calls from concurrent methods may occur simultaneously and will be treated as soon as concurrency on data permits its execution. A parallel RT method is a method whose actions do not use any data of the database in reading mode nor in writing mode.

5.3 Constraints

The fourth component of a RT object, C, is a set of constraints. Constraints permit the specification of a correct object state. Our constraint specification is based on the model introduced in [8], where a constraint is characterized by < N, AttrSet, Pred, Er >.

• N (Name): denotes the name of the constraint.

- AttrSet (Attribute Set): is a subset of attributes of the RT object.
- Pred (**Pred**icate): is a boolean expression which is specified using attributes from the AttrSet field. The Pred property can be used to state the logical and timing consistency constraints of the RT data stored in the RT object by referring to the value (i.e. Va), timestamp (i.e. Ti), validity duration(i.e. Vd), and maximum data error (i.e. Mde) fields of the RT attributes in the set.
- Er (Enforcement rule): is executed when the predicate evaluates to false. The enforcement rules are of the form < *Exc*, *Op*, *Mc* >. As with RT methods, Exc is a set of exceptions which the enforcement rule can signal, Op is a set of operations which represent the implementation of the enforcement rule, and Mc is a set of method constraints on the execution of the enforcement rule.

For example, as mentioned earlier, the Aircraft RT object has a *speed* RT attribute which is is updated with the latest sensor reading every 30s. To maintain the temporal consistency constraint of this attribute, the following constraint is defined:

 $N: speed_Avi$ AttrSet: speed Pred: speed.timestamp <= currentTime - 30s Er: if Missed <= 3 then speed.timestamp = currentTime; Missed = Missed + 1; signal speedWarning;else speedAlert();

The enforcement rule specifies that if only one, two or three of the readings have been missed, a counter is incremented stating that a reading was missed and a warning is signaled using the exception *speedWarning*. If more than three readings have been missed, the *speedAlert()* RT method is called which might lead to a message being sent the the Aircraft operator.

5.4 Local Controller

Because of the dynamic nature of the real world, more than one operations may send requests to the same RT object. Concurrent execution of these operations allows several methods to run concurrently within the same object. To handle this essential property of RT database systems, we associate to each RT object a local concurrency control mechanism, that manages the concurrent execution of its operations. Thus, the RT object receives messages (or requests) awaking its local controller that checks the timing constraints attached to messages and selects one message following a special scheduling algorithm. The local controller verifies the concurrency constraints with the already running methods of the RT object. Then, it allocates a new thread to handle the message when possible. When an operation terminates its execution, the corresponding thread is released and concurrency constraints are relaxed [4].

As illustrated in Figure 2, the local controller is made of four components: state controller, deadline controller, freshness controller, and concurrency controller.

- State controller: It is used to avoid RT object overload. In fact, the RT object has a pool of schedulable resources. When no schedulable resources are available, the RT object may be overloaded. When receiving a message, the state controller checks if it can be handled. Otherwise, the message is rejected.
- Deadline controller: It controls RT operation validity. If the current time is greater than operation deadline, operation will be aborted. Otherwise, if the verification step succeed, then the operation is transferred to the freshness controller.
- Freshness controller: It checks the freshness of acceded data just before an operation commits. This way, the data accessed by committed operations are always fresh at commit time. If the accessed data is fresh, operation can be executed. Otherwise, the operation will be blocked.
- Concurrency controller: The main objective of this component is to verify the concurrency constraint between operations. If it detects a conflict, it aborts the operation having the lowest priority.

5.5 Behavior

A RT database is a collection of RT objects which are used to manage time critical dynamic systems in the real world. Each RT object may own one or several behaviors. For each of these behaviors is defined a message queue for saving the messages received by the RT object. Figure 2 depicts a state-machine diagram that provides a simplified view of the *Aircraft* behavior. The *Aircraft* has four states: *TakeOff*, *Flying*, *Update*, and *Landing*. Periodically, it enters in the *Update* state for updating the *Aircraft* sensors values. The sensors values update has to be done with a period of 3 s and lasts 2 s. Then it returns to the state which activated the update transition.

6. Implementation of the RTO-RTDB data model

Throughout this section, we will use a running example to illustrate a concrete implementation of our RTO-RTDB data model. We illustrate our proposal on a freeway traffic control system which consist of a large collection of data describing the current traffic state. We focus precisely on modeling the RT database used in storing the data collected from sensors.

6.1 Application Modeling

As depicted in Figure 3, the freeway traffic control system architecture consists of seven entities respectively dedicated to: acquire data from the environment (*Sensor*),



Fig. 3: Freeway Traffic Control System Class Diagram.

indicate the positions of sensors (*Location*), monitor and analyze the traffic flow (*Controller*), indicate the notified incidents (*Incident*), represent a transport infrastructure road that links two conurbations (*RoadLink*), depict a part of a route (*RoadSegment*), represent physical entity (*Vehicle*). *RoadSegment* and *Vehicle* represent the description of two physical elements that are supervised by the controller. In addition, they are characterized by one or more RT data which could determine theirs evolution. These RT data are classified into either sensor data or derived data. In fact, each vehicle has two sensor data (i.e. *speed* and *length*) which are periodically updated to reflect its state and each road segment has two derived data (i.e. *trafficVolume*, and *trafficOccupancy*) which are calculated from sensor data.

6.2 Implementation

This section provides code implementation of the freeway traffic control system case study. This code includes using new concepts defined in the RTO-RTDB data model and those imported from the RT Java language as well [16]. Typically the RTO-RTDB object model transformation is described by the following steps. First we give the Java code defining the RTAttribute class that factorizes the properties of a RT data (i.e. Va, Ti, and Vd).

abstract class RTAttribute
{private int Va;
private Timestamp Ti;
private int Vd;
public void setVa(int v){Va = v;}
public int getVa(){return Va;}

```
public void setTi(Timestamp t){Ti = t;}
public Timestamp getTi(){return Ti;}
public void setVd(int v)
{Vd = v;}
public int getVd()
{return Vd;}
public RTAttribute(int v)
{Vd=v;}}
```

Second, we define the Sensor and Derived classes which represent sensor and derived RT data.

```
class Sensor extends RTAttribute
{private int Mde;
public void setMde(int m)
{Mde = m;}
public int getMde()
{return Mde;}
public void periodicUpdate(int v)
{super.setVa(v); super.setTi(//current time)}
public Sensor(int v,int m)
{super(v);Mde=m;}
}
class Derived extends RTAttribute
{public void sporadicUpdate(int v)
{super.setVa(v); super.setTi(//current time)}
public Derived(int v)
{super(v);}
```

Third, we define the classes that compose the local controller component: **D**eadlineController (DC), **S**tateController (SC), **F**reshnessController (FC), and ConcurrencyController (CC).

public class DC implements Runnable{
public void run(){setupDC();}
private void setupDC(){//...}
private boolean checkDeadline(){//...}
//...}
public class SC implements Runnable {
public void run(){setupSC();}

```
private void setupSC(){//...}
//...}
public class FC implements Runnable{
public void run(){setupFC();}
private void setupFC(){//...}
//...}
public class CC implements Runnable{
public void run(){setupCC();}
private void setupCC(){//...}
//...}
```

Fourth, we give the RT Java code defining the LocalController (LC) component.

```
public class LC implements Runnable{
private DC dc;private CC cc;
private FC fc;private SC sc;
public void run()
{//Create Runnable controllers
dc = new DC(); cc = new CC();
fc = new FC(); sc= new SC();
// Create and start RT Threads
RealTimeThread dcThread = new
RealTimeThread(null,null,null,null,dc);
RealTimeThread ccThread = new
RealTimeThread(null,null,null,null,cc);
RealTimeThread fcThread = new
RealTimeThread(null,null,null,null,null,fc);
RealTimeThread scThread = new
RealTimeThread(null,null,null,null,sc);
dcThread.start();ccThread.start();
fcThread.start();scThread.start();}
```

Fifth, we give RT Java code defining the RealTimeClass class. An instance of this class will encapsulate an instance of the LocalController defined above.

```
class RealTimeClass
{private LC lc;
public RealTimeClass()
{LC lc=new LC();
RealTimeThread lcThread = new
RealTimeThread(null,null,null,null,null,lc);
lcThread.start();}}
```

Sixthly, we give the Java code defining the Vehicle class of our freeway traffic control system application as it can be written by a Java programmer.

```
public class Vehicle extends RealTimeClass {
  private Sensor speed;private Sensor length;//...
  public Vehicle(){super();}
  public void setSpeed(Sensor s){speed=s;}
  public Sensor getSpeed(){return speed;}
  //...}
```

7. Conclusion

The design of a RT database, which is by definition a database system, has to take into account the management of many components such as data, transactions, scheduling policy, and concurrency control protocol. In order to deal with time-constrained data, time-constrained operations, parallelism, and concurrency property inherent to RT databases, we introduced the RTO-RTDB data model. This model combines features of RT databases and object-oriented databases. It specifies that instances of a class will encapsulate RT attributes, RT methods and a local concurrency control mechanism. RT attributes are divided into two types: sensor attribute and derived attribute. RT methods are classified into three classes: periodic methods, sporadic methods, and

aperiodic methods. The local controller manages concurrents execution of RT methods.

We are currently undertaking a research work to establish a RT query and transaction specification languages that allow users to specify time constraints for query (or transaction) processing, and that are based on our RTO-RTDB data model.

References

- [1] K. Ramamritham, "Real-time databases," *Distributed and Parallel Databases*, vol. 1, no. 2, pp. 199–226, 1993.
- [2] W. Kim, Object-Oriented Database Systems: Promises, Reality, and Future. Modern Database Systems. Addison Wesley, 1995, pp. 255– 280.
- [3] A. Bestavros, K.-J. Lin, and S. Son, *Real-Time Database System: Issues and Applications*. Kluwer Academic Publishers, 1997, ch. Advances in Real-Time DataBase Systems Research, pp. 1–14.
- [4] N. Louati, C. Duvallet, R. Bouaziz, and B. Sadeg, "RTO-RTDB: A real-time object-oriented database model," in *In Proceedings of the International Conference on Parallel and Distributed Computing and Systems*. ACTA Press, 2011.
- [5] N. Idoudi, C. Duvallet, B. Sadeg, R. Bouaziz, and F. Gargouri, "Structural model of real-time databases: An illustration," in *ISORC*, 2008, pp. 58–65.
- [6] N. Idoudi, N. Louati, C. Duvallet, B. Sadeg, R. Bouaziz, and F. Gargouri, "A framework to model real-time databases," *International Journal of Computing and Information Sciences (IJCIS)*, vol. 7, no. 1, pp. 1–11, 2010.
- [7] K. Ramamritham, S. H. Son, and L. C. DiPippo, "Real-time databases and data services," *Real-Time Systems*, vol. 28, no. 2-3, pp. 179–215, 2004.
- [8] J. Prichard, L. DiPippo, J. Packham, and V. Fay-Wolfe, "RTSORAC: A Real-Time Object-Oriented Database Model," in *Proc. of the* 5th *Intl. Conf. on Database and Expert Systems Applications (DEXA'94)*, Springer-Verlag, Ed., London, UK, 1994, pp. 601–610.
- [9] J. Lee, S. H. Son, and M.-J. Lee, "Issues in developing Object-Oriented Database Systems for Real-Time Applications," in *Proceeding of the IEEE Workshop on Real-Time Applications*, vol. 26. Washington, DC, USA: IEEE Computer Society, 1994, pp. 136–140.
- [10] M. L. Perkusich, M. de Fatima, Q. Turnell, and A. Perkusich, "Object-oriented real-time database design based on petri nets," in *In Proceedings of the International Workshop on Active and Real-Time Database Systems*. Springer, 1995, pp. 104–121.
- [11] J. Taina and K. Raatikainen, "Rodain: a real-time object-oriented database system for telecommunications," in *Proceedings of the* workshop on on Databases: active and real-time, ser. CIKM '96. New York, NY, USA: ACM, 1997, pp. 10–14. [Online]. Available: http://doi.acm.org/10.1145/352302.352306
- [12] J. A. Stankovic and S. H. Son, "Architecture and object model for distributed object-oriented real-time databases," in *ISORC*, 1998, pp. 414–424.
- [13] S. Kim, S. H. Son, and J. A. Stankovic, "Performance evaluation on a real-time database," in *IEEE Real Time Technology and Applications Symposium*, 2002, pp. 253–265.
- [14] M. Amirijoo, J. Hansson, and S. H. Son, "Specification and management of qos in real-time databases supporting imprecise computations," *IEEE Trans. Computers*, vol. 55, no. 3, pp. 304–319, 2006.
- [15] E. Bouazizi, C. Duvallet, and B. Sadeg, "Management of qos and data freshness in rtdbss using feedback control scheduling and data version," in *Proceedings of 8th IEEE International Symposium* on Object-oriented Real-time distributed Computing (ISORC'2005), Seattle, United State, May 18-20 2005, pp. 337–341.
- [16] E. Bruno and G. Bollella, *Real-Time Java Programming With Java RTS*. Prentice Hall, 2009.

The Case for Meta-modeling Frameworks Specialisation

S. Temate¹, L. Broto¹, and D. Hagimont¹ ¹IRIT/ENSEEIHT, 2 rue Charles Camichel - BP 7122 31071 Toulouse cedex 7 (France)

Abstract—Domain Specific Modeling Languages (DSMLs) are increasingly used today as they allow users to express strategies without being programming experts. This is particularly true for graphical DSMLs, inspired from UML. Implementing a DSML means providing specific editors and an execution machine.

Many experiments showed that implementing specific graphical editors is much manpower consumming. Our analysis is that most framework for building such editors (e.g. EMF/GMF) are generic, i.e. aim at fulfilling the requirements of any field, which leads to increased complexity. If domain specialization was successful for languages, why don't we apply it to frameworks ?

In this paper, we argue for Domain Specific Modeling Frameworks (DSMFs) for building DSML editors in specific application fields. We describe our experiment in implementing such a DSMF devoted to component-based DSML. We then generalize this proposal: we believe that DSMFs could be metamodeled in the same way as DSMLs are metamodeled.

Keywords: Meta-modeling, DSML, autonomic management, components

1. Introduction

In the context of the TUNe project [1], [2], we are investigating the design and implementation of an autonomic administration system, relying on a component-based middleware. In order to help the definition of administration policies, we explored the design of graphical Domain Specific Modeling Languages (DSMLs) for describing software to deploy, configure, monitor and reconfigure. From our experiments, we noticed that there is a need to design different DSMLs according to the considered administration facets and application domains. To implement the editors associated with these DSMLs, we initially relied on well known frameworks like EMF/GEF/GMF, but we found that they are overly complex. Therefore, we decided to implement our own framework. This framework called yTUNe allows designing a DSML in terms of its metamodel and provides a fully automatized generation of the associated editor. This is possible because we are focussed on component-based DSMLs (in the context of TUNe), thus allowing to wire in yTUNe how the concepts from the DSML metamodel are represented in the editor.

From a wider point of view, we consider that:

- yTUNe is a Domain Specific Modeling Framework (DSMF) devoted to the design of component-based DSMLs (and generation of their editors)
- it could be possible to implement different DSMFs for different domains
- DSMFs could be metamodeled in the same way as DSMLs are metamodeled
- the DSMF implementation could rely on a common graphical framework.

The rest of the paper is structured as follows. Section 2 presents the motivations which led us to the design of a DSMF for building component-based DSML editors. Section 3 describes the design and implementation of this framework. Section 4 explains how we plan to generalize this approach. We conclude in Section 5.

2. Motivations

In this section, we present the motivations which led us to the design of a DSMF for building DSML editors.

2.1 Application needs

The research described in this paper takes place in the context of a project which aims at designing and implementing an autonomic administration middleware called TUNe [1]. An autonomic administration system allows the deployment, supervision and reconfiguration of software in a distributed environment. TUNe is a successor of the Jade [3] system which is an autonomic management system based on the Fractal component model [4]. Each software managed by Jade is encapsulated in a Fractal component which provides a generic interface allowing its management (essentially starting, stopping and configuring it). However, with Jade we observed that Fractal interfaces were too low level and difficult to use. The management behavior had to be implemented in Java and had to invoke Fractal's API. The main motivation in TUNe was to introduce higher abstraction formalisms (DSML) to hide the details of Fractal. To enable such an autonomic administration, we had to introduce several DSMLs:

 for node diagrams. We defined a DSML to describe hardware environments in which applications are deployed and administrated. This DSML allows a graphical description of the topology of the environment as a set of interconnected clusters (groups of machines sharing the same characteristics).

- for deployment diagram. We defined a DSML to describe the initial deployment of applications. This DSML allows a graphical description of a set of interconnected software elements, which forms a software architecture.
- for reconfiguration diagram. We defined a DSML to describe reconfigurations which are triggered by monitoring probes. Such a reconfiguration is defined as a state diagram where states are annotated with administrative actions.

Moreover in the TUNe project, we experienced that it could be very convenient to specialize the above DSMLs for particular application contexts:

- intensional diagrams. Grid applications generally require hundreds of servers and it is not practical to describe each of them in extension. Therefore we designed a language which allows describing a software architecture in an implicit way, we called it *description in intension*. The principle is to describe the application deployment pattern by using cardinalities as in a class diagram. Therefore, a described component or link can lead to the deployment of many (component or link) instances.
- application specific deployment diagrams. It is often interesting to issue a DSML specific to a given application as it captures constraints from this application. For instance, we designed such a DSML for clustered JEE applications [5] composed of web (e.g. Apache), application (e.g. tomcat or Jboss) and database (e.g. MySQL) tiers. This DSML only allows defining coherent JEE architectures (e.g. it is not possible to interconnect an Apache server with a MySQL server).

In summary, our experience in the TUNe project revealed the need for administration languages which are specific to the considered administration facets and specific to the application domains. If not all, many of these DSML have a graphical syntax and the implementation of editors associated with these DSML is the main issue addressed in this article. An important common denominator of the above DSML is that they all address the description of architectures composed of components, connectors, bindings and attributes.

Notice that in other previous research projects related to the definition of component models (such as Fractal [4] or SCA [6]), we also had to implement such graphical editors for defining component architectures, which emphasizes the need for tools that help implementing such editors.

2.2 Existing tools

In the MDE context, the *abstract* syntax of a DSML is defined by its *meta-model* and an instance expressed with this DSML is called a *model*. Each model must conform to the corresponding meta-model. A *concrete* syntax associated with a DSML provides users with a textual or graphical

formalism to manipulate the DSML's concepts. The abstract syntax (meta-model) and the concrete syntax are used to generate an editor for editing models which conform to the language definition.

The Eclipse Modeling Project (EMP) stands in this context. It is an open source project which is composed of several sub-projects centered on different MDE's aspects. Some of these sub-projects such as EMF [7] (Eclipse Modeling Framework), GEF [8] (Graphical Editing Framework) and GMF [9] (Graphical Modeling Framework) address DSML design and implementation.

EMF is a modeling framework for specifying meta-models and managing models. From the specification of the metamodel (.ecore), a set of Java classes are produced, where each element described in the meta-model is represented as a Java class. These Java classes can later be used as the foundation for developing tools. One interesting aspect of EMF is the possibility to generate an editor from the meta-model in order to manage instances (models). But this editor uses a Tree-viewer to display these instances, and typically doesn't provide an architectural view, thus lacking expressiveness.

GEF provides support for building graphical editors on top of the EMF framework. It allows the association of a graphical view with each element of the meta-model. This graphical view is implemented in the lower-level Draw2D framework which is a standard 2D drawing framework based on SWT from Eclipse. Designing a graphical view using GEF consist in writing Java code which can be a very painstaking work and this is why GMF came into existence.

GMF provides the ability to generate graphical representations of DSMLs, on top of EMF and based on the infrastructure provided by GEF. To implement a graphical editor for a DSML with GMF, developers need to construct a number of intermediate models which are used to generate the Eclipse plug-in which implements the editor. The domain model defines the abstract syntax of the language(.ecore meta-model). The graph model specifies the shapes that will be used in the editor. The tooling model specifies which tools will be available in the editor palette (e.g. a button to create an element in the language). Finally, the mapping model specifies how the shapes of the graph model are mapped with the concepts of the meta-model. GMF also provides a wizard which can derive first rough versions of the intermediate models (graphical, tooling and mapping) from the domain model. Then these generated intermediate models can be refined.

GMF is presented as the most complete framework for generating model's graphical editors. However, similarly to experiments reported in [10], our experiments revealed that GMF is quite complex, and not as user friendly as we had like it to be. Evaluations reported in [11] show that only 12% of users find GMF easy to use. The wizard which is responsible for generating the initial versions of the intermediate models works well only for very simple DSMLs, and when the user diverts from the wizard default solution, he has to hand-code what he really wants (with GEF).

Our analysis is that most framework for building such editors (e.g. EMF/GMF) are generic, i.e. aim at fulfilling the requirements of any field, which leads to increased complexity. Specializing such a framework according to the constraints of a domain (in our case component-based DSML) would allow keeping the definition of a specific editor simple, while fulfilling the requirements of the considered domain. Therefore, in the next section, we describe the design of a Domain Specific Modeling Framework which allows generation of DSML editors in the area of component-based languages.

3. A Domain Specific Modeling Framework

3.1 Design principles

We designed and implemented a modeling framework called yTUNe. This framework provides a meta-modeling environment for editing meta-models. From the edited metamodels, users can generate graphical editors for models edition. The yTUNe meta-modeling language basically uses MOF concepts and is dedicated to the design of componentbased languages. The aim of yTUNe is then to provide a user-friendly support allowing fully automatized generation of graphical editors, for modeling component-based architectures.

Regarding fully automatized generation of model editors, our idea consists in designing the yTUNe metamodel (which allows defining component-based DSML meta-models) so that a common concrete syntax can be applied to the targeted languages. It means that the yTUNe meta-model defines all the component-based concepts which shall be used in a DSML meta-model, and the common concrete syntax defines the types of graphical representations these concepts should have. The definition of the yTUNe meta-model is based on our knowledge of the type of concepts which are usually used in our component-based domain, and we also know the usually used graphical representations. However, the graphical representations can be configured/adapted when designing a particular DSML with yTUNe. Then the editor generation is achieved by simply interpreting the graphical elements associated with the concepts included in the DSML meta-model.

3.2 The yTUNe Modeling Framework

3.2.1 Meta-modeling environment

The meta-modeling environment provides specific component-based constructions which are:

- *Component*: It is a cohesive unit of attributes and properties, which represents the structuring unit that we call component in our domain. The graphical representation of a component is a shape which can be specialized in the yTUNe meta-modeling editor.
- *Connector*: This concept corresponds to a connection point (also called *interface* in the component domain). The graphical representation of a connector is a shape which can be specialized in yTUNe.
- *Binding*: In the component domain, components' connectors can be linked together to constitute a component architecture. The *Binding* concept represents these links. The graphical representation of a binding is a link (which can also be graphically specialized).
- *Aggregation*: an aggregation is a relationship which allows combining a component to a connector. Graphically it is characterized by the shape of the connector stuck with the shape of the component. This relationship always has a 1..n cardinality where only the n is settable. This cardinality is a constraint which enforces that in the generated editor a component can be stuck to n connectors but a connector can be stuck to only one component.
- *Association*: an association is a relationship with two ends, where one end is on a connector and the other end on a binding. A binding can only be associated with two connectors. If the user wants (in a meta-model) to associate more than two connectors with a binding, e.g. he wants in a model to allow a binding of type B which links connectors of type C1, C2 or C3, he can define in the metamodel an abstract connector C from which C1, C2 and C3 inherit and associate B with C.

The graphical representations described above are statically defined in yTUNe, but they can be specialized There are two *Components* properties which are concrete syntax elements:

- *Schematic property*: It is one of the graphical elements added in the meta-model definition. This property allows to assign a graphical representations to a component thanks to a drawing editor provided in yTUNe. The user can draw its desired representation or import an existing image.
- *Abstract property*: This graphical property allows to specify whether a component can be instantiated or not. According to its value, a creation tool will be available in the intended editor palette or not.

Besides these specific concepts, the meta-modeling environment also offers some elementary constructions provided by most of the existent meta-modeling environments: attributes, composition and inheritance relations.

We use cardinalities to restrain our languages, but to allow complete definition of DSL, we are working on integrating an OCL like constraint definition.

3.2.2 Editor generation

Once the DSL (meta-model) is specified and saved, generation of the models editor is achieved in only one click. In fact at this step the meta-model includes all the necessary information for generating the editor on the fly. The metamodel is interpreted and a model editor is generated according to the semantic assign to each concept defined in the meta-model (semantic described above) and their graphical representation.

The generated editors are enough user-friendly, and the edited models are conform to their corresponding metamodel definition. Next we illustrate the yTUNe editor generation process using two running examples.

3.3 illustrative examples

The two running examples we present here are representative enough of our specific requirements (components, binding, connectors, attributes ...).

a) J2EE: The first example is a configuration language for a J2EE platform. The J2EE platform defines a model for developing web applications in a multi-tiered architecture. Such applications are typically composed of a web server an application server and a database server. This use case consist in designing a configuration language for a J2EE platform, and then generating the corresponding J2EE graphical model editor. This configuration language designed with yTUNe is depicted in Figure 1. The web server is the *apache* component with its *mod_jk* connector, the application server is the *tomcat* component with its *ajp13* and *datasource* connectors and the Database server is the *mysql* server with its WAClients connector.



Fig. 1: J2EE Configuration language modeled with yTUNe



Fig. 2: Fractal ADL modeled with yTUNe

b) Fractal: The second example is a description language of the Fractal component model. Components have functional (client and server) and non-functional (control) interfaces, they are linked together thanks to *Links* to constitute a software architecture. The Fractal component model enables composite components which are components that contains other components. A basic meta-model of Fractal's ADL (Architecture Description Language) designed with yTUNe is depicted in Figure 2 and shows all the basic concepts of Fractal, which are: *components*, functional interfaces (*Client*, *Server*) and non-functional interfaces (*Control*).

Figure 3 shows the two editors generated from the above meta-modeled DSML.



Fig. 3: (a)J2EE editor. (b) Fractal editor.

Notice that in the meta-models described in Figures 1, all the components which describe servers has been assigned the corresponding server's icon as graphical representation. Then models edited in the J2EE generated editor (Figure 3(a)) are very expressive. Also the defined J2EE configuration language enforces that only a consistent j2EE architecture can be defined (e.g. it is not possible to interconnect an Apache component with a MySQL component).

Regarding the Fractal meta-model, Figure 2 shows that the *Interface* connector is aggregated with the *Component* component. In the generated Fractal model editor (Figure 3(b)) it is represented with *Interface* instances stuck on *Component*

instances.

Fully automatized generation of such editors has been made possible because we have specialized the meta-model of our modeling framework (for the area of componentbased languages) and associated a graphical concrete syntax with it (even though this syntax can be specialized for a given language). We believe that such an approach can be generalize in order to allow generation of different types of yTUNe frameworks according to different domain requirements. This vision is explained in the next section.

4. Generalization

4.1 Vision

We believe that this approach (implementing a DSMF for component-based DSMLs) can be generalized to issue different DSMFs associated with different fields. For instance, we are targeting a DSMF which could be devoted to the design of DSMLs for describing timing diagrams such as Gantt¹ charts. Therefore, in the same way as DSMLs are defined in terms of a DSML metamodel (from which an model editor is generated thanks to a DSMF), a DSMF can be defined in terms of a DSMF metamodel (from which an DSML metamodel editor is generated). Figure 4 illustrates this generalized approach. The generation of a



Fig. 4: Our vision

DSML metamodel editor from a DSMF metamodel cannot be fully automatized. At this step, the designer of a DSMF must associate with the DSMF metamodel a specification of the edition behaviour (a concrete syntax). For instance, the component-based DSMF metamodel defines Component, Connectors, Binding, Association and Aggregation concepts. The edition behaviour specification must define how these concepts are managed graphically. We are currently refactoring the yTUNe prototype to support both a component-based DSMF and a timing-based DSMF, while factorizing libraries and interfaces as much as possible. This should led us to the design of a framework for implementing DSMFs, which would allow combining a DSMF metamodel and reusable graphical libraries.

4.2 Positioning

This vision can be positioned compared to the Eclipse Modeling Project as follows.

From the point of view of the generation of the DSMF, we believe that this generation cannot be fully automatized, but that a implementation framework can help developing DSMF and promote software component reuse (libraries of graphical behaviors). This approach can be compared to GEF as the level of abstraction is quite similar. However, this comparison only holds from a narrow part of our proposal, as GEF does not target DSMF but DSML.

From the point of view of a particular DSMF which allows fully automatized generation of DSMLs' editors, this approach can be compared to GMF which aims at generating graphical editors with minimal development. However, since GMF is generic and does not target specific domains, it fails addressing the needs of these specific domains.

Concerning the related work, we did not found any similar proposal in the literature. The closest work to our is that of Nordstrom et al. [12] who propose to describe DSL meta-models with UML and OCL, and to associate presentation specification with such a meta-model to obtain a modeling environment (editor). However, this association of the graphical syntax with the concrete syntax is not detailed and will probably require development. This proposal is similar to our regarding the generation of DSMFs (which are modeling environments for DSMLs). We believe that this DSMF generation will require much expertise, which should be only required to generate DSMF, DSMFs being dedicated to a specific domain and much simpler to use for generating a DSML editor.

5. Conclusions

We are conducting research on a key issue in the MDE community which is the development of editors associated with DSMLs. Many experients showed that the standard tools in the domain (i.e. EMF/GEF/GMF) are overly complex and difficult to use. Our believe is that these framworks are too complex because they aim at being generic. Instead, we claim that we should apply the principles of DSML (Domain Specialization) to these frameworks. By doing so, DSML designers in a given field will be able to easily design DSMLs and generate the associated editors.

¹time-space diagrams designed for project management

Regarding this orientation, the goal of this paper was twofolds:

- to describe our experience in designing and implementing one domain specific modeling framework (DSMF) devoted to (i) the design of component-based DSMLs and (ii) the generation of dedicated editors.
- to present a vision of the ideal framework which would allow generation of DSMF.

References

- L. Broto, D. Hagimont, P. Stolf, N. Depalma, and S. Temate, "Autonomic management policy specification in Tune," in *Annual ACM Symposium on Applied Computing (SAC), Fortaleza, Ceara, Brazil.* ACM, Mars 2008.
- [2] B. Combemale, L. Broto, X. Crégut, M. Daydé, and D. Hagimont, "Autonomic management policy specification: From uml to dsml," in Proceedings of the 11th international conference on Model Driven Engineering Languages and Systems, ser. MoDELS '08, 2008.
- [3] S. Bouchenak, F. Boyer, S. Krakowiak, D. Hagimont, A. Mos, S. Jean-Bernard, N. de Palma, and V. Quema, "Architecture-based autonomous repair management: An application to j2ee clusters," in SRDS '05: Proceedings of the 24th IEEE Symposium on Reliable Distributed Systems. Washington, DC, USA: IEEE Computer Society, 2005, pp. 13–24.
- [4] E. Bruneton, T. Coupaye, M. Leclercq, V. Quéma, and J.-B. Stefani, "The FRACTAL component model and its support in java: Experiences with auto-adaptive and reconfigurable systems," *Softw. Pract. Exper.*, vol. 36, no. 11-12, pp. 1257–1284, 2006.
- [5] S. Microsystems, "Java 2 platform enterprise edition (j2ee)," 200x, http://java.sun.com/j2ee/.
- [6] OSOA, "Service component architecture (sca)," 200x, http://www. osoa.org.
- [7] D. Steinberg, F. Budinsky, M. Paternostro, and E. Merks, *EMF: Eclipse Modeling Framework*, 2nd ed. Boston, MA: Addison-Wesley, 2009.
- [8] T. Modica, E. Biermann, and C. Ermel, "An eclipse framework for rapid development of rich-featured gef editors based on emf models." in *GI Jahrestagung*, ser. LNI, S. Fischer, E. Maehle, and R. Reischuk, Eds., vol. 154. GI, 2009, pp. 2972–2985. [Online]. Available: http://dblp.uni-trier.de/db/conf/gi/gi2009.html#ModicaBE09
- [9] R. C. Gronback, Eclipse Modeling Project: A Domain-Specific Language (DSL) Toolkit. Upper Saddle River, NJ: Addison-Wesley, 2009.
- [10] A. Evans, M. A. Fernández, and P. Mohagheghi, "Experiences of developing a network modeling tool using the eclipse environment," in ECMDA-FA '09: Proceedings of the 5th European Conference on Model Driven Architecture - Foundations and Applications. Berlin, Heidelberg: Springer-Verlag, 2009, pp. 301–312.
- [11] V. Pelechano, M. Albert, J. Muñoz, and C. Cetina, "Building tools for model driven development. comparing microsoft dsl tools and eclipse modeling plug-ins," in DSDM, 2006.
- [12] G. G. Nordstrom, B. Dawant, D. M. Wilkes, and G. Karsai, "Metamodeling - rapid design and evolution of domain-specific modeling environments," in *Proceedings of the IEEE ECBS'99 Conference*, 1999, pp. 68–74.

Modeling and Analysis of NASA's Mission Software Development Archetypes using Petri-Net Graphs

Amanda Pavlicek and Tai-Chi Lee Department of Computer Science & Information Systems Saginaw Valley State University 7400 Bay Rd., University Center, MI 48710

Abstract:

The purpose of this research is to improve the validation and performance of mission safety software within the Mission Control room, as well as achieving financial objectives and fulfilling governmental regulation while utilizing the best software engineering and project management practices. The implementation of this project will be represented in an analysis paper that utilizes Petri nets to portray the enhanced relationships within the Mission Control software infrastructure and the underlying issues of the previous software infrastructure within Mission Control. The results of this research will contribute to our understanding of designing a more efficient software infrastructure. With this newly acquired knowledge, NASA organization, once the organization is given the opportunity to restart manned space missions, can push the barriers of space exploration and aeronautical science within the current market's reinforced limitations in future NASA space missions.

1. Introduction

The NASA space program has undergone a series of innovations for more than half of century since Russia launched Sputnik, the first satellite to orbit the Earth in 1957. Through the years, NASA has made great

accomplishments in space exploration using advanced Mission Control software. While some mission software systems are successful in their objectives, there are other systems that need to be improved upon so the risk of failure during a NASA mission can decrease. In addition, the world economy's unpredictable state should be considered in determining how and where the NASA software team obtains or creates Mission Control and vehicle software.

Mission software is based on bettering all phases of a NASA space mission, in which examples of such are ground and flight data systems and on-orbit performance management. This category of software, which is composed of workflow and data-processing applications, is accountable for the scientific progress of a mission and the safety of NASA astronauts, space shuttle, and other aeronautical equipment. Research will revolve around not only NASA's current state as a governmental agency, but also the installation and performance of NASA mission software systems and applications within the Mission Control Room prototype, individually as independent applications, and collectively as a functional, vital component of a NASA mission.

2. Modifications in Software Methodologies

Norman Kluksdahl, a Systems Engineer at Johnson Space Center, noted that the most monumental change with NASA's Mission

Control software infrastructure and enhancement is their software development methodology, the work-related outline executed to plan and implement software applications, in creating software for both space vehicles and the Mission Control consoles [6]. All areas of the software development process, including programming, management, and application users, will be affected by this prevalent modification as it allows the development team to produce validated, well-received software before the deadline and with minimal cost.



Figure 1. A pictorial representation of the Waterfall Development Methodology that the Mission Control Systems team relied on for their software-engineering framework in employing Mission Control software (The Smart Method). [12]

In the past, the Mission Control Systems team followed a very straightforward process that is known as the Waterfall development, a sequential software-engineering process [10]. According to my research in project management, the programming team converses with management and other stakeholders to obtain details about the application being developed. After the application's first version is fully advanced, an abundance of Quality Assurance and Quality Control procedures occur in the application development premission phase known as the Baseline Release Process [6]. Succession with module tests allows the Simulation team to proceed with integration testing, which provides a set of tools to test subsystems [6]. The "Test Rig" or "Simulation" practices occurs when changes suggested by programming team are made, which gives the team an opportunity to watch the system and observe how it interacts with environment [6]. To be truly efficient, it is necessary to conduct lifecycle analysis upon all pieces of software to ensure that each piece doesn't have tight interdependency within one another. The team loads the new software into the infrastructure after the selected applications pass the validation testing, but the implementation process varies depends on how long simulation and runtime should take [6]

The greatest flaw with this waterfall process, according to Kluksdahl, is that the progress of the project is compromised if a stakeholder wishes for a change to be made in the new application [6] [10]. Therefore, all of the completed work, whether it was a sliver of a program functionality to an entirely new version, frequently has to be scrapped and the developers have to start all over again. There is no limit to how many times this could occur, and this major disadvantage has misused precious time, money, and solid relations within the NASA organization. Unfortunately, a stakeholder's influence on the project's final result will always remain a constant in any substantial endeavor because not everyone will know what they exactly desire, especially if it revolves around a concept that was never physically created before [10]. NASA, along with other corporations, have witnessed what the waterfall development cycle can do for an software development-driven project, where 70 percent of objectives in an average project are not met [11]. Furthermore, the waterfall development approach has demonstrated that it is not the most fiscally-conservative methodology as project costs are, on average, much higher than the estimated costs approximated during project initiation [11].

In addition, space operations and management officials rarely observe concrete, effective software development in their project preparation stage or incorporate room for flexibility within the finished product for future software-engineering trends [6]. Incorporating future trends within obtaining the project's end objectives allows advancing technology and strategies to take root. However, today's software development tends to revolve around the latest programming fad, such as object-oriented design, as a cure-all for efficient development. Often, the effect is to reduce productivity because engineers are required to learn a new methodology. As exemplified in several information technology-related projects, industry experience shows that the first use of a new methodology is often very inefficient, and thus more costly, than the methodology it replaced. This causes issues with maintainability and adaptability.



Figure 2. A pictorial representation of the Agile Software Development Cycle that Dr. Kluksdahl and other software system engineers are incorporating into their work methodology to ensure stress-free software development. [1]

What the Mission Control System group needs to build a better project framework for NASA developers and management to follow is one that embraces compliance in the sense of team communication, project objectives, and resources. During this time of mission inactivity, Kluksdahl expressed that this lack of productivity gives the team an opportunity, enough time and resources to "decide on a methodology and adapt to its standard practices"[6]. The Mission Control System group selected the Agile Software Methodology to integrate into their project objectives. Created by seventeen prominent software developers within a document called the "Agile Manifesto" in 2001, the Agile Software Methodology embraces adaption, cooperation, swift delivery, and client goal-orientation as the crucial elements of successful software development [3]. As stated in the Manifesto, programmers that function under an Agile programming settings value "individuals and interactions over processes and tools, working software over comprehensive documentation, customer collaboration over contract negotiation [and] responding to change over following a plan", which allows teams to actually provide a product on time and within the flexible requirements set by the client [3].

The Agile Methodology thrives on four basic stages of software development: Project Requirements, Software Design, Development or Implementation, and Testing and Debugging, the four stages that are present in the waterfall model [2]. One of the most notable discrepancies between the agile and waterfall development approaches is that when project requirements change, resources are modified, and other factors previously analyzed confront a major alteration, the entire project doesn't start over. Instead, the project leaders and key stakeholders integrate the changes into the project plan accordingly as the developers interact with each phase of the cycle until the current issue is resolved [2]. How various developers and other stakeholders are able accomplish distinctive tasks from different stages of the software development process is through "iterations", which allows software developers to complete a full software cycle for numerous tasks within a short amount of time [2]. A fully operative software development team is engaged in the planning, development, testing, and execution included in every iteration, where a deliverable or a segment of the finalized deliverable is produced [2]. Each project member must be heavily involved in selected Agile practices, such as face-to-face conversations, intricate team meetings involving daily reports, and comprehensive documentation [2]. In addition, users and developers alike must review and approve changes before these project modifications are merged into the baseline, developers must review each other's code and programming skill sets, and each piece of code must undergo unit tests. Customer representatives are introduced into the project team to act upon the stakeholder's behalf, allowing the project team sufficient access to the client's needs. The iteration approach contrasts with the waterfall development approach, which focuses on employing distinct phases with checkpoints and deliverables at every stage, a nonflexible, risky method of project planning [2] [11]. Kluksdahl described that through these phases allow developers and project leaders to acquire the knowledge necessary to fully implement an unfamiliar system, in which the waterfall methodology greatly lacks, by completing easier tasks in the earlier phases and focusing on the challenging assignments later [6] [11]. Also, verification of project requirements occurs much earlier in the development process than it would with the waterfall development approach, which permits stakeholders to adjust requirements while they are still relatively painless to modify [11].

In his discussion concerning the transferring of software methodologies, Kluksdahl stressed that the Mission Control Systems team is only restructuring the way the team operates within a software development environment and recycling the indispensable components (software development phases, requirements, communication means, etc) to fit the Agile structure [6]. Currently, Kluksdahl and other NASA system engineers are fundamentally pushing a development system into operation that does not employ all of the most modern software practices or is without any of the elements of the preceding software methodology. According to Kluksdahl, "What we are doing is we are making software development better, more sufficient for everyone involved in future continuing, and complex projects. Everything still works, so why should we get rid of it" [6]. The factors that needed to be improved upon, such as flexibility, team collaboration, and project integrity were the only portions that needed to be emphasized to build an advanced, enduring development methodology [6] [11].

3. Mathematical Approach of Petri Net Models

A Petri net is exemplified by these mathematical structures provided by James Peterson, in his book: *Petri Net Theory and the Modeling of Systems* [9]:

DEFINITION 3.1: A Petri net structure *C*, is a four-tuple, C = {P, T, I, O}, where P = {P1, P2...Pn) is a finite set of places, $n \ge 0$. T = {t1, t2... tm) is a finite set of transitions, m ≥ 0 . The set of places and the set of transitions are disjoint, P \cap T = Ø. I: T \rightarrow P^{∞} is the input function, a mapping from transitions to bags of places. O: T \rightarrow P^{∞} is the output function, a mapping from transitions to bags of places.

The cardinality of the set P is n, and the cardinality of the set T is m. We denote an arbitrary element of P by pi, i = 1...n, and an arbitrary element of T by tj, j = 1...m. I(ti) is a bag of input places for the transition ti.

O(ti) is a bag of output places for the transition ti [9].



Figure 3. An example of a Petri net graph with four places, two enabled transitions, and four markings, such that $I(p1) = \{t2\}$, $O(p1)=\{t1\}$, $I(p2)=\{t1\}$, $O(p2)=\{t2\}$, $I(p3)=\{t1\}$, $O(p3)=\{t2\}$, $I(p4) = \{t2\}$, and $O(p4)=\{\}.[4]$

This research endeavor will specifically concentrate on a branch of the Petri net theory that is called Applied Petri net Theory, where a Petri net is represented by special graphs called Petri net graphs [9]. Petri net graphs are also described in mathematical terms below, establishing Petri net graphs as bipartite directed multigraphs, where the multigraph allows various arcs or functions from one node to another [9]:

DEFINTION 3.2: A Petri net graph *G* is a bipartite directed multigraph, G=(V, A), where V ={v1, v2, ..., v5} is a set of vertices and A = {a1, a2, ..., a,} is a bag of directed arcs, ai = (vj, vk), with Vj, Vk is a subset of V. The set V can be partitioned into two disjoing sets P and T such that V = P U T, P \cap T = Ø, and for each directed arc, a is a subset of A, if ai = (vj, vk), then either Vj is a subset of P and vk is a subset of T or Vj is a subset of T and vk is a subset of P [9].

A Petri Net is composed of five elements: a set of places (represented by the letter 'P' in formulas and by circles in graphical representations), a set of transitions (represented by the letter 'T' and by bars in graphical representations), an *input* function (represented by the letter 'I' and by incoming arrows or arcs in graphical representations), an output function (represented by the letter 'O' and by outgoing arrows or arcs in graphical representations), and markings located inside place (are represented by the marking *u* and by dots in graphical representation) [8] [13]. The input and output functions creates relations between transitions, which controls the arrival and departure of resources, and places, which is a component that holds resources [8] [13]. A mapping from a specific transition (tj) to a specific place (pi), within a collection of input places I(tj) is an input function I[9] [13]. If transition (tj) is mapped to a specific place (po), within a collection of output places, O(tj), then the result is an output function [8] [9].

Markings are an essential component of a Petri net graph because they represent resources used to simulate the overall functionality of the replicated system. Described below is the mathematical description of a marking:

DEFINITION 3.3: A marking μ of a Petri net $C = \{P, T, I, O\}$ is a function from the set of places *P* to the nonnegative integers *N*. μ : $P \rightarrow N$.

The marking μ can also be defined as an n-vector, $\mu = (\mu 1, \mu 2, \dots, \mu n)$, where n = |P| and each μ i is a subset *N*, i =1, ..., n. The vector μ gives for each place pj in a Petri net the number of markings in that place. The number of markings in place pi is μ i, i = 1...n. The definitions of a marking as a function and as a vector are obviously reflected by $\mu(pi) = \mu i$. The functional notation is somewhat more general and so is more commonly used [9] [13].

In other words, markings are used to define the execution of a Petri net, where the number and positions of markings may alter during the execution and controls the execution's duration [9] [13]. In order for execution to occur within a Petri net graph, a transition must first be enabled for firing, which requires each of the transition's input places to have at least as many markings in it as arcs from the places to the transition [8] [9]. Next, the transition is fired by removing all of its enabling markings from its input places and then depositing into each of its input places and then depositing into each of its output places one marking for each outgoing arc from the transition to the place. Thus, multiple markings are produced for multiple output arcs and transition firings can continue as long as there is at least one transition enabled. Otherwise, the execution is terminated and the system is presumed to be idle [8] [9].

The most straightforward view of a Petri net representation of a particular system focuses on two concepts: events and conditions. Actions that occur within the system are events, in which the instigation of an event is influenced by the current state of the system. The varying state of the system



Figure 4. A Petri-net graph that demonstrates how the NASA Mission Control Systems team applied the basic outline of Waterfall Methodology in their software development process.

is derived by a set of predetermined conditions, "predicate[s] or logical description[s]" to the system's status [9] [13]. In order for an event to occur, the event's preconditions must be completed. The repercussions event. post of an or conditions, may cause the preconditions of other events to end or commence [9] [13]. Conditions are modeled by places in a Petri net while events are represented by transitions [9] [13]. The inputs of a transition are considered preconditions of

the consequent event, where an enabled transition indicates that the preconditions are completed [9] [13]. Thus, the outputs of a transition are the post conditions of an event [9] [13].

4. Petri Net Graphical Comparisons of Software Methodologies

Four Petri net graphs were created using an open-source application that creates Petri net graphs called Pipe 2. The resulting Petri net graphs fully resemble how the NASA Mission Control Systems team conducted and will conduct their software development utilizing the precedent approach: Waterfall Development Methodology and the impending approach: Agile Software Development. Figure 5 displays the progression of an archetypal project where the project team must endure the requirements phase, design phase, implementation phase, verification phase, and maintenance phase (which is an ongoing process until the software is eliminated from the system). So, when the project starts, the resources of the first version (which is symbolized by the marking) will undergo the requirements phase and when it is completed, a transition will be enabled, fired, and the marking will proceed to the design phase. This will continue until the software is successfully released.

Nevertheless, when project requirements change or a major project issue arises within the software development or the deliverable is being produced, another marking will be added to the "Project Hindrance" place, as exemplified in Figure 4. This particular graph is built in a manner in which if there are two markers next to each other in separate places (which represents the encounter or realization of a significant issue), a specific transition will be enabled (representing the event in which the team realizes that there is an issue at hand), firing one of the markings to the beginning of the project process and firing the other to one of the "Resources Cache" place. The project team must repeat the project process again with new resources, a revised project-related blue print, while the dilapidated resources and wasted time will be accumulated "Resource Cache" places until the deliverable is released into the system. Evidently, this Petri net graph (Figure 5) illustrates how the Waterfall development methodology exhibits an uncomplicated organization, but this limits the amount of resources and adaptability of the project team. The following Petri net graph (Figure 5) epitomizes the Mission Control System team's current development process, which employs Agile Software Development practices and six delivery iterations. The marking (current state of the specific project) begins at Iteration 1 at the beginning of the project and the project team undergoes its own software development cycle (represented in Figure 7) to manufacture a deliverable. When that iteration's particular deliverable is generated, the precondition of the next event is fulfilled and the nearest transition in Figure 6 (T0) will become enabled, transporting the ticket from the preceding iteration to the subsequent iteration, which is called Iteration 2. This process continues until the marking reaches the maintenance place, where the marking will always be repositioned by a transition that will permanently be enabled (T7).



Figure 5. The constructed Petri Net that represents how the Waterfall Methodology was used within the NASA software development environment and the multiple problems that the methodology possessed. The Project Hindrance place represents any kind of condition or problem which may delay the project, such as management disagreement, misallocated resources, and drastic change in project requirements. When a resource (project resource) is found in a Resource Cache place, it is considered a misplaced resource (a condition to an imperfect project) that increases the final cost of the project. In addition, if there is more than one marking within the Maintenance place, then the number of markings represents the number of versions in which the delivery has been through.



Figure 6. A Petri net graph that represents how Kluksdahl's Mission Control System team is going to integrate the basic concept of Agile Development Methodology. Each place represents an iteration or project phase and each marking represents allocated project resources. If an iteration possesses a marking after a transition is fired, then it resembles the current state of the specific project. If a problem occurs in an iteration, it is taken care of within that iteration right away.

The Mission Control System team must consider if any problematic difficulties will arise during the development of an iteration's deliverable. Within the world of Petri nets and graphical mathematical models, there is no way to logically relocate the marking to the first place of the Petri net due to



Figure 7. This Petri net graph represents the Agile software development cycle that is assigned to each project iteration. This closely resembles the flowchart example found in Figure 7 where the present location of the marking represents the current phase in which the project team is undergoing. However, the team is able to move to any other phase if need be.

the fact that the advancement of the project is almost linear, from beginning to end. If an issue occurs that would characteristically cause the project team to retrace back to Iteration 1, or the establishment of the project, the issue will be resolved within that current iteration before the next transition enables. Within the iteration, the project team can move to any phase of the software development cycle to disentangle the issue (whether it be using more face-to-face communication with clients, or revising the requirements of the iterations, or conduct more conscientious coding during the implementation phase) at hand. After the issue is resolved, the project team can go forth in completing the originally-planned objectives within the iteration and the team can make an endeavor to prevent a similar issue from happening again.

5. Conclusion

The Petri net investigations of both software development methodologies revealed that the Mission Control System team is making significant, beneficial progress in renovating the Mission Control Room by establishing Agile Software Development characteristics into their own infrastructure. Through the project's deliveries, six delivery phases in span of three planned years, the team expects to repeat each delivery process (update, design, work from high-level tasks to specific tasks) in order to ensure the best results. The deliveries and project's development is undergoing a "loading bar" process, or an incrementing progression towards a project's completion [6]. This allows system analysts and other researchers to measure a project's development and their ability to learn how to operate the technology using key phases. It permits the project team to approach the major issue concerning the Mission Control software infrastructure through a divide and conquer approach: learn the easy tasks first and put the most difficult and/or incompletely-specified tasks on the side to be completed later in the project [6]. Although the software production, resource verification hardware acquisition, and other tasks are still under construction, I believe that the verification of how the Agile Methodology attributes will be integrated into the infrastructure propelled the Mission Control Systems team a great length towards their ambition: the project's final delivery of a resourceful, powerful Mission Control Room.

As expressed before by Dr. Kluksdahl, along with several members of the Mission Control Systems team, experience is their main advantage in this project, a statement that I agree with whole-heartedly [6]. Gained knowledge gives the Mission Control Systems team, along with other stakeholders, the best results under an optimal testing environment that is facilitated by the combination of multiple project and development methodologies [6]. I found this component intriguing since knowledge gained from the previously-completed phases helps provide operators and testers sufficient knowledge to evaluate past error recognition and trends. When a situation is fully corrected based on the team's strict standards, the confidence of the team and crew in achieving the flight mission increases significantly, which allows the team to work harder towards a successful future. The more each team member grows mentally and personally, the more their own motivation will develop and they would finally realize the importance of their work to the entire organization, and, perhaps, the entire human race. Who would compromise the promise of a better future and the opportunity of obtaining knowledge in order to maintain the archaic environment that has failed NASA before. If the team, along with management and all other departments, tenaciously continues this endeavor in the determined "loading bar" pace and able to create resourcesaving solution, NASA will be triumphant once again.

6. Acknowledgements

I would like to express my gratitude to Dr. Norman Kluksdahl, whose contributions and constant encouragement inspired me to push my own boundaries in this project and understand the significance of not only NASA's role in society, but also my contributions to the academic world. I would also like to thank the rest of the Mission Control Systems team, and my research mentor, Dr. Tai-Chi Lee, who provided extremely valuable input concerning the conducted research and my own conclusions.

References

[1] Agile Methodology,

http://karthiksangi.files.wordpress.com/2009/07/agile3.gif [2] "Agile Programming", http://agileprogramming.org/, 2012

[3] K. Beck et al., Agile Manifesto.

http://agilemanifesto.org/, 2001

[4] Detailed Petri-Net.

http://en.wikipedia.org/wiki/File:Detailed_petri_net.png

[5] T. England, "NASA's New Future." University of Michigan, Michigan Space Grant Consortium Conference, Ann Arbor, MI, November 12, 2011.

[6] N. Kluksdahl, Systems Engineer at Johnson Space Center.

[7] N. Heath, (2010). "Space Exploration: The Computers That Power Man's Conquest of the Stars." *Silicon.com*. http://www.silicon.com/management/publicsector/2010/09/25/space-exploration-the-computers-that-

power-mans-conquest-of-the-stars-39746245/print/

[8] T.C. Lee, "Some Properties of Petri Nets and Their Graph Models," Proc. 27th ACM Annu. Conference of Southeast Region, April 1989, pp. 657-659. Interest

[9] J.L. Peterson, "*Petri Net Theory and the Modeling of Systems*," Englewood Cliffs, N.J.: Prentice-Hall, 1981.

[10] W.W. Royce, "Managing the Development of Large Software Systems," Proc. IEEE Wescon, 1970.

[11] Serena (2007). "*Agile in the Enterprise*" http://www.serena.com/docs/repository/products/teamtrack/a gile-in-the-enterprise.pdf

[12] Waterfall Model, http://www.learnaccessvba.com

[13] M. Zhou and K. Venkatesh, "Modeling, Simulation, and Control of Flexible Manufacturing Systems—A Petri Net Approach," World Scientific, vol 6, 1999.

A new trigonometric Method for automatic visualization of metro map layout

Somayeh Sobatimoghadam¹, Ahmad Absetan²

¹Departement of electronic and computer engineering, Hakim Sabzevari

University, Sabzevar, Iran, sobati58@yahoo.com

²Department of electronic and computer engineering, Hakim Sabzevari

University, aabsetan@yahoo.com

Abstract In this paper we describe an algorithm to automatically generate metro map layout. This automatic method use to move from the initial geographic layout of the map to a schematic layout. We use trigonometric relationship in five steps. The proposed method observes metro map criteria to produce a nice metro map. We describe the steps of algorithm on real world metro.

Keywords: public transport schematics, graph drawing, metro map layout Automation.

I. INTRODUCTION

The Metro map is a type of diagram that is used for illustrating transportation networks. Metro map visualize the public transport and illustrates interconnections of rail road networks. Today it is familiar to people because it is easy to read and understanding and many people are able to use it quickly. It is easy for travelers to find their routes from start to destination point, and they know how to read route maps for metros and buses [1]. Schematic map is used to represented underground metros or above ground tram networks. Designing a schematic map today is still a challenge, because it is created manually. Traditionally, metro maps are drawn manually [2] and there was no automatic way to produce this diagram; the cartographer must decide where to put the stations and how to draw the lines in the diagram. In complex and large networks it is not simple and takes many efforts and time consuming. Therefore designing a schematic map today is still a challenge. The main objective of this paper is provide solutions for metro map layout problem. The metro map can be presented as a graph; it consists of a set of lines which have intersections or overlaps. In all metro maps it is common to consider the graph in which, metro stations are considered as vertices and their interconnections as edges. The lines are straightened and restricted to horizontals, verticals and diagonals at 45[°]. For forming a good metro map layout, a set of

1

aesthetic criteria should be defined. We describe a method to automatically generate metro maps quickly. Our method uses some simple trigonometric relationships with a suitable preprocessing step layout and observes a set of aesthetic criteria for good metro map, The experiment results show that our method produces good metro map layouts.

II. RELATED WORK

In this direction there are some related work. In this section we review these methods and we compare them. Nollenberg and Wolff [3] have investigated an algorithm for automatically drawing metro maps. A list of soft and hard constraints has specified. They have presented a mixed- integer linear program (MIP) which finds a drawing that fulfills all hard constraints. This program optimizes a weighted sum of costs corresponding to the soft constraints. They have shown how to include vertices labels in the drawing [3]. Stott and Rodger [1] have proposed an automatic method for drawing metro maps. They have considered the metro map as a graph. The graph is embedded on an integer square grid. They have implemented a total of eight criteria. Each criterion measures some geometric property of the map such as the length of edges or edges crossings and is weighted. The nodes and labels are repositioned such that the total of the weighted criteria is always reduced. F. Bertault [4] has proposed an algorithm based on a force-directed approach. They considered three kinds of forces: the attraction forces between nodes, the repulsion forces between each pair of nodes, and the repulsion forces between nodes and edges [4]. Their algorithm produces drawings which are more aesthetically pleasing that the initial layout and the final layout of a given graph has the same embedding as the original graph [2].Frick and Ludwig [5] have presented an algorithm to compute a layout for undirected graphs based on local temperature. They have calculated a local temperature for each vertices and a global temperature that is the average of the local

temperatures. The vertices positions have updated until the global temperature is lower than a minimal temperature. This algorithm can handle large and complex graphs. Hong and Merrick [2] have investigated the new problem of automatic metro map layout. They have defined a set of aesthetic criteria and presented a method to produce the lavout automatically. They have combined several different graph layout methods. They have designed layout methods based on the first four these criteria. They have investigated five different methods.In comparison Hong and Merrick have reduced the running time, but labels overlaps are avoided, and sometimes labels and edges intersect. The result of their methods is not very similar to the maps that are drawn by graphic designers. The edge lengths are not uniform; which means the final layout is unbalanced. Their approach is severely limited by not taking the topology of the metro systems into account. The running time of Stott and Rodgers, is higher in comparison to Hong et al. The edge lengths are uniform and edges are octilinear. But there exist many label overlaps. The Nollenberg and Wolff method draws octilinear metro map but it is complex to use and the length of edges are not uniform.

III. MODELING METRO MAP LAYOUT

For illustrate the conception of metro map we define existing criteria for a good metro map layout and then we modeled the metro map as a graph.

A.GRAPHS

Graphs are important because they can be used to represent essentially *any* relationship between entities. Graphs visualize the information for the users, and provide important information about the objects. For example, graphs can model a network of roads, with cities as vertices and roads between cities as edges. A nice layout of graph aids user to find immediately the information that he is looking for.

A graph G=(V,E), consists of a set of vertices V, |V|=n, and a set of edges E, |E|=m. An edge is a pair e=(u, v), $u, v \in V$. Two adjacent vertices are connected by an edge. The degree of a vertices is the number of edges that incident to this vertices A path in G is a sequence of distinct vertices of G like (v1, v2, ..., vt) such that $vivi+l \in E$ for 1 < i < t.

A graph is called planar if it can be drawn in the plan without edge crossing.

B.Graph layout

Graph drawing applies topology and geometry to draw two and three dimensional representation of graphs. Very different layouts can correspond to the same graph. There are different graph layout strategies. In straight line drawing each edge is drawn as the straight line between the vertices. Orthogonal layout the edges are drawn as polygonal chains of horizontal and vertical line segments. We will use the octilinear layout for metro map drawing. In octilinear layout all edges are horizontally, vertically or 45° . The advantage of using this layout in comparison of orthogonal layout is that the maximum possible vertices degree increases from 4 to 8. Some drawings are better than the others. Aesthetic criteria attempt to characterize readability of a layout. Various attempts have been made to specify the readability of a layout that comprises:

- minimize crossings
- minimize area
- minimize bends (in orthogonal drawings)
- minimize slopes (in polyline drawings)
- maximize smallest angle
- maximize display of symmetries.

In the other hand there are graphical properties for readability of a layout. For example the shape and size and color of vertices; the shape (curve, line), thickness and color of edges. In any algorithmic method of graph drawing these properties are left to the decision of user. Another property that must be taken into account is the graph labeling. The size and place of labels are important to minimize overlaps.

C.The metro map problem

The metro map can be modeled as a graph, formally.

Input: The metro map graph in which stations appear as vertices and and relations between two station appears as an edge.

Output: a suitable metro map graph layout in which observed following criteria.

D. metro map criteria

First, we have to know what kind of map is suitable for metro map and what is a "nice" metro map? The map should be as readable and clear as possible without displaying unnecessary details. metro map graph must be include following criteria:

- 1- Each line drawn as straight as possible.
- 2- No edge crossings.
- 3-No overlapping of labels.
- 4- Lines mostly drawn horizontally or vertically, with some at 45° .
- 5- Each line drawn with unique color.

We discover a method based on these criteria, in next sections we describe our method .

V.TRIGONOMETRIC METHOD

We implemented trigonometric relationships to acquire metro map layout. Each step is continuance of previous step. To illustrate the performance of our algorithm, we used it performs on Washington metro map, Initial layout of the Washington Metro Map is shown in fig.1. Our algorithm consists of five steps:



Fig. 1. Initial layout of the Washington Metro Map(G).

Step 1: Minimizing the graph by eliminate the nodes that their degrees are equal two.

Step 2: Planing the graph if it is not planar.

Step 3: Producing octilinear layout using trigonometric relationships.

Step4: Inserting the removed nodes.

Step 5: Labeling

Step 1 : Minimizing graph:

The initial graph G is shown in fig. 1, and for applying next steps easily, the vertices of gegree two is eliminated. The result G' is shown in fig. 2.

Step 2: Smoothing The Graph G'

We considered that the grapg is plannar. If it is not planar, it is plannared using following formulas . New coordinations of vertices is calculated by following formulas:

$$x(v) = \frac{1}{deg(v)} \sum_{w \text{ adj } v} x(w) \qquad (1)$$

$$y(v) = \frac{1}{deg(v)} \sum_{\substack{w \text{ adj to } v}} y(w) \quad (2)$$

Where v and w are two vertices that are adjacent and deg(v) is of degrees of node v.

step 3:octilinear layout

In this step the angle between each egde and horizontal axis is calculated by formula (3).In other word the angle between the edge that is between each pair of adjacent vertices and horizontal axis.

Where ang(v,u) is angle between each pair of adjacent vertices and horizontal axis, x(v) and x(u) are x-components of two adjacent vertices and y(v) and y(u) are y-components too. In this step the deference between the ang(v,u) and three degrees, 0^0 , 90^0 and 45^0 is calculated .A angle θ is selected which have minimum deference by ang(v,u). The new coordinates of node u is yielded by (4) and (5).

$$x(u) = x(v) + r\cos(sign(x) \times \theta)$$
(4)

$$y(u) = y(v) + r\sin(sign(y) \times \theta)$$
(5)

Where *r* is the distance between *v* and *u* vertices, sign(x) is the sign of (x(u)-x(v)) and sign(y) is the sign of (y(u)-y(v)). Our method in this step is traversed all vertices of the graph as BFS (Breadth First Search) and is coordinated each node.





Fig. 3. G", octilinear layout of graph G'

Step 4 : Inserting removed nodes

In this step the nodes that was removed in inserted to graph. For each edge (u, v) if n is number of nodes that were deleted and θ the angle between this edge and horizontal axis :

1-if *θ*=0

$$\sum_{k=1}^{n} x(k) = x(v) + \frac{k \times r}{n} \cos \theta \quad , \quad y(k) = y(u)$$
 (6)

2- if
$$\theta = 90$$

$$\sum_{k=1}^{n} x(k) = x(v) , \quad y(k) = y(u) + \frac{k \times r}{n} \sin \theta$$
(7)
3aifg θ = $\tan^{-1} \frac{|y(v) - y(u)|}{|x(v) - x(u)|}$ (3)



Fig. 5. The final result of Stott and Rodger method [1].

The metro map layout of Stott and Rodger [1] is shown in fig. 5.

5- labelling

In last step the lables is inserted. For each node 8 position is considered for labling as is shown in fig 8.



Fig. 6. The position of lables.

The position of lables is considered according to the angle beetwin inserted nodes and edge (u,v), which u and v are beginning and end of edge. For 0 and 180 degree the position 2 and for the other, the position 3 is recommended.

VI. CONCLUSION AND FUTUR WORKS

In this paper we presented an algorithm using trigonometric relationship to automatically generate good layouts of metro maps. The results show that final layouts of the metro maps have significant improvement on the original geographic layout and they are similar to official maps. We believe that our method for labeling can be extended relatively easily in order improve further on the quality of the map labling.



Fig. 7. The final layout with labelling

REFERENCES

[1] J. M. Stott and P. Rodgers. Metro Map Layout Using Multicriteria Optimization, Proc. of International Conference on Information Visualisation (IV04), 2004, pp. 355-362.

[2] S. Hong, D. Merrick. Automatic Visualisation of Metro Maps. Journal of Visual Languages & Computing, 2006, PP. 203-224.

[3]M. Nollenburg and A. Wolff. A Mixed Integer Program for Drawing High Quality Metro Maps. In Proc 13th International Symposium on Graph Drawing. Verlag Berlin Heidelberg, Springer, 2005.

[4] F. Bertault. A Force-Directed Algorithm that Preserves Edge Crossing Properties, Graph Drawing99, LNCS 1731, 1999,pp. 351-358.

[5] A. Frick, A. Ludwig and H. Mehldau. A Fast Adaptive Layout Algorithm for Undirected Graphs, Proc. of Graph Drawing 94, LNCS 894, 1995, pp. 388-403.

[6] R. Davidson and D. Harel. Drawing Graphs Nicely Using Simulated Annealing. Technical ReportCS89-13, Department of Applied Mathematics and Computer Science, The Weizmann Institute of Science, Rehovot, Israel, ACM, 1993.

[7] H. de Frayssex, J. Pach, and R. Pollack. How to Draw a Planar Graph on a Grid. Combinatorica,1990, pp. 41-5.

[8] T. M. J. Frachterman and E. M. Reingold. Graph Drawing by Force-Directed Placement. Software Practice and Experience, 1991, pp. 1129-1164.

[9] S. Hong, D. Merrick, H. A. Nsacimento. The Metro Map Layout Problem. Australasian Symposium on Information Visualisation, Christchurch, 2004.

[10] W. He and K. Marriott. Constrained graph layout. An international journal,1998,pp.289-314.

[11] JavaScript Vector Graphics Library, http://www.walterzorn.com/jsgraphics/jsgraphics _e.htm.

[12] S. Kirkpatrick, C. D. Gelatt and M. P. Vecchi, ,'Optimization by simulated annealing', Science, 1983, pp. 671–680.

[13] R. Laurini. Information Systems for Urban Planning. A Hypermedia Cooperative Approach, Taylor and Francis, 2001.

[14] R. Laurini, D. Thompson, Fundamentals of Spatial Information Systems. Academic Press, 1993.

[15] P. Lacomme, C. Prints, M. Sevaux, Algorithmes de Graphes. Ed. Eyrolles, 2003 .

[16] T. Masui, Evolutionary Learning of Graph Layout Constraints from Examples. ACM, November 2-4, 1994, pp. 103-108.

[17] D. Merrick and J. Gudmundsson. Path simplification for metro map layout. In Pro Graph drawing, 2007, pp. 258-269.

[18] R. Tamassia. On Embedding a Graph in the Grid with the Minimum Number of Bends. SIAM Journal of Computing. Vol. 16. No. 3. 1987. pp. 421-444.

BUSINESS OPTIMIZATION THROUGH A PORT E-COMMERCE

Georgeta Şoavă¹, Mircea Alexandru Răduțeanu² ¹Faculty of Economics and Business Administration, University of Craiova, Romania ²Faculty of Economics and Business Administration, University of Craiova, Romania

Abstract:

Today's companies are faced with the need to exploit technology changing computer environments, in order to improve customer satisfaction and reduce costs. A successful approach to electronic portals is an effective demonstration of new ways of relating to the client. Thus, we conducted this work, first trying to emphasize the need to develop appropriate software for e-business, creating a portal for electronic commerce and finally to introduce it, and that to highlight the main features that arose logically from the objectives set. We completed work showing some of the advantages of implementing Electronic Commerce Portal applications.

Keywords: e-commerce, e-business, e-Portal

1 Introduction

From a business perspective, e-business offers companies a way, to grow their businesses in an environment where technology meets user demand collaboration applications developed for the Internet. Thus, by using e-Business (the effective use of new information technologies in business, by developing an alternative sales channel with relatively low costs) is made an individualized approach to client relationships and manages IT using establishing relationships with incomparably greater number of clients, from traditional approaches.

2 Need to develop e-business software

The climate information society, in terms of the technological infrastructure is continually improving, the ability to specialize and to reconfigure the functionality is distributed over a field and heterogeneous computing resources is one of the basic problems is solved with the ability to package the functionality so that it can be used by other applications. Additional business will benefit from consistency rules, the more speed the implementation of a rapid response marketing and change management.

These needs - specialization and reuse of software functionality - have led to new models of software structure. An application is a set of services required to solve an economic problem. Whenever two applications need the same service, they will use a common service, that same piece of software that implements the designated service. Logical analysis and design process involves the discovery service (and use of its character) needed to solve a specific business problems. The physical aspect of the process is to decide which services will be developed based on new components (ideally less) and that service will be handled by existing components in use (ideally more).

Component-based software development refers to techniques and tools that allow software from prefabricated components. These applications are, in a sense, the next generation of client-server applications and the tendency to reduce end-user involvement in the control flow processing, computing power is distributed between client and server.

These types of technologies have a number of features, of which the most important can be summarized as follows:

- are directed toward the goal enable the company to focus on a target, based on their skills - which makes them successful and to be unique;
- provide prompt responses have the ability to provide prompt responses to customers or market opportunities or other elements in the external environment;
- have the advantage of flexibility at the operational and business processes;
- robust offering the ability to respond to any changes at the level of decision making in business and the technological environment.



Fig.1. The key components of e-business application

Companies can achieve and needs using new technologies developed on the experience taken from

existing architectures, and those who develop e-business technologies must be based on good architecture, meaning a good definition of it. Attributes of these technologies that offer flexibility, rapid response and efficiency in demand by organizations that implement them are: integration, virtualization, automation and standardization open (permissive) [1].

3. Presentation application ePortal

3.1. Application objectives

Need for ePortal application in terms of business development is required by the following economic aspects:

- reduce costs for both the entity and its business partners;
- reduce the time of sale and thus improve planning;
- standardization processes and expand enterprisewide scalability;
- to obtain competitive advantage in the market.

In developing e-commerce portal ePortal we considered the following objectives[2] (Figure 2):

- □ structured communication, effective collaboration;
- □ complete and closed circuit of orders and deliveries;
- □ planning and inventory optimization;
- effective management of customers and suppliers.



Fig.2. Business optimization

Business optimization is achieved by:

- real time processing;
- > monitoring the contractual and payment;
- friendly interface, which one requires specialized training;
- automatic identification of exceptions;
- increase data accuracy;
- proactive and effective communication with business partners;
- increase efficiency orders / deliveries;
- lower accounting costs;
- eliminate the storage of products and reduce telephone costs;
- collaborative planning;
- automatically update inventory and automatic volume control needed to buy;

> make electronic payments, electronic bidding prices.

Due to more complicated problems that managers have to solve in the current period, in an environment, internal and external, more complex, it is necessary that management have to be made available in a coordinated system extended to all levels of companies. In order to act effectively in pursuit of the main functions of management planning, organization, coordination and control, current manager needs information, which - to be useful - must be of good quality and available in time (sometimes instantly).

In this context we have considered implementing a solution with open architecture-oriented applications available on-line services (Figure 3) that will extend business process automation to the entire internal value chain of business partners: suppliers, producers and customers.



Fig.3. Implementing a solution that enables online trading

Thus, full integration of processes and value chain visibility will attract the quality of service differentiation (Figure 4).



Fig.4. Extending process automation to the entire value chain

3.2 Optimization of planning by using the portal

Planning, one of the most important benefits of such a solution approach to integrate business processes aims to decrease inventory costs and commitments to the client, referring particularly to: forecast accuracy, inventory levels, planning cycles the accuracy commitments, value chain visibility.

Planning reflects the major issues of planning and budgeting processes, made mainly as training activities of a financial year. The example application created, we considered the following planning cycle: it generally starts with sales planning, based on their defining the production plan, resource plan, and plan and then supply the income and expenses. Plans may be updated during the year to highlight the modifications to the strategy or tactical approach to organization and changes dictated by market developments.

To obtain information flow is necessary to obtain historical information from sales, supply, capacity, production, and external information - kind of market trends, raw material etc. Making plans and budgets is through successive refinements, simulations, scenarios "what happens if" etc. In terms of system for this type of applications suitable type systems Data Warehouse and Decision Support - the multidimensional technology.

Objectives of planning activity within the application ePortal reflected the following benefits:

- lead to a correct forecast demand through collaborative planning;
- reduces inventory by optimizing inventory deliveries increasing accuracy;
- reduce planning cycles with full optimization and planning;
- increased accuracy commitments to customers by including information from production, transport and capabilities of suppliers;
- identify and resolve exceptions occurred at all levels through integrated management company.

Systemic approach to the solution components and advanced planning are identified in Figure 5, may be observed that this solution will:



Fig.5. Advanced planning solution between business partners

- simultaneous material and capacity planning;
- commitment to Partners (available / able);
- advanced simulation capabilities;
- > materials and capacity planning in real time;
- > allocation of purchase orders based on historical data, the delivery plans of providers, provider-

specific command modifiers cycle times and providing specific article.

<u>Planning sales</u> next period, the application ePortal is based on sales in previous years, current contracts, market developments, new products, new markets, new technologies etc. Sales can be broken down by product, geography, distribution channels, sales channels, sales agents, etc. Objective achieved by planning sales, aimed to implement a multidimensional technology planning subsystem, with scenario simulation capability. Such a system can be extended with a foresight (forecasting) and control account (online sales), which in turn can be integrated with the order (order entry).

We considered also the possibility of providing direct information to a marketing system that uses information to optimize sales training sales campaigns.

<u>Materials supply planning</u> and materials allow the following functionality:

- scoreboard scheduler use to determine resources overused and tracking that can create any production interruptions;
- re-planning orders, change quantities or to remove places narrow limits;
- re-planning components and load resources in real time.

The system must receive data from the production system (for other raw materials) of planning production, existing stocks - by split materials in systems maintenance and repair management (for materials, spare parts), the planning of repairs and statistics unplanned repairs and investment system (investment material) of investment planning, investment budgets etc. The system will automatically generate purchase requests (requisitions) when it comes to a planned period (Figure 6).



Fig.6. Automate business processes of the beneficiary

Of particular importance we attached to the supply cycle where contracting, delivery, transportation, receipt, unloading and storage takes place long intervals, there are risks of exceeding the standard limits set.

Strategic products (raw materials) are defined framework agreements between business partners,
including delivery terms, terms, quality, etc., negotiated terms being found in the purchase orders. This process includes the marketing supply (requests for proposals, bids, offers analysis, supplier analysis, preferred suppliers, catalogs, prices, etc.).

Objective achieved for planning to manage the supply and suppliers, was to implement a contract management subsystem (frame, dates, amounts, payment methods, etc.). Given the particular enterprise customer, such a system aims to optimize the short-term planning work related to production planning and supply of sales will ultimately lead to improved regularity of production, minimize inventory of raw materials / finished products and thus optimize financial exposure.

In this context, supply optimization aims to:

- ✤ reduce overall cost of procurement;
- suppliers sorting strategy;
- ✤ full automation;
- ✤ collaboration with suppliers;
- performance indicators (level of expenditure/ supplier).

As can be seen from Figure 7, in terms of customer ePortal application involves:

- □ identifying the best deals through unified Internet search tools, support for complex negotiations;
- □ automation of the entire flow of purchases by processing all the necessary flows payment;
- □ increase efficiency by working with suppliers through the portal access real-time messages in XML format;
- □ identify opportunities to reduce costs through integrated analysis of costs and performance of suppliers.



Fig.7. Simplify logistics by eliminating the storage

Furthermore the correlation of supply processes of dissolution will lead to substantial improvements in performance and profitability of enterprises, as seen in the examples shown below the business value chain is improved due to automated data management.

Optimizing sales and delivery activities aims to improve services, delivery speed, and lower costs through the following facilities:

- demand multi-channel;
- configurable products;
- collaboration;
- dynamic business practices;
- integrated logistics processes.

Management solution chosen sales and deliveries in the application has the following features:

- Supports flow control completely from the collection by multi-channel operation, configuration, configurable price, extended delivery promises;
- Encourage collaboration through automated processing, various delivery options, self-service;
- Use dynamic business processes workflow-based architecture, friendly user interface, control at the command line;
- > Processes have provided integrated logistics warehouse management and transport. Activities in this area relates to sales, marketing and logistics. It is important that in the end all sales channels (sales contract with strategic customers, direct sales based on ad hoc orders, retail sales in your network) and customer to converge into a unified system of reporting. It is possible to track existing and potential customers, to make predictions (forecast-s) to be compared with budget plans and sales channels, agents, etc. products. Marketing plans are organized according to the company's strategy for expanding its product range, market. Are reports of all types - sales, products, geographies, channels, agents, and client - to compare the achieved / planned.

Objective we have considered on sales activity was the implementation of a subsystem customer groups, sales channels, agents, products to short-term sales forecast and sales plan update. Forecast had related to the plan and achievements and thus tracks the effectiveness of sales channels, agents, and update plans for sales, production, supply, resources and budgets. Also, in this situation and recommend implementing a marketing professional, who can do, besides launching a campaign, and campaign results management, resource management, customer campaigns etc. to be addressed. Is excluded and implementation in a given time horizon, a customer relations center on different channels (call / interaction center).

As can be seen from Figure 8 in terms of stock management, application ePortal provides the following features:

- complete integration;
- multiple implementations, warehouses, addresses;
- flexible structure code article;
- groups of products as user-defined criteria;
- follow-up version, batch, serial number;
- cost methods: standard weighted average periods weighted average, FIFO, LIFO;
- regular inventory and annual inventory;
- forecasts;
- planning and flexible supply of stock.



Fig.8. Inventory management solution fully integrated

4 Functional structure of the application ePortal

Electronic Commerce Portal is a collaborative application that allows client companies and partners, to communicate via the Internet. This allows them to have real time information such as orders placed in the system and delivery schedule, to meet the client company providing order confirmations, shipping notifications, requests for amendments or planning details. On the other hand, the application allows the company to make acquisitions, to seek information about orders, shipments, receipts, invoices and other payment information on all suppliers and subunits of its business and to respond to requests for modification issued.

The interface is intuitive portal to facilitate communication with business partners, the principal objective is selling the markets both domestic and foreign markets and create an alternative channel to sell on the Internet, from traditional sales channels. Thus, it contains visual elements in English to facilitate communication with foreign partners, but in addition the application has the option for English, which can be selected to register as external users of the portal.

Thus, business partners can remove geographical barriers between the places of their work by using this flexible, fully integrated, significantly reduced cost compared to traditional operation to improve their own business and services, and default to increase market competitiveness by using the-art work (Figure 9).

EPortal application supports transactional documents and documents open read only (read-only):

- the transaction documents can be sent confirmation or request for change order or you can initiate transactions that advance notification of delivery, delivery notification or pre-payment of bills;
- read the documents you can view information about delivery schedules and forecasts of payment forecasts; Thus, a provider can access information about inventory, invoices and payments; it can send



solution

For this, everyone must register first as a business partner, then as the user application to access the Internet. The procedure for authorizing access site begins with an invitation to visit the site, sent to the e-mail provided by pressing the button "Call" and the type of business relationship with your partner had, access is granted or not specific features of the application. Thus, domestic buyers and the company and the system administrator will have specific access rights to the application menu.

After receiving the address of visiting the site, which is the address link to the next step of the consent procedure, the user will access the recording, which will be entered particulars of the company for which he works. Access to Electronic Commerce Portal will be thus secured the ID and password in the login screen, then through the home page of the application where features are displayed according to the authorizations that the user has the system (eg access to the administration menu is further restricted). Index Browser is actually the main portal and search tool and information retrieval transactions.

At the top of the screen are tabs to access menus as follows: residence, Orders, Shipments, Planning, Finance, Product, Decision Support and Administration. On the right side of the screen is displaying a high level diagram of data flow within the portal, which can be accessed for quick reference to the desired page, link to it containing addresses key application functionality.

In the center of the home page of the application are displayed most recent five notices containing system messages pending to be revised, the latest five orders in the system can view details of which will address about the number of command and supply most recent five full list of which can access all of this page by clicking on the "list".

Some notifications are only for viewing, others require user intervention, but to see the details necessary to access your connection to the subject.

* Orders

Real-time data provided by Electronic Commerce Portal allows exchange of information on the sale of the provider and the client company during flow tracking purchase orders: view sales orders and contractual arrangements; order confirmations; submitter change requests; tracking stocks to third parties; can share supplies in convenient quantities; cancel delivery orders; can view the revision history of the procurement documents and agreements established with the supplier; information and retrieve commands can be followed throughout the procurement process-payment; planning collaborator type.

Deliveries

Electronic Commerce Portal lets you view deliveries on this basis can be created or canceled the notifications in advance of shipping the products. Furthermore you can view programming schedules deliveries can be sent receptions / payment receipts or notices on bills that have not been fully paid. Access Menu Deliveries will be made on the home page of the application by going to the Shipping tab of the main menu, submenu appeared we can: wanted delivery; hierarchical registration supplies (no transport, no. Lot, no. Recording); optimized charging by volume; check returns; view inspections; internal logistics.

If the purchasing company has implemented a quality assurance system and has established a plan for testing products from suppliers, test results will be included in e-Commerce Portal. They will be evaluated before the goods are delivered, so may be canceled deliveries of those products which do not prepare quality standards set.

Planning

EPortal allows forecasting application development, using existing information in the database. This is based on consistent data that are maintained throughout the product development cycle, such as product number, or the registration of orders, etc. modifiers. In terms of inventory, the inventory management environment allows the seller collaborationist and maintaining custody of spare parts distributors or data products for sale: retrieve information on registered products; update command modifiers; updating capacity constraints; item tracking in custody; listing details of items to be contracted; listing of delivery forecasts; the vendor inventory management; reports tracking transactions; graphical statistical evaluation of certain values. Communication through e-Commerce Portal allows not only a modern and competitive way to manage a partnership relationship between the company and its business partners, on the operation (issuing documents, inventory on-line), negotiation and electronic payments, but and better plan activities based on existing resources, which are known as partners.

Billing and payments

Invoices and payment information can be viewed by accessing the Financial tab, in the corresponding menu is available all necessary information about the status of your invoices and view real-time debt or payment (partial or full) done: online billing; management accounts; monitoring payment schedules; reconciliation of accounts.

Management partners

Management profiles partners involved in online transactions enables the key attributes of business relations, enabling them to provide the company with accurate and timely purchasing. These include their address, major contacts, bank details, type of business and the work it carries information about the products they sell and the services they provide. Access the menu administration is through Admin's tab (Administration) only authorized and personalized update information is quickly realizing, be detailed: attract business partners; self check suppliers and users; setting security levels in the system; set of contact information: address, phone number, e-mail; setting category bullets d.p.d.v. business activity; identify services and products provided; definition of bank details.

Tracking performance can: tracking performance indicators: price, quality, delivery services; comparing the results; forecasts.

Links between the main features of the application are illustrated in Figure 10 which shows that e-Commerce Portal nucleus consists of order management; data flow is explained in detail below.



Fig.10. General diagram of the application ePortal

EPortal application implies that the sale / purchase are made from a command that can be launched under a contract system or directly, depending on conventions established a priori between buyer and seller. Order establishes specific conditions of purchase (quantities, prices firm, delivery, payment terms, quality, etc.), and follows a stream of approval.

The next step is acceptance of quantity, quality, and storage.

Starting from the supplier, once received, materials and inventory management enters the recipient.

Adopting this solution replaces a process often conducted manually, which starts at the reception and not the order is dictated by legislative requirements of the accounting.

After the information gathered in studies we observed that generally is originally paper reception (NIR) and then is introduced into the system and supply management process (request for proposals, quotations, etc.) is conducted manually without the support an integrated management at customers and suppliers.

We recommend extending the system by managing the supply requests (requisitions) and provider orders (purchase orders). Orders should reflect exactly the conditions of the contract - quantity, firm prices, delivery dates, terms and payment terms etc. mandatory and must be controlled through an approval cycle. You can also make and supply budget control, in order to control purchases by budgets.

The reception should be made directly to the system - by linking them with approved order - under which purchase was made. Also, for certain products are necessary and a reception quality. Storage needs to be done directly from reception, to avoid operating errors, also the price of storage to be correlated with price control and invoice (if it exists at the time of storage). In addition, the integration of a system of marketing orders (requests for proposals, quotations, etc.).

The invoice must reflect the amounts received and prices and conditions defined in the contract. Therefore, an essential operation is the pairing bills and orders receptions. Further, payment of invoices is treated in the Financial module. The aim was to correlate the financial system (billing and payment) to the trading system (command reception). Each invoice must be correlated with: (amounts, dates, prices, etc.) and reception (quantity, time). Also be addressed possible differences in price between the invoice and order. It is necessary, in addition, implementation of the hierarchy for approval.

5 Benefits of implementing e-Commerce Portal application

Benefits of implementing a software solution as the portal for e-Commerce are given some specific advantages:

- □ exchange of information the massive increase in the amount of information and need to exchange information quickly between different points in geographically distant locations are necessary to connect between autonomous computers;
- □ share resources the cost of increasing the capacity of a distributed system is much smaller than for resources connected to a single-server computer at a time will be exceeded, and in terms of investment, many organizations prefer to buy more computers around reasonable cost and power than to buy one, much stronger, but more expensive;
- □ increased safety in operation if a computer system consists of a single computer malfunction made it impossible to use the whole system, whereas in a distributed system the failure of a node does not disrupt operation of the other, but in most cases they take tasks of the unavailable;
- □ increased safety in operation if a computer system consists of a single computer malfunction made it impossible to use the whole system, whereas in a distributed system the failure of a node does not disrupt operation of the other, but in most cases they take tasks of the unavailable;
- □ performance improvements the presence of multiple processors in a distributed system makes it possible to reduce computing time to achieve a

massive, this is possible by dividing tasks among different processors, subsequent collection of partial results and determine the final outcome, this process is known as name of parallelization of the calculation;

□ specialization nodes - designing an autonomous computer system with more functionality can be very difficult for practical reasons, so this design was simplified by dividing the system into modules, each module implements some functionality and communicates with other modules, revealing two aspects: the first covers the hardware and computing machines that are seen as many autonomous entities and the second part concerns the software and means that users need to collect all programs as a single system.

6 Conclusions

The study which resulted in achieving portal may conclude that the use of leading commercial portals first solve the premises requirements of information society, the changing:

- complex representation of reality (company, customers, products, services, etc..);
- information managed in a system tends to increase in complexity and manipulation must be in a readily perceived by the end user;
- informatics systems must be flexible in relation to changing data structures and must evolve naturally over time, thus following the evolution of the organism it serves;
- ✤ IT systems evolve to broad areas of application approaches to meet the growing needs of users.

Proposed Electronic Commerce Portal allows the introduction of a high level of management and control of online sales, providing cost savings, production planning and most importantly, smooth operation of the circuit information. The computer system is so fully integrated resulting distributed and flexible to support the company's main trading activities.

7 References

[1] Şoavă G., Răduteanu M, Business Intelligent Agents for Enterprise Application, Annals of the University of Craiova - Economic Science, Nr.38, 2010, pag.312-321

[2] Şoavă G., Răduteanu M, Electronic commerce portal, Annals of the University of Craiova - Economic Science, Nr.39, 2011, pag. 232-239

[3] Şoavă G., Răduteanu M, Modeling an electronic portal application, Annals of the University of Craiova - Economic Science, Nr.39, 2011, pag. 240-250

An Implementation of the 16-ary Grid Graphs for the Multiply Layered Rectangular Dissections

Koichi Anada¹, Koushi Anzai², Shinji Koka³, and Takeo Yaku³

¹Waseda Research Institute for Science and Engineering, Waseda University, Tokyo, Japan ²Department Economics, Kanto Gakuen University, Gunma, Japan ³Department Computer Science and System Analysis, Nihon University, Tokyo, Japan

Abstract - Heterogeneous rectangular dissections are frequently used in information processing such as multiple paged books in spread sheet languages and multiple layered image data. In previous studies, a hexadecimal grid graph model was proposed for multiple layered rectangular dissections and certain algorithms were provided.

In this paper, we propose a 32-ary list structure to implement their algorithms for the hexadecimal grid graph model. The list structure has limited number of fields in a record, so the computation time is low.

Furthermore, it has one record for each node in the given hexadecimal grid, and 48 fields for each record. We also show a data format of the whole structure of the list.

Keywords: modeling of spreadsheets, hexadecimal grids, list structures, multiple layered rectangular dissections

1 Introduction

Heterogeneous rectangular dissections are frequently used in visualization of information such as multiple paged books in spread sheet languages and multiple layered image data. In previous studies, a hexadecimal grid graph model was proposed for multiple layered rectangular dissections and certain algorithms were provided [2, 3]. The formalization of business documents as in Figure 1 has become an important subject with the progress of e-commerce and e-government (see, e.g. [1]). In order to formalize financial statements, we have to specify the spatial order of items, and specify calculation methods of categorized items. A context sensitive graph grammar is proposed in [6] that specify financial statements.

In section III, we propose a 32-ary list structure to implement their algorithms for the hexadecimal grid graph model. The list structure has limited number of fields in a record, so the computation time is low. Furthermore, it has one record for each node in the given hexadecimal grid, and 48 fields for each record. We also show a data format of whole structure of the list in section IV.

2 Multilayer rectangular dissections and hexadecimal grids

Figure 1 illustrates a k-layered multilayer rectangular dissection D (left) and its corresponding hexadeci-grid G_D (right).



Figure 1. A *k*-layered multilayer rectangular dissection D (left), its corresponding hexadeci-grid G_D (right).

In a hexadeci-grid, two nodes are horizontally linked if they correspond nearest cells with the ruled line in common, and two nodes are vertically linked if they correspond nearest cells with their corner in common. Figure 2 shows the links around an inner node in a hexadeci-grid, and corresponding cells in a *k*-layered multiply layered rectangular dissections.



Figure 2. Links in a hexadeci-grid.

3 32-ary list structures for hexadecimal grids

In this section, we represent the hexadeci-grids by list structures. We assign one record in the list structure for one node in the hexadeci-grid. The structure of each record is illustrated as Figure 3.



Figure 3. A record in H4CODE list.

Whole structure of the list is illustrated as Figure 4. Each link in the hexadeci-grid is represented by bidirectional links in the H4CODE list.



Figure 4. An H4CODE list corresponding to Figure 1.

We note that the degrees of the record in a H4CODE list are bounded by 32.

4 Data formats

The data format of H4CODE consists of the following header block and repetition of list blocks [5].

1. Header Block

The header block indicates the specification of an H4CODE list.

Table 1. The header block in H4CODE format.

| Field Number | Name | Content | |
|--------------|-------------|--|--|
| 1 | Version | the version of H4CODE | |
| 2 | Row Size | the size of the multilayer rectangular dissections | |
| 3 | Column Size | | |
| 4 | Layer Size | | |
| 5 - 8 | Not Used | | |

2. List Block

Each record in H4CODE list in represented by the following 48 field data. The H4CODE list is represented by a repetition of those 48 field data.

| Table 2. | A data for | r each record | in an | H4CODE | list |
|----------|------------|---------------|-------|--------|------|
|----------|------------|---------------|-------|--------|------|

| Field Number | Name | Content | |
|--------------|-----------|-------------------------|--|
| 01 | node id | id for each node | |
| 02 | cell type | perimeter / inner | |
| 03 | new_right | the link to the nearest | |
| 04 | new_left | common | |
| | | | |
| 42 - 48 | Not Used | | |

5 Concluding remarks

In this paper, we designed a list structure called an H4CODE list for hexadecimal grids, and also designed a data format called H4CODE. An H4CODE list consists of records as illustrated in Figure 3. And, the H4CODE format consists of a header block with 8 fields and a list block with 48 fields. We are implementing several processing systems such as tabular form editing systems, and terrain map systems using H4CODE (*cf.* [3, 4]). We would like to thank Profs. Goro Akagi, Kensei Tsuchida, and Mr. Kenshi Nomaki for their valuable suggestions.

6 References

[1] M. Burnett, A. Sheretov, and G. Rothermel, Scaling Up a "What You See Is What You Test" Methodology to Spredsheet Grids. *IEEE Symp. on Visual Languages* 1999, 30-37.

[2] K. Nomaki, S. Koka, T. Arita, K. Tsuchida, and T. Yaku, A Hexadecimal Grid Graph Model for the Multiply Layered Tabular Forms, *Proc ICCSM* 2010.

[3] S. Koka, K. Anada, K. Nomaki, and T. Yaku, Tabular Form Editing with a Hexadecimal Grid Graph Model, *Proc. IEEE VL/HCC* 2011, 253-254.

[4] S. Koka, K. Anada, K. Nomaki, Y. Shindo, and T. Yaku, Row Manipulation in the Heterogeneous Tabular Forms with a Hexadecimal Grid Graph Model, *Proc. ACM SAC* 2012, 710-711.

[5] K. Anzai, S. Koka, and T. Yaku, H4CODE 1.2 Reference Manual, Working Group of Automata and Its Applications Research Report, 12-002, 4p, May, 2012. URL: http://www.waap.gr.jp/waap-rr/waap-rr-12-002/index.html

[6] T. Yaku, K. Anada, K. Anzai, S. Koka, M. Shimizu, and Y. Shindo, A Graph Grammer Model of Financial Statements with Heterogeneous Perts, *ibid, to appear.*

A study of the magnetic permeability of ferromagnetic thin films for evaluating the GMI effect

Driton Rustemaj and Debashis Mukherjee Engineering Design Department, FESBE, LSBU, London, UK.

Abstract - An investigation of the Giant Magneto Impedance (GMI) effect on planar thin films of soft ferromagnetic material with different compositional ratios has been reported. The report contains the background theory leading to the explanation of the underlying Landau-Lifshitz equation for evaluating the magneto-impedance behavior of the material, in relation to the relevant measurement parameters, i.e. the damping parameter and the resistivity. The findings of this study have been useful for calculating complex magnetic permeability, and the ferromagnetic resonance frequency of soft ferromagnetic materials. The results of the investigation are in good agreement with the previously published data on the magnetic permeability of the soft ferromagnetic materials, demonstrating the validity of the model within the frequency range that has been under investigation.

Keywords: GMI effect, Soft Ferromagnetic materials, Complex permeability constant, Magnetic sensors, Nano-magnetic devices

1 Introduction

The Giant Magneto-Impedance (GMI) effect [1] in soft ferromagnetic materials has attracted significant interest for multimedia applications, such as high sensitivity nano/micro magnetic sensors, filtering devices, and microwave communication circuits [2]. The effect can be attributed to a change in the impedance of the material under the influence of an external magnetic field. It is caused due to the penetration of the field through the skin of the softferromagnetic material. The depth of penetration (δ) depends on the angular frequency (ω =2 π f) of the ac current flowing through the material, being governed by the expression:

$$\delta = \sqrt{\frac{2\rho}{\omega\mu}} \tag{1}$$

where ρ is the resistivity of the material and μ is the complex magnetic permeability of the material which is expressed in the form $\mu=\mu'+j\mu''$.

A variation of the frequency dependent δ in softferromagnetic wires and thin films causes the eddy current losses [3], which in turn decrease the GMI effect in the material. These losses occur due to the existence of magnetic domain structures, which are dependent on the method of growth and related processing techniques [4]. The μ is the most convenient parameter to describe the GMI-effect which in theory, at high frequencies (up to few GHz), can be attributed to the changes mainly caused by the motion of the magnetic moments [5,6].

Theoretically, the GMI effect can be extracted, through obtaining μ by simultaneously solving the already established Maxwell's equations in the form of wave equations, and the linearized form of dynamic equation of motion or Landau-Lifshitz equations [1]:

$$\nabla^2 H - \frac{\mu_0}{\rho} \overset{\bullet}{H} = \frac{\mu_0}{\rho} \overset{\bullet}{M} - graddivM \tag{2}$$

$$\frac{dM}{dt} = -\gamma M \times H_{eff} + \frac{\alpha}{M} M \times \frac{dM}{dt}$$
(3)

where *H* is the magnetic field, μ_0 is permeability in vacuum, γ is the gyromagnetic ratio, α - the Gilbert damping parameter, M - the magnetization vector and H_{eff} is the dc component of the effective field [7, 8].

The analytical solution of these equations is far from straight forward and has only been obtained previously by several authors under simplifying assumptions [7, 8] such as:

i. an in-plane uni-axial anisotropy

- ii. the initial magnetisation is uniform
- iii. the film thickness is very small and
- iv. the only existence of dynamic demagnetization field.

Tanaka et al [9] has, for example, shown the solution to be of the form:

$$\mu_{i}(\omega) = \mu_{i}(\omega)' - j\mu_{i}(\omega)'' = 1 + \frac{\gamma 4\pi M_{s}}{\gamma H_{k} + j\alpha\omega} \times \left[1 + \frac{\omega^{2}}{(\gamma H_{k} + \gamma 4\pi M_{s} + j\alpha\omega)(\gamma H_{k} + j\alpha\omega) - \omega^{2}}\right]^{(4)}$$

where H_k and M_s are the anisotropic field and the saturation magnetization, respectively. They used this equation to find the real and imaginary parts of the magnetic permeability. We have used the same equation in our calculation with an extended frequency range of up to 10 GHz and in combination with studying the effects of the variation of the damping parameter and the resistivity. This enabled us to simulate the waveforms of the magnetic permeability for various sample dimensions and composition, with a view to exploring device structures for a range of applications

The evaluation of the GMI effect is involved with the measurement of the μ value in the element, which is primarily dependent upon its size, geometry and the composition for both homogeneous and inhomogeneous types. The measurement is also influenced by the fabrication method and the associated annealing process, which can affect the internal alignment of the easy axes of the magnetic material. There are other numerous factors depending on the type of materials that have to be considered in evaluating the GMI effect.

When an ac current (I = $I_o e^{-j\omega t}$) flows through the planar thin film, it brings about changes in its magnetic domain structure, causing a variation in the anisotropic field within it. This will directly affect the permeability value. which will also be influenced by the strength of the external magnetic field and the frequency of the current [10].

2 Numerical simulation and discussion

Figure 1 shows the frequency dependent variation of the complex permeability values due to Tanaka et al [9]. Their calculations took into account the thickness

700 20 Complex permeability, $\mu_{e}(\omega)', \mu_{e}(\omega)''$ 50 and 100 nm 200 nm 15 Quality factor, 300 nm J_{010} $\mu(\omega)$ 0 0 2 5 3 4 6 7 1 Frequency (GHz)

Fig. 1 (as Fig. 2 in Ref. [9]) – Complex permeability spectra



Fig.2(a) The complex permeability spectra for the given values of ferromagnetic thin film ($\alpha = 0.015$)

of the sample and a fixed value of 0.015 for the Gilbert damping constant α , as can be seen in the diagram.



Fig. 2(b): The complex permeability spectra for the given values of ferromagnetic thin film ($\alpha = 0.01$)

In obtaining these results, the programming code was written in Matlab, using the appropriate parameter values, taken from Tanaka et al [9] and converted for SI system compatibility. As expected, the figures show good similarity with the plot of figure 1 upto a frequency of about 7 GHz. The test for an alternative damping constant α (0.01), was also carried out in this study in compliance with the published experimental data. The initial value of the real part of the magnetic permeability corresponding to the lower end of the frequency [8, 9], is also strikingly similar to that in figure 1 with a value of the order of 100. As can be seen in the diagrams of figure 2, the damping factor α inversely affects the magnitudes of both the real and imaginary parts of the magnetic permeability. The effect on the bandwidths around the Ferro-magnetic resonance (FMR) frequency is also similar, although appearing to be relatively more pronounced for the imaginary part.

3 Conclusions

Numerical calculations have been done using Matlab, to evaluate the complex permeability of Ferromagnetic thin films using published data and theoretical work. The values of the real and imaginary parts of the complex magnetic permeability have been predicted through assuming damping parameter (α) values of 0.01 and 0.015 with a frequency range of upto 10 GHz. The results show good agreement with previously published results, both in terms of the shape of the plots and the ranges of values obtained for the FMR. The effect of the variation of resistivity was also studied. The results provides a useful basis for predicting the GMI effect of thin films for various applications.

4 References

 L.V. Panina and K. Mohri, "Magneto Impedance Effect in Amorphous Wires", Appl. Phys. Lett., 65, 1189-1191, 1994.

[2] P. Martin and A. Hernando, "Applications of amorphous and nanocrystalline magnetic Mater", J. Magn.Magn. Mater., 215-216, 729-734, 2000.

[3] Y.F. Li et al., J. Magn.Magn. Mater. 268, 57-61, 2004.

[4] S. Atalay, Physica B, 368, 273-278, 2005.

[5] L Kraus, Sensors and Actuators A, 106, 187-194, 2003.

[6] M. Knobel, K.R. Pirota, J. Magn. Magn. Mater. 242-245, 35-40, 2002.

[7] Chengyan Dong et al., J. Magn. Magn. Mater., 250 288-294, 2002.

[8] V. Becker, K. Seeman and H. Leiste, J. Magn. Magn. Mater. 296, 37-45, 2006.

[9] Terumitsu Tanaka et al., IEEE Transactions on Magnetics, Vol. 40, No. 4, pp. 2005-2007, July 2004.

[10] N.A. Usov et al., J. Magn. Magn. Mater. 185 159-173, 1998.

8k-ary Grid Graph Modeling of the Rectangular Dissections

Takeo Yaku¹, Koichi Anada², Koushi Anzai³, Shinji Koka¹, and Kensei Tsuchida⁴
¹Department Computer Science and System Analysis, Nihon University, Tokyo, Japan
²Waseda Research Institute for Science and Engineering, Waseda University, Tokyo, Japan
³Department Economics, Kanto Gakuen University, Gunma, Japan
⁴Department Information Sciences and Arts, Toyo University, Saitama, Japan

Abstract - Rectangular and rectangular piped dissections are commonly used as structures in information processing. Among them are tabular forms in document processing, tables in spread sheet processing, land forms in GIS and raster data in CG. We note that such transformation often include ruled line preserving transformations.

This paper deals with graph modeling of rectangular and rectangular piped dissections with respect to ruled line oriented transformation. We survey octal, hexadecimal and tetraicosa grid graph representation of rectangular and rectangular piped dissections. Moreover, each cell unification of these grid graph representations runs in O(1) time, since the number of links around a node is bounded by 8, 16, 24, respectively.

Keywords: modeling of spreadsheets, rectangular dissections, rectangular piped dissections, ruled line oriented transformations

1 Introduction

This paper deals with graph modeling of rectangular and rectangular piped dissections with respect to ruled line oriented transformation such as the tabular form editing in document processing and LOD problems (see e.g. [6]). While, rapid rendering requires rapid transformation among multi resolution graphics.

Quadtrees and octrees are well known graph models for rectangular and rectangular piped dissections, respectively. Those data structures are originally introduced for information retrieval in CG. So, those require rather large computation time and rather complicated program structures for ruled line oriented transformation. In this paper, we survey octal, hexadecimal and teraicosa grid graph representation of rectangular and rectangular piped dissections.

2 Octgrids for the rectangular dissections

We introduced octal degree heterogeneous grid graphs, called *octgrids* (see e.g. [2, 5, 6]), that represent heterogeneous rectangular dissections, and provide efficient

algorithms for ruled line preserving transformation of CG objects.

The octgrid for a rectangular dissection D is defined informally as follows: Each node in octgrid corresponds to one rectangle (cell) in D. Two nearest nodes are linked if the corresponding two cells in D are nearest and have a ruled line in common as in Figure 1. Figure 2 shows a rectangular dissection and the corresponding octgrid.



Figure 1. Links around a node in an octgrid.



Figure 2. A rectangular dissection (left) and its corresponding octgrid (right).

Accordingly, we obtain a cell unification algorithm that runs in O(1) time. From these properties, we obtain efficient resolution reduction algorithms that provide 3D maps with the appropriate resolution [6].

3 Hexadeci-grid for the multi-layer rectangular dissections

We generalized octgrids in order to represent multi-layer rectangular dissections and introduced hexadecimal grid graphs called *hexadeci-grids* (see, e.g. [4]), that are applied to multi-page books in spread sheet languages (see, e.g. [1]) and stratum maps in GIS for example. Two nodes in a hexadecigrid are linked horizontally as in an octgrid, and linked vertically if corresponding two cells are nearest and have a corner in common. Figure 3 shows links around a node in hexadeci-grid. Figure 4 shows a multi-layer rectangular dissection and its corresponding hexadeci-grid [4]. Figure 5 shows a concept of the multi-layer rectangular dissection and a multi-page book.



Figure 3. Links around a node in a hexadeci-grid.



Figure 4. A multi-layer rectangular dissection (left) and its corresponding hexadeci-grid (right).



Figure 5. A concept of multi-layer rectangular dissection of multi-page book.

4 Tetraicosa-grids for the rectangular piped dissections

Furthermore, we generalized hexadeci-grids and introduced 24-ary grid graphs called *tetraicosa-grids* (see, e.g. [3]) which represent the rectangular piped dissections for voxel representation. Two nodes in a tetraicosa-grid are linked if the corresponding two voxels are nearest and have a beam in common. Figure 6 shows links around a node in a tetraicosa grid. We show a rectangular piped dissection and the corresponding tetraicosa-grid in Figure 7.



Figure 6. Links around a node in a tetraicosa-grid.



Figure 7. A rectangular piped dissection (left) and its corresponding tetraicosa-grid (right).

5 Properties

Each cell unification of octgrids, hexadeci-grids and tetraicosa-grids runs in O(1) time, since the number of links around a node is bounded by 8, 16, 24, respectively. From these properties, we can show that 8k-ary grid graph model provide rapid ruled line preserving algorithms (e.g. [3, 4, 6]).

6 Conclusion

We surveyed 8*k*-ary (k=1, 2, 3) grid graph models for the rectangular and rectangular piped dissections. The authors would like to thank Professors Kimio Sugita of Tokai University, Goro Akagi of Kobe University and Kensei Tsuchida of Toyo University for their valuable suggestions. They also thank Mr. Kenshi Nomaki, Mr. Yuki Shindo, and Ms. Chiaki Arai of Nihon University.

7 References

[1] M. Burnett, A. Sheretov, and G. Rothermel, Scaling Up a "What You See Is What You Test" Methodology to Spredsheet Grids. *IEEE Symp. on Visual Languages*, pp.30-37, 1999.

[2] T. Yaku, Representation of Heterogenenous Tessellation Structures by Graphs. In: *Memoir of WAAP Meetings* 108, 6p, Dec. 2001. http://www.waap.gr.jp/waap-memoir/waap108/waap108_02-yaku/011201waap108table-rep-doc.doc

[3] T. Arita, S. Kishira, T. Motohashi, K. Nomaki, K. Sugita, K. Tsuchida and T. Yaku, Implementation of 24-ary grid representation for rectangular solid dissections, *Proc. 4th GRAPP*, pp. 103-106, 2009.

[4] S. Koka, K. Anada, K. Nomaki, Y. Shindo, T. Yaku, Row Manipulation in the Heterogeneous Tabular Forms with a Hexadecimal Grid Graph Model, *Proc. ACM SAC*, pp.710-711, 2012

[5] T. Yaku, K. Anada, K. Anzai, S. Koka, M. Shimizu, Y. Shindo, A Graph Grammar Model of Financial Statements with Heterogeneous Parts, *ibid, to appear*.

[6] G. Akagi, K. Anada, S. Koka, Y. Nakayama, K. Nomaki and T. Yaku, A Resolution Reduction Method Multi-resolution Terrain Maps, *SIGGRAPH* 2012 *Poster*, 2012, *to appear*.