

Apr 20th, 12:00 AM

Proceedings, MSVSCC 2017

Old Dominion University, Department of Modeling, Simulation & Visualization Engineering

Old Dominion University, Virginia Modeling, Analysis & Simulation Center

Follow this and additional works at: <https://digitalcommons.odu.edu/msvcapstone>



Part of the [Engineering Commons](#)

Recommended Citation

Old Dominion University. Department of Modeling Simulation & Visualization Engineering, & Virginia Modeling Analysis and Simulation Center. (2017, April 20). Proceedings, MSVSCC 2017. 11th Annual Modeling, Simulation & Visualization (MSV) Student Capstone Conference, Virginia Modeling, Analysis & Simulation Center, Suffolk, VA. 211 pp. <https://doi.org/10.25776/zjc7-vp17>

This Other is brought to you for free and open access by the Virginia Modeling, Analysis & Simulation Center at ODU Digital Commons. It has been accepted for inclusion in Modeling, Simulation and Visualization Student Capstone Conference by an authorized administrator of ODU Digital Commons. For more information, please contact digitalcommons@odu.edu.



MODELING. SIMULATION. VISUALIZATION.

THE 11TH ANNUAL
STUDENT
CAPSTONE
CONFERENCE
April 20, 2017
VMASC, SUFFOLK, VA 2017



PROCEEDINGS



VMASC



The Virginia Modeling, Analysis and Simulation Center (VMASC) is a university-wide multidisciplinary research center that emphasizes modeling, simulation, and visualization (MS&V) research, development and education.

VMASC is one of the world's leading research centers for computer modeling, simulation, and visualization. The mission of the Center is to conduct collaborative MS&V research and development, provide expertise to government agencies and industry, and to promote Old Dominion University, Hampton Roads and Virginia as a center of MS&V activities. Working with more than one hundred industry, government, and academic members, VMASC furthers the development and applications of modeling, simulation and visualization as enterprise decision-making tools to promote economic, business, and academic development. Annually, the Center conducts approximately \$10M in funded research.

Old Dominion University is a state-assisted institution and one of only four Virginia schools in the Carnegie Research Universities (high research activity) category. The University offers a complete range of Modeling & Simulation degree options from Bachelor's to Ph.D.

The Hampton Roads region is home to the Joint and Coalition Training (JCW), the US Army's Training and Doctrine Command, the Military Transportation Management Command, NATO Allied Command Transformation, the Armed Forces Staff College, the U.S. Navy's Commander Operational Test and Evaluation Force, the Naval Sea Systems Command, and the Space and Naval Warfare Center. In addition, the Department of Energy's Jefferson Lab, NASA-Langley Research Center and numerous regional industries are important users of MS&V technology. The economic value of MS&V-related business activity in Hampton Roads is estimated to be over \$500M.

VMASC concentrates on eight core modeling and simulation applied research areas:

- Transportation
- Homeland Security and Military Defense
- Virtual Environments
- Social Sciences
- Medicine & Health Care
- Game-based Learning
- M&S Interoperability
- System Sciences

STUDENT CAPSTONE CONFERENCE

MSVE

VMASC relocated to its new facility located at 1030 University Boulevard, Suffolk, VA on September 24, 2007. The Center also maintains an additional location on the main ODU campus in the Engineering and Computational Sciences building.

The Suffolk facility includes a multi-purpose computer laboratory with both Windows-based and Linux-based workstations for computer simulation development, labs dedicated to transportation, computational science, game-based learning, and a virtual simulator lab supporting live, virtual, and constructive simulation integration.

The main VMASC facility primarily supports military, defense and homeland security, social science, medical, computational science, and enterprise engineering research efforts plus administrative and program support. The facility on campus consists of approximately 6,000 square feet of lab space including two general-purpose labs, a visualization lab, a human factors lab, and a 74-seat virtual reality theater supporting both research and teaching requirements. This facility primarily supports the Center's visualization and medical research and development. It also serves as the on-campus presence for VMASC with linkages to the other departments throughout the university.

The MSVE Department

Engineering and Computational Sciences Building on the Old Dominion University Norfolk Campus. In addition to the department and faculty offices, this facility also houses several instructional and research laboratories, a virtual reality theater, and a four-walled C.A.V.E. (Cave Automatic Virtual Environment). A significant resource to the department is the Virginia Modeling, Analysis and Simulation Center located adjacent to the University's Tri-Cities Higher Education Center in Suffolk, Virginia. VMASC occupies a two-story 60,000 square foot building designed to support state-of-the-art research in modeling, simulation and visualization. Some of the center's facilities are used in the department's educational programs; in addition, VMASC researchers teach courses and mentor students in the department's academic programs.





Robert Moorhead

received his PhD in Electrical and Computer Engineering (ECE) from North Carolina State University in 1985. He is a Billie J. Ball Professor of ECE at Mississippi State University (MSU). Dr. Moorhead is the Director of the Northern Gulf Institute (NGI), a NOAA Cooperative Institute, and the Director of the Geosystems Research Institute at MSU, which focuses on understanding Earth's natural and managed

systems. GRI focuses on spatial technology, visualization of complex datasets, and computational modeling in agriculture, forestry, water resources, climate, weather, and oceanography. Dr. Moorhead has developed a nationally-recognized research and outreach program focused on the use of small UAS to advance our understanding and management of agriculture and natural resources.

- Dr. Moorhead has published over 200 articles.
- He has been an Associate Editor of IEEE Trans. on Visualization and Computer Graphics.
- Served as chair of the IEEE Computer Society's Technical Committee on Visualization and Graphics.
- He was summer faculty at the Navy for 8 summers in the 1990s.
- Research Staff Member at IBM's T.J. Watson Research Center from 1985-1988.
- He has received competitive research funding from a multitude of federal agencies and industrial concerns.
- He has a patent related to an algorithm used in the MPEG standard.
- Served on the international committee that developed the JPEG image coding standard.



The tracks that the papers were divided up into included the following:

- Gaming & Virtual Reality
- Infrastructure Security, Military Application & Transportation
- Education & Training
- Business & Industry
- Agent Based Modeling
- Medical Simulation
- General Sciences & Engineering

For each track there were awards given out: The Best Paper award, and the Best Presentation Award. Those recipients for best paper are listed below by track.

Agent-based modeling:

Wessam Elhefnawy, "Make my neighborhood safe again: an agent based model for burglary crime prediction and capture its patterns"

General Sciences & Engineering

Nathan Li, Debrup Banerjee, and Jiang Li, "A comparative study of classification schemes in transfer learning for PTSD diagnosis"

Michelle Pizzo, and Fang Hu "Simulation of sound absorption by scattering bodies treated with acoustic liners and the assessment of its high-performance parallel computing capabilities"

Medical Simulation:

Jing Xu and Andrey Chernikov, "Homeomorphic Tetrahedral Tessellation for Biomedical Images"

Business & Industry:

Paul Delimarschi, Jonathan Griffith, Mitchel Howard, Sean McBryde, Gene Lesinki, "Simulation and analysis of the aircraft corrosion control facility at the Corpus Christi Army Depot"

Overall Best Paper- The Gene Newman Award

The overall best paper is awarded the Gene Newman award. This award was established by Mike McGinnis in 2007; the award is given for overall best presentation, best paper, and research contribution. The Gene Newman Award for Excellence in M&S Research is an award that honors Mr. Eugene Newman for his pioneering effort in supporting and advancing modeling and simulation. Mr. Newman played a significant role in the creation of VMASC by realizing the need for credentialed experts in the M&S workforce, both military and industry. His foresight has affected both the economic development and the high level of expertise in the M&S community of Hampton Roads. The Students receiving this award will have proven themselves to be outstanding researchers and practitioners of modeling and simulation.

For the 2017 Student Capstone Conference, The Gene Newman Award went to: ***Paul Delimarschi, Jonathan Griffith, Mitchel Howard, Sean McBryde and Gene Lesinki***, from the U.S. Military Academy Westpoint for their paper entitled '*Simulation and Analysis of the Aircraft Corrosion Control Facility at the Corpus Christi Army Depot.*'





STUDENT CAPSTONE CONFERENCE

2017

PAGE 1

GAMING & VIRTUAL REALITY

PAGE 6

INFRASTRUCTURE SECURITY, MILITARY & TRANSPORTATION

PAGE 13

EDUCATION & TRAINING

PAGE 30

BUSINESS & INDUSTRY

PAGE 68

AGENT BASED MODELING

PAGE 94

MEDICAL SIMULATION

PAGE 128

GENERAL SCIENCES & ENGINEERING



STUDENT CAPSTONE CONFERENCE

2017

GAMING & VIRTUAL REALITY

- Page 2 Katherine Smith, John Shull, Yuzhong Shen, Tony Dean, and Jennifer Michaeli
Old Dominion University
Captivate: Employing Classic Game Mechanics in a Serious Stem Game
- Page 4 Zinat Afrose and Yuzhong Shen
Old Dominion University
Point Cloud Denoising Using Adaptive and Order Statistic Filters

CAPTIVATE: EMPLOYING CLASSIC GAME MECHANICS IN A SERIOUS STEM GAME

Katherine Smith, John Shull, Yuzhong Shen, Tony Dean, and Jennifer Michaeli
Department of Modeling, Simulation and Visualization Engineering
Old Dominion University
5115 Hampton Boulevard, Norfolk, VA, USA
k3smith@odu.edu

ABSTRACT

This paper discusses development of CAPTIVATE, a serious game for calculus and physics veterans education as part of the Stern2STEM project at Old Dominion University. Building on lessons learned from previous development efforts, CAPTIVATE reworks classic game mechanics from popular games and combines them with STEM content in order to engage the student in a familiar narrative while keeping the game content focused. In this work, two of the games developed so far are discussed.

Keywords: STEM education, game development, educational games, symbolic mathematics.

1 INTRODUCTION

Building upon previous work that resulted in the development of MAVEN, a serious game to assist veterans in learning precalculus (Smith et al., 2016), a new game to help veterans learn calculus and physics is underway. So far, the development efforts have been focused on calculus and the topics highlighted in this abstract include function behavior classification and derivatives.

The remainder of the abstract will be organized as follows. First, the Methods section will provide an overview of how classic gameplay can be reemployed to create engaging educational games. Next, the Results section will provide an overview of two games that have been developed so far. Finally, the Conclusions and Discussion section will provide conclusions and summarize future work.

2 METHODS

Good game design focuses not only on game mechanics, but also engages players through a compelling story (Fullerton, 2014). One of the challenges when designing serious games is keeping players engaged while conveying content. When showcasing MAVEN at a variety of events, many individuals responded positively to games that reemployed classic game mechanics from popular games. This is likely due to the fact that the players have a familiarity and emotional connection with the story lines from these games which consequently evoke a positive response. The goal for CAPTIVATE was to develop a series of games for mobile platforms that players could play at any time and in any order to strengthen their calculus and physics skills. Since players can play in any order, having a cohesive narrative that runs through all the games is not feasible. Instead, popular classic games were reinvented by incorporating calculus and physics content.

3 RESULTS

The two games that will be introduced in this work are MineSweeper and FunctionHero (Fig. 1). Game play in MineSweeper begins with the familiar MineSweeper layout updated so that the player is searching through the desert for improvised explosive devices (IEDs). The player explores by tapping a square. In

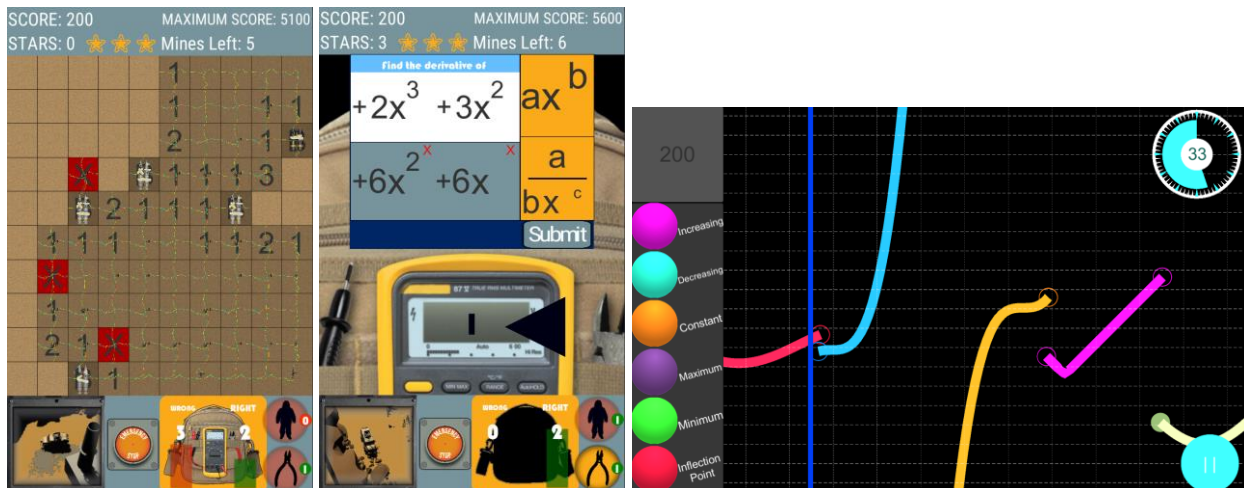


Figure 1: Screen captures from the MineSweeper (left) and FunctionHero (right) games.

traditional MineSweeper, the player would automatically receive information about the number of mines in the surrounding squares. In the CAPTIVATE version, the player is first presented with an expression that they must differentiate. Some questions require application of derivative rules in a single step while others require multiple steps, first applying derivative rules and then simplifying the resulting expression. If the player answers correctly, they are rewarded with the information and the chance to receive a life-saving bomb suit or wire cutters that can be used to skip a question. If they answer incorrectly, a red “X” is placed in the square which penalizes the player by withholding necessary information. Game play ends when the player accidentally triggers an IED or uncovers all the squares. Similar to the classic game, players are rewarded for using given information to mark the squares they determine to contain IEDs.

The game mechanics in FunctionHero are borrowed from the popular GuitarHero games. The player is presented with a piecewise function that scrolls from right to left. They are tasked with correctly classifying the behavior of the function at the point where it intersects the blue target line. For behaviors that occur over a range such as increasing, decreasing, and constant, the player holds down the corresponding button over the entire interval where this behavior is exhibited. Critical points are classified by tapping the minimum, maximum or inflection point buttons.

4 CONCLUSIONS AND DISCUSSION

So far, the gameplay mechanics from a variety of classic games have been adapted with the addition of STEM content to produce games that are engaging and instructive. Future work includes the development of additional games as well as an efficacy study to ascertain the effect of the games as a learning resource.

ACKNOWLEDGEMENT

This paper, and its associated research, was made possible through the Office of Naval Research STEM under ONR GRANT11899718.

REFERENCES

- Fullerton, T. “Game design workshop: a playcentric approach to creating innovative games.” *CRC Press*. 2014.
- Smith, K., Shull, J., Dean, A., Shen, Y., and Michaeli, J. “MAVEN: A Serious Mathematics Game for Veteran Education.” *Proceedings of the 2016 MSVE Student Capstone Conference*: 30-31. Suffolk, VA 2016.

POINT CLOUD DENOISING USING ADAPTIVE AND ORDER STATISTIC FILTERS

Zinat Afrose

Department of Modeling, Simulation and
Visualization Engineering
Old Dominion University
zafro001@odu.edu

Yuzhong Shen

Department of Modeling, Simulation and
Visualization Engineering
Old Dominion University
yshen@odu.edu

ABSTRACT

Complete and clean data are mostly required for various applications that involve point cloud data. The applications include reverse engineering, CAD model generation, modeling and rendering models of 3D data. But in practice, we have to deal with incomplete, unclean data contaminated by noise or missing important features. This paper proposes an approach for noise removal and calibration of noisy point cloud data based on adaptive and order statistic filters. This paper demonstrates the performance and noise-robustness of the proposed methods (vector median filter, fuzzy vector median, adaptive mean, adaptive median, adaptive vector median) applied to standard models as well as synthetic models, and real scenes.

Keywords: Point Cloud Data, Noise reduction, Order-statistic filters.

1 INTRODUCTION

A point cloud is a set of data points in 3D space which is defined by X, Y and Z coordinates. It can be generated by any type of 3D scanners or laser scanners. Since the scanning devices introduce noise inevitably, the point clouds of a 3D object cannot accurately represent object's original shape. Recent studies have been focused on the necessity of robust methods for denoising point cloud while preserving important features. In this paper, image processing filtering methods are extended to remove noise from point cloud data. In addition, we develop a new type of filter, namely, adaptive vector median, for point cloud denoising, which can be utilized for other types of applications as well.

2 METHOD

Six filtering techniques are utilized in this paper: median, vector median, fuzzy vector median, adaptive mean, adaptive median and adaptive vector median filter. The median filter considers each point and its neighborhood and replaces the center value with the median of its neighborhood. The vector median filter is the extension of median filter that considers each point as an inseparable 3D vector. The FVM filter is done in the following way: for a given point cloud dataset, it first computes the normal at each point, then the final outcome is the weighted sum of input point sets, where the weights are determined by the fuzzy relation between each point and the vector median. As an adaptive filter, adaptive mean filter changes its behavior based on the statistical characteristics of the point cloud inside the window with a specified radius. The values of the noisy points, variance of the noise, local mean of the points and local variance of the points are taken into account. Adaptive median filter works in two stages. In the first stage, if the median of the intensity value is between the minimum and maximum intensity values in the window, then it moves forward to second stage else increases the window size. If the window size is less than or equal to the maximum allowed size, then it repeats the first stage, otherwise the output will be the median of intensity value. In the second stage, if the intensity value of the center point of the window is between the minimum and maximum intensity values in the region then it outputs the center value, otherwise it outputs the median value. It changes the window size depending on the max and min intensity values. The adaptive vector median is an extension of the adaptive median filter. The size of the window surrounding each point is variable. This variation depends on the vector median of the points in the present window. If the vector median value is between the max and min intensity value, then the size of the window is expanded. Otherwise, further processing is done on the part of the data within the current window specifications.

3 RESULTS

Fig. 1 shows the visual representations of the results of the filters applied to a gear model. Additive Gaussian noise ($\sigma=0.01$) and impulse noise was added to the original gear model. As can be seen, the edges of the gear (zoomed in) are well preserved. All filters removed noise to some extent. The adaptive vector median filter performs best in handling noise and preserving edges. Fig. 2 demonstrates the performance of the filters based on RMSE and Hausdorff distance. In terms of complexity, median, vector median, FVM, adaptive mean, adaptive median, adaptive vector median have time complexity of $O(N)$, $O(N^2)$, $O(MN^2)$, $O(MN)$, $O(MN)$, $O(MN)$ respectively, where M is number of points and N is the window size. The computational time (in sec) for Gear model is 0.425, 0.512, 0.647, 0.568, 0.599, 0.610 respectively.

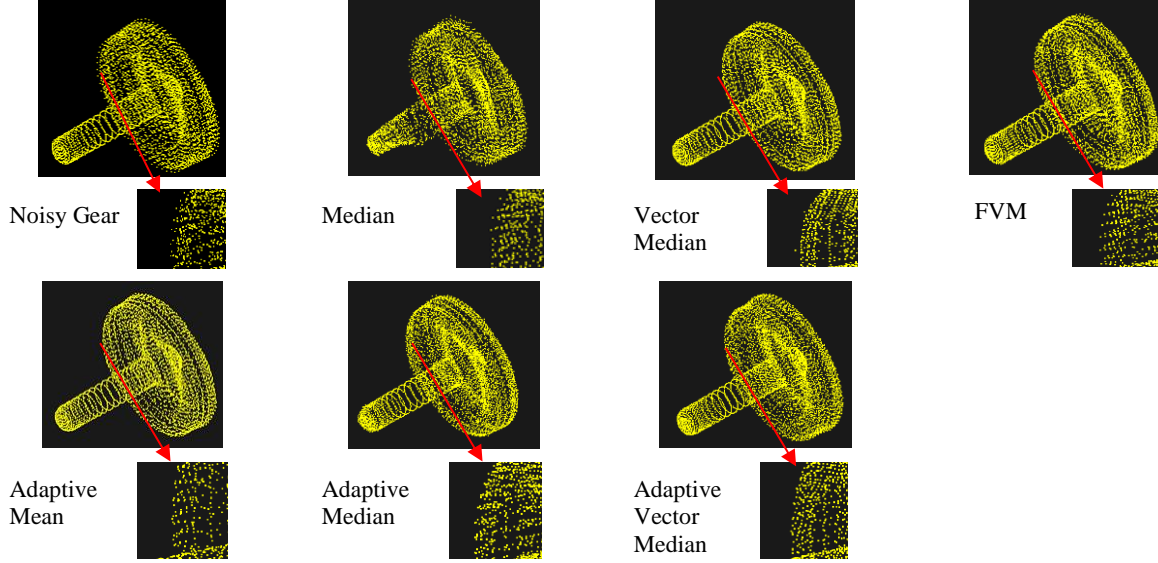


Figure 1: Noise Removal of a Noisy Gear.

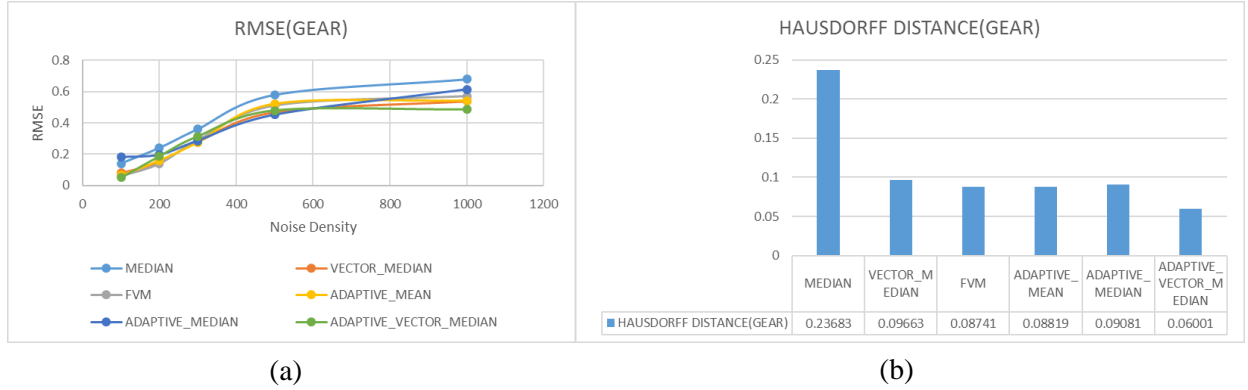


Figure 2: (a) Performance comparison, (b) Hausdorff distance for the model Gear.

4 CONCLUSION

This paper presents six methods for point cloud denoising based on adaptive and order statistic filters adapted from image processing technique. The adaptive vector median, a brand new filter is developed in this paper. All the methods presented in this paper were able to remove noise from point cloud data to some extent. Adaptive median and adaptive vector median filter achieved best performance. A possible extension for this work would be to use high performance computing to speed up the computation time for larger datasets.



STUDENT CAPSTONE CONFERENCE

2017

INFRASTRUCTURE SECURITY, MILITARY & TRANSPORTATION

- Page 7 Mahmud Hasan, Daniel Perez, Yuzhong Shyen and Hong Yang
Old Dominion University
Enhancing Microscopic Traffic Simulation with Virtual Reality
- Page 9 Zhenyu Wang and Hong Yang
Old Dominion University
Integrating Sensor Data for Identifying Secondary Crashes

ENHANCING MICROSCOPIC TRAFFIC SIMULATION WITH VIRTUAL REALITY

Mahmud Hasan, Daniel Perez Ibanez, Yuzhong Shen, Hong Yang
Department of Modeling, Simulation & Visualization Engineering
Old Dominion University
Norfolk, VA

ABSTRACT

Microscopic traffic simulation is an important tool for planning, designing and analyzing transportation systems. Although it has been widely investigated, there is a factor that has not been incorporated yet: human interaction. Human behaviors are overly simplified and all vehicles are controlled by computer algorithms in microscopic traffic simulations. In this paper, we present a framework that introduces human interaction into microscopic traffic simulation using gaming and virtual reality technologies. In particular, some vehicles in microscopic traffic simulations are controlled by humans in a realistic 3D environment that represents various traffic simulation scenarios. To enhance user experience and improve simulation fidelity and accuracy, virtual reality devices are utilized as both outputs and inputs. A prototype framework has been developed to integrate Paramics, a microscopic traffic simulation tool, and Unity, a leading game engine. Various virtual reality devices such as Oculus Rift and dome, have been utilized.

Keywords: Microscopic Traffic Simulation, Virtual Reality, Paramics, Unity3D, Dome

1 INTRODUCTION

Traffic simulation models various transportation systems through the application of computer software to better help plan, design, analyze and operate transportation systems. When the detailed motion of each individual vehicle is taken into account, it is called microscopic traffic simulation. From research to industry, traffic simulation software is widely used. However, the driver behavior models in microscopic simulations are overly simplified, and so are movements of individual vehicles (e.g., lane changing does not take time). Meanwhile, video games, such as car racing and flight simulation games, are extremely popular as they provide highly realistic 3D environments that simulate the real world and engaging user interactions. Gaming technology can be used to complement and improve microscopic traffic simulations by providing users enhanced 3D environments that represent various traffic scenarios and enable human in the loop (HITL).

This paper presents a system architecture that incorporates multiple user controlled vehicles into microscopic traffic simulation and enhances the user experiences by using VR technologies. In the following sections, we describe the system architecture, discuss its implementation and finally draw the conclusions.

2 METHODOLOGY

The proposed system architecture is shown in Figure 1. The system consists of two types of entities: server and clients. The server hosts the microscopic traffic simulation software and resides in the cloud (Internet). A client allows a user or player to view the microscopic traffic simulation 3D environment and

control a selected vehicle in the traffic simulation. The system allows multiple clients so that several or many vehicles can be controlled by real users. Each client receives user inputs such as acceleration or braking and sends them to the server through network communications. The server then runs the traffic simulation, updates vehicle positions at each simulation step, and sends the information back to the clients, which then update vehicle positions in their 3D environments accordingly. Gaming technologies can be utilized by the clients to provide highly realistic 3D environments and some physics based responses well. In addition, direct client-to-client communications are also possible to model and simulate various applications and scenarios, such smart road network and driver distractions (e.g., texting during driving). Microscopic simulation handles the major analysis tasks. In addition, the built-in simulation models manage all non-human-controlled vehicles as well as their interactions with human-controlled vehicles, which were often simplified in conventional driving simulator.

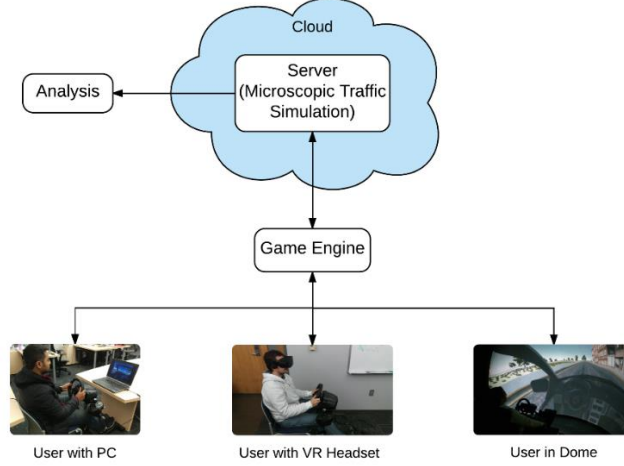


Figure 1: Architecture of the System

3 RESULTS

Quadstone Paramics is a microscopic traffic simulation software package and it is being used to build a prototype of the system discussed in Section 2. Unity is a leading game engine and has been used to develop the client software. A simple highway road with 20 vehicles was created as a testbed. The positional and movement data of those vehicles were collected through Paramics API and fed to Unity3D in real time.

Unity3D has the same scene that runs on Paramics. The positional and movement data received from Paramics are used to create vehicle movements in Unity. Multiplayer networking in Unity was performed by using Photon Unity Network (PUN). The code in Paramics is written in C, and the Unity program is in C#. The communication between these two main programs is done through network sockets. Various virtual reality devices have been utilized, including Oculus Rift, Logitech G920 gaming steering wheel, a customized seat, and a visualization dome.

4 CONCLUSION

This extended abstract presented a framework for enhancing microscopic traffic simulation using gaming and virtual reality technologies. A prototype is being developed using Paramics and Unity. Initial results demonstrated the feasibility of the proposed approach. Several issues need to be addressed, such as time synchronization between the server and client, and final vehicle position computation based on microscopic traffic simulation and physics simulation in the client.

INTEGRATING SENSOR DATA FOR IDENTIFYING SECONDARY CRASHES

Zhenyu Wang, Hong Yang
Dept. of Modeling, Simulation & Visualization Engineering
Old Dominion University, Norfolk, VA 23529, USA
Zwang002@odu.edu

ABSTRACT

The presence of secondary crashes on highways induces not only traffic delays but also safety issues. The prevention of secondary crashes is of great importance to transportation agencies. Prior to the deployment of any effective countermeasure, the first task needed is to find out the secondary crashes. Thus, this paper aims to develop an automated approach to determine secondary crashes. We proposed the piecewise shockwave-based procedure to estimate the impact area (IA) of primary incidents. The proposed procedure fused probe vehicle data with traditional loop detector data for enhancing speed estimation, which in turn led to improved IA estimation. A refined angle summation algorithm was then applied to automatically identify secondary crashes within the IAs. A case study was performed to test the performance of the proposed approach. The results show that the proposed approach can efficiently determine potential secondary crashes with relatively high accuracy rates.

Keywords: secondary crashes, probe vehicle, sensor data, loop detector, data fusion.

1 INTRODUCTION

Secondary crashes often cause additional traffic congestion and safety problems. In order to mitigate the impact of secondary crashes, it is necessary to explore the underlying mechanism of secondary crash (SC) occurrence. However, one cannot easily explore these crashes as they were not regularly documented. To assist the analysis of SC occurrence, the main objective of this study is to develop an approach that helps identify potential SCs based on fused sensor data.

2 PROPOSED METHODOLOGY

This paper proposes the shockwave-based methodology to identify secondary crashes by leveraging the use of probe vehicle data. The proposed framework consists of three main components as shown in Figure 1. First, the loop detector data and the probe vehicle data are fused to improve speed estimation (component in blue color). Second, the piecewise shockwave-based approach is developed to estimate the boundary of the impact area (component in green color). Finally, the refined angle summation algorithm is introduced to automatically classify potential secondary crashes (component in grey color).

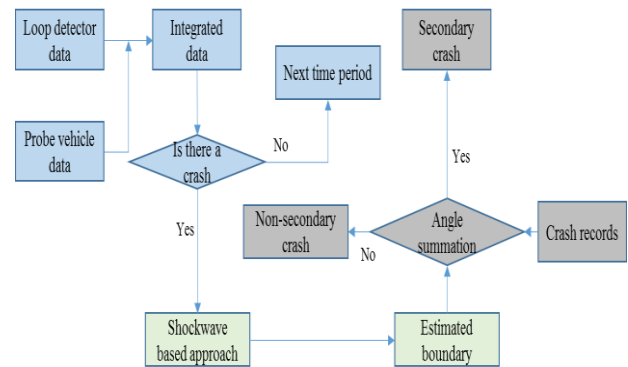


Figure 1 Framework of the proposed methodology

We have developed a data integration procedure to fuse data collected from spaced loop detectors and sparsely sampled probe data. As shown in Figure 2, blue lines represent the probe vehicle trajectories and green dots denote probe data points. The speed contour map was developed based on speed measurements from loop detectors A, B, and C (red zones denote congested area). The original speed contour map was reconstructed with higher spatiotemporal resolutions (i.e., smaller cell size) based on the integrated data. Two scenarios were specifically considered: (a) cells without probe points; and (b) cells with probe points. Depending on the spatial relationship between each cell and the original detector, the speed estimation for the refined cells are obtained. For example, each cell in Figure 2 (b) to (c) represents the coverage of a virtual loop detector. With the estimated speed based on these virtual loop detectors, the improved speed contour map can be developed.

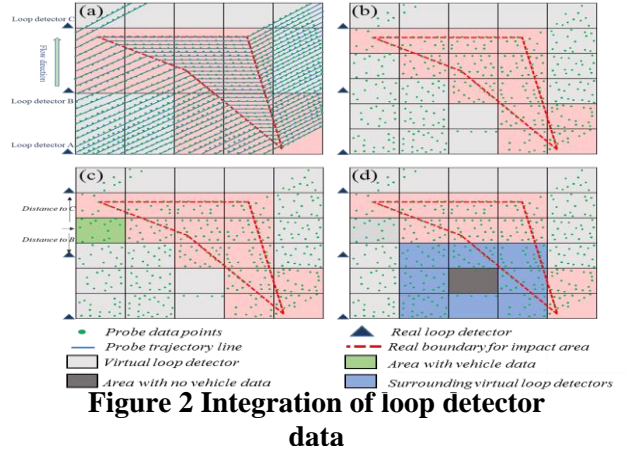


Figure 2 Integration of loop detector data

With the improved speed contour map, we proposed the piecewise shockwave-based approach to estimate the impact area as illustrated in Figure 3. The first piece of shockwave starts upon the occurrence of the primary crash, and its propagation slope was assumed to be linear. Yellow star icons represent the estimated state changes of virtual loop detectors, and yellow circles represent the actual stage changes for virtual loop detectors. By connecting all yellow circles and stars, the purple line segments between virtual loop detectors were constructed. Then, linear regression models were used to smooth the estimation of the boundary (i.e., the dash green line in Figure 3). Once the boundary was estimated, the angle summation algorithm was introduced to identify any potential secondary crashes with the boundary.



Figure 3 Piecewise shockwave approach

3 CASE STUDY AND RESULTS

Actual probe vehicle data from the Mobile century project (<http://traffic.berkeley.edu/>) is used as a case study. Figure 4 shows the speed contour map based on the original loop detectors and on our proposed approach that fused the loop detector data and probe data. It can be seen that Figure 4(b) offers a higher resolution with more smoothed transition between different traffic states.

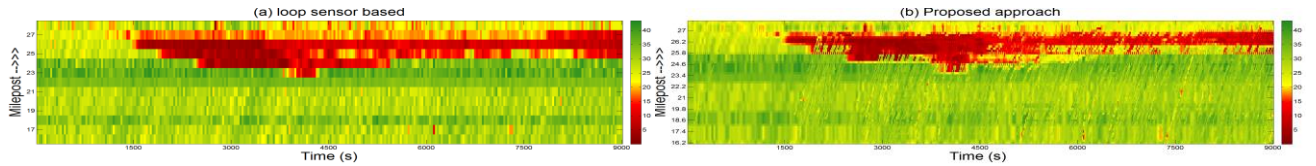


Figure 4 Speed contour map (a) loop detector based (b) proposed approach

4 CONCLUSION

Our proposed approach can improve the construction of speed contour map based on the integration of loop detector data and probe vehicle data. The piecewise shockwave-based approach offers better estimation of the impact area. The refined angle summation algorithm was able to automatically determine the potential secondary crashes. The case study confirmed the efficiency and the performance of our proposed approach.

MODELING FAILURES AND REPAIRS OF A SMART REDUNDANT SYSTEM USING CONTINUOUS TIME MARKOV CHAIN

Manoj Banik

Dept. of Modeling Simulation &
Visualization Engineering (MSVE),
Old Dominion University (ODU),
Norfolk, VA. USA
E-mail: mbani003@odu.edu

Bharat B. Madan

Dept. of Modeling Simulation &
Visualization Engineering (MSVE),
Old Dominion University (ODU),
Norfolk, VA. USA
E-mail: bmadan@odu.edu

ABSTRACT

Data redundancy is very important for Intrusion Tolerant Systems (ITSs) and replication based redundancy has been used to implement such systems. However, plain replication, though effective against intrusions designed to compromise availability, can in fact be detrimental against confidentiality and integrity intrusions. Smart redundancy is based on fragmentation, coding, dispersion and reassembly (FCDR) to eliminate single point of security failure. A data block or packet is fragmented into n fragments, which are augmented by k additional fragments using *Erasur Codes* (e.g., Reed-Solomon codes). Erasure coding ensures that any n (out of $n+k$) fragments suffice to reproduce the original data block. Consequently, the system can survive k availability and integrity attacks, n confidentiality attacks. The paper shows that these models take the form of Continuous Time Markov Chains (CTMSs), which are analyzed to quantify the effectiveness of smart redundancy in tolerating confidentiality, integrity and availability intrusions.

Keywords: intrusion tolerant systems, security, fragmentation, Markov chains, smart redundancy.

1 INTRODUCTION

Cyber systems like other real systems experience failures, which generally occur randomly. Therefore, modeling of failures and subsequent repairs is an important part of analyzing and quantifying a system's performance. This paper's focus is on modeling and analysis of security failures resulting from cyber-attacks and subsequent repairs. Cyber-attacks seek to compromise a cyber-system's confidentiality (C), integrity (I) and availability (A) security attributes (Madan and Banik 2014). Since recovering from an attack can take unpredictably long time, cyber systems used in safety critical applications, need to be inherently capable of tolerating failures resulting from security intrusions, known as ITSs which depend on incorporating redundancy in their design (Reed and Solomon 1960). The attacks and repairs of the data fragments formed a continuous time Markov chain (Trivedi 2001) which is shown in Figure 1. State s_i represents a state with i compromised fragments which means the system has been attacked i times successfully. So, state $(k + 1)$ is availability/integrity compromised state and state n is confidentiality compromised state. We can assume that the system takes $t_{i,i+1}$ random time to make a transition from a state s_i to another state s_{i+1} and $t_{i,i-1}$ random time to transit from s_i to s_{i-1} . Furthermore, all such transition times are assumed to be exponentially distributed.

From quantitative results, it was found that a smart redundant system stays in uncompromised state with higher probability as long as the repair rate is higher than the attack rate. The rest of the paper is discussed in two following sections.

2 PROBLEM FORMULATION

In an ITS, typically a failure resulting from an intrusion has to be dealt with first detecting the underlying intrusion using intrusion detection system (IDS) and after successful detection, a repair process initiates to bring the system back to its pristine state. We model an intrusion as a random event such that k^{th} type of intrusion is assumed to take random interval of time with $EXP(\lambda_k)$ to achieve its objective of causing some kind of damage and IDS is modeled as a Bernoulli process with P_{d_k} as the probability of successful detection of a k -type of intrusion. Conversely, $\widetilde{P}_{d_k} = 1 - P_{d_k}$ denotes the probability of an intrusion remaining undetected. On detecting a failure, the repair process starts instantaneously and the time taken to complete the repairs is modeled as continuous time random process. So, we can model this failure, detection and repair processes using CTMC (Figure 1).

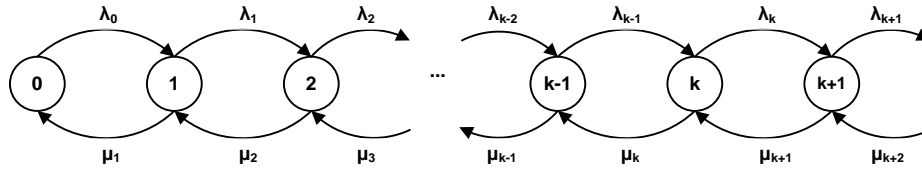


Figure 1: Continuous Time Markov Chain model of attacks and repairs processes in smart redundancy.

3 RESULTS AND CONCLUSION

The marginal probability for the states of a steady state smart redundant system has been found for n, k, λ & μ and is shown in the Figure 2. From the graph it is clear that when the repair rate is higher than the failure rate ($\lambda/\mu < 1$) then the state with highest marginal probability is a good working state which ensures service availability.

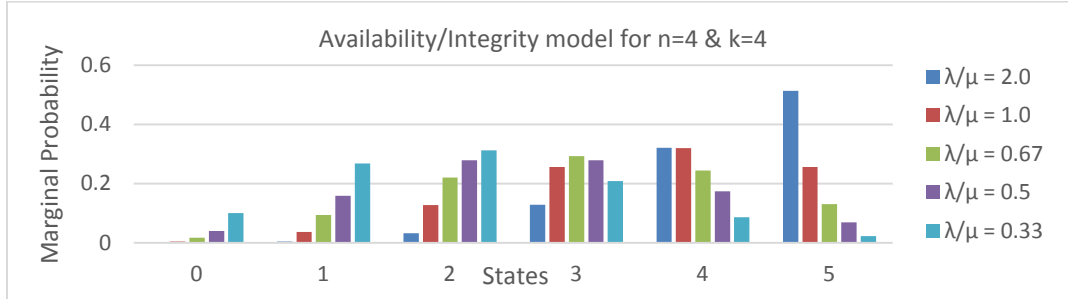


Figure 2: Marginal probability chart for Availability/Integrity model ($n = 4$ & $k = 4$).

REFERENCES

- Trivedi, K.S. "Probability and Statistics with Reliability, Queuing, and Computer Science Applications", John Wiley and Sons, New York, 2001.
- B. B. Madan, M. Banik, 2014. "Attack tolerant architecture for big data file systems". ACM SIGMETRICS Perform. Eval. Rev., 41 (4) , pp. 65–69.
- I. S. Reed and G. Solomon. Polynomial codes over certain finite fields. Journal of the Society for Industrial and Applied Mathematics, 8:300– 304, 1960.
- M. Banik, B. B. Madan, "Modeling and analysis of multistage failures of a system ", Spring Sim 2017, April 23-26, Virginia Beach, VA, USA; Accepted for publication.

STUDENT CAPSTONE CONFERENCE

2017

EDUCATION & TRAINING

- Page 14 Laura Welsch, Johanna Hoch, Andrea Parodi, Rebecca Poston and Muge Akpinar-Elci
Old Dominion University
A Web-based Interprofessional Education Program For School Nurses And Athletic Trainers: A Pilot Study
- Page 17 Levi Warvel, Mark Scerbo, Samantha Zybak, Rebecca Kennedy, Yannis Papelis, Menion Croll, and Hector Garcia
Old Dominion University
Enhancing Speech Recognition in a Virtual Operating Room
- Page 19 Joshua G. Stubbs and Ginger Watson
Old Dominion University
The Teachlive Simulator as a Tool for Learning Versus Tool for Assessment of Teaching and Classroom Management Skills
- Page 23 Khadijeh Salimi
Old Dominion University
A System Dynamic Modeling of Iranian women Movement

A WEB-BASED INTERPROFESSIONAL EDUCATION PROGRAM FOR SCHOOL NURSES AND ATHLETIC TRAINERS: A PILOT STUDY

Lauren Welsch
College of Health Sciences
Old Dominion University
3008 Health Sciences
Building
Norfolk, VA 23529
lwels001@odu.edu

Johanna Hoch
Physical Therapy & Athletic Training
Old Dominion University
3121 Health Sciences
Building
Norfolk, VA 23529
jhoch@odu.edu

Andrea Parodi
Virginia Modeling, Analysis & Simulation Center
Old Dominion University
VMASC
Suffolk, VA 23435
aparodi@odu.edu

Rebecca Poston
Nursing
Old Dominion University
3126 Health Sciences
Building
Norfolk, VA 23529
rdeal@odu.edu

Muge Akpınar-Elci
Center for Global Health
Old Dominion University
3134 Health Sciences
Building
Norfolk, VA 23529
makpinar@odu.edu

ABSTRACT

Improvements in teamwork and communication in the healthcare system have been proposed as a means to improve patient safety and outcomes. Education programs designed to teach these skills are notably absent in less traditional healthcare settings such as school health. This project describes one educational program designed to teach teamwork and communication skills to athletic trainers and school nurses.

KEYWORDS: TeamSTEPPS®, communication, interprofessional education

1 INTRODUCTION

As the healthcare system becomes increasingly complex, healthcare professionals from multiple disciplines are often required to work together in interprofessional collaborative practice (ICP) (Kohn, 1999). Effective teamwork and collaboration among healthcare teams can improve the delivery of care and positively impact patient outcomes, while the absence of communication has been linked to increased medical errors and adverse health outcomes (Kohn, 1999). Interprofessional education (IPE), defined as two or more healthcare groups learning with, from and about each other (Horsburg, 2001), is often necessary to teach the knowledge and skills necessary for collaboration. Because IPE programs are tailored to meet the needs of the participants and the resources available, there is great diversity amongst programs in regards to participants, length, content taught, delivery mode, outcomes and results (Reeves, 2012).

While IPE programs are widely implemented in traditional healthcare settings, high functioning teams are needed in other, less common, practice settings. In school healthcare, students-athlete's care falls to school nurses (SN) and athletic trainers (AT). ATs and SNs face additional challenges to ICP such as working at different times, differing employment institutions and separate education. Therefore, the purpose of this study is to develop and implement an IPE program tailored to meet the unique needs of ATs and SNs.

2 METHODS

This pilot study is mixed method, quasi-experimental, one-group pretest-posttest design. This study will utilize a convenience sample of ATs and SNs currently employed in one district in the state of Virginia. No exclusion(s) will be made from this group.

A variety of instruments will be used to assess program success. A Roles and Responsibility Knowledge Survey (RRKS-SN/AT) will be used to determine the level of knowledge each participant has of the roles and responsibilities of the other profession. The General Self-Efficacy Scale (GSE) will be used to examine changes in general self-efficacy. The TeamSTEPPS® Teamwork and Attitude Questionnaire (T-TAQ) was designed to assess changes attitudes toward the role of teamwork in the delivery of healthcare. The

System Usability Scale (SUS) assesses participant's response to the usability of the online program. Lastly, the Participant Response Survey (PRS) was designed by the researchers to examine the participant's views of the program. 1 month following the program, a qualitative interview will be used to examine communication changes which may have occurred following the program

The learning content is divided into 4 parts; Roles and Responsibilities, TeamSTEPPS® Team Structure, TeamSTEPPS® Communication and a Simulation. Parts 1-3 are delivered via voiced over PowerPoints and will instruct on the roles and responsibilities of the other profession and the skills necessary for improved team forming and communication. Part 4 involves previously recorded simulations; one poor example of a handoff between a SN and AT and one good example. The participant will be required to view each simulation and identify the errors in the poor example.

Descriptive statistics will be completed for each instrument (overall and each item and separated by profession. Significance level will be set a-priori at $p=0.10$ for all analyses. A Wilcoxon signed-rank test will be used to examine differences in scores pre-and post-learning intervention (dependent variable=RRKS score and independent variable= time). GSE and SUS analysis will compare the present study data to previously reported industry norms (Schwarzer 1995, Bangor 2008). Changes in T-TAQ scores will be compared through a Wilcoxon test (dependent variable=T-TAQ score and independent variable= time). The researchers will highlight areas of interest in the PRS instrument such as the lowest and highest scoring items or the percent agreement of each statement. Phenomenological data analysis will provide the framework for the analysis of the qualitative interviews which involves identifying themes, horizontally coding and then clustering themes to create larger related themes (Moustakas, 1994).

3 RESULTS

The reaction to an IPE program as well as changes to participant attitudes towards teamwork and knowledge of the other profession's role will be measured. In addition, behavioral changes surrounding communication practices will be documented one month following the program.

4 DISCUSSION

Failures in communication and teamwork in the healthcare system present an occasion for numerous adverse events. Because adverse events are physically, psychologically and financially costly, it is salient to prevent them. One suggested avenue to prevent adverse events is through IPE which teaches healthcare providers the importance of forming and maintaining healthcare teams and the skills to do so. An IPE program designed for use in school healthcare could help to ensure nontraditional healthcare providers are provided the opportunity to learn to work interprofessionally and improve patient outcomes.

REFERENCES

- Bangor A, Kortum PT, Miller JT. An Empirical evaluation of the System Usability Scale *International Journal of Human-Computer Interactions*. 2008;24(6):574-594.
- Horsburg M, Lamdin R, Williamson E. Multiprofessional Learning: The Attitudes of Medical, Nursin and Pharmacy Students to Share Learning. *Med Educ*. 2001;35:876-883.
- Kohn LT, Corrigan JM, Donadlson MS. *To Err is Human: Building a Safter Health System*. Washington, D.C.: National Academy of Sciences; 2000.
- Moustakas C. *Phenomenological Research Methods*. Thousand Oaks, California Sage Publications; 1994.
- Reeves S, Tassone M, Parker K, Wagner SJ, Simmons B. Interprofessional Education: An Overview of Key Developments in the Past Three Decades *Work*. 2012;41:233-245.
- Schwarzer R, Jerusalem M. *Generalized Self-Efficacy Scale*. Windsor, England: NFER-NELSON; 1995.

Welsch

ENHANCING SPEECH RECOGNITION IN A VIRTUAL OPERATING ROOM

Levi Warvel
Department of Psychology
Old Dominion University
Norfolk, VA, USA
lwarv001@odu.edu

Mark Scerbo
Department of Psychology
Old Dominion University
Norfolk, VA, USA
mscerbo@odu.edu

Samantha Zybak
Department of Psychology
Old Dominion University
Norfolk, VA, USA
szyba001@odu.edu

Rebecca Kennedy
Department of Psychology
Old Dominion University
Norfolk, VA, USA
rkenn014@odu.edu

Yannis Papelis
Virginia Modeling, Analysis, and Simulation Center
Old Dominion University
1030 University Blvd
Suffolk, VA, USA
ypapelis@odu.edu

Menion Croll
Virginia Modeling, Analysis, and Simulation Center
Old Dominion University
1030 University Blvd
Suffolk, VA, USA
mcroll@odu.edu

Hector Garcia
Virginia Modeling, Analysis, and Simulation Center
Old Dominion University
1030 University Blvd
Suffolk, VA, USA
hgarcia@odu.edu

ABSTRACT

The current study compared current speech recognition program performance in a surgical training simulation relative to earlier benchmarks. Surgical residents performed a simulated laparoscopic cholecystectomy while verbally communicating with virtual members of the operating room. Results indicated that modifications to speech recognition grammar rules produced some new challenges that impacted objective performance measures. However, subjective measures of participant satisfaction with vocal control were more positive than those obtained in earlier design iterations. Findings suggest that the new grammar rules improved user experience but have not yet satisfied intended system performance goals.

Keywords: medical training, surgical simulation, speech recognition, finite state machine

1 INTRODUCTION

Procedural medical simulators have risen in demand over the past decade because of their ability to enhance both patient safety and learning opportunities for medical personnel. Despite the value of procedural medical simulators, most do not address the context in which the skills will be applied which limits the ability for more experienced physicians to learn non-procedural skills, such as critical thinking or communication. To address this limitation, developed a context-relevant medical simulator called the Virtual Operating Room (VOR) that allowed trainees to practice decision-making and interpersonal communication skills in a safe and realistic environment (Baydogan, Belfore, Scerbo, & Saurav, 2009; Scerbo et al., 2007; Papelis et al., 2014). In the VOR, trainees interact with virtual members of a typical OR to gain situational insight or delegate tasks.

To support vocal control in the system, a speech recognition system was developed (Papelis, 2014) using a Dragon Naturally Speaking client and an extended Finite State Automata (FSA) converted from a semantic interpretation grammar model established from the Speech Recognition Grammar Specification (W3C, 2004). The extended FSA included NULL and ANY transitions to make recognition more flexible and detect a wider range of utterance phrasings. However, initial task analyses underrated the variability of surgical terminology used during the training task, limiting speech recognition performance. Additional research and testing indicated that many of the utterances that could occur at each procedural step shared a few common terms unique to that step alone. In addition, the Dragon client often misinterpreted utterances in consistent ways, (e.g., misinterpreting “trocar” as “true car”). Accordingly, we attempted to improve system performance by further extending the FSA to include each observed variation

while reducing the minimum phrase length in most procedural steps to include only the most unique predicted utterances. We predicted that this new extended FSA would improve accuracy in speech recognition and enhance the user's experience.

2 METHOD

Sixteen surgical residents (6 women and 10 men) were recruited from Eastern Virginia Medical School (EVMS) to perform a laparoscopic cholecystectomy (gall bladder removal) in the VOR. Each participant was at least in their second year of residency school and had assisted in, performed, and/or observed a laparoscopic cholecystectomy. The VOR was rendered in a C.A.V.E. on the EVMS campus. Each of the three C.A.V.E. walls featured one of the three virtual team members: circulating nurse on the left, anesthetist on the right, and the attending surgeon directly in front of participant. The physical patient model was placed near the front wall. Trocars and simulated laparoscope were placed prior to each participant's session and remained relatively fixed. The body cavity contained a customized version of the Simulab, Inc. LapTrainer system (Seattle, WA). The LapTrainer system was composed of a hard plastic model of the stomach and liver bed upon which a soft rubber gallbladder was adhered. Each gallbladder was modified to detect cuts of the main structures. Cutting the correct or incorrect structures would inform the state machine and advance the scenario. Except for the modified gallbladder, all system interaction was done verbally via a headset microphone. Each participant first filled out a demographic form, performed the simulated procedure, and completed a brief survey regarding their opinions of the experience.

3 RESULTS

Performance was measured by the number of times speech recognition failed and the researcher needed to manually intervene. Mean number of manual transitions was 2.75 compared to 3.9 found in earlier testing. Participants largely indicated that the virtual team members in the VOR sounded and responded like real team members ($M = 5.4$ out of 7). Thirty-three percent of participants also indicated that open communication with the virtual team members and their ability to respond naturally was the most valuable aspect of the simulation. No participants identified the speech recognition system as the worst part of the simulation in the current build, representing a significant subjective improvement in performance compared to the 20% of participants found previously.

However, 25% of participants found the worst aspect of system performance was a delay in system response, which caused them to be uncertain that the system actually heard them. Some participants felt that the system responded in inappropriate ways when it did recognize their utterances. Many participants considered some procedural steps to be obvious and did not communicate them (e.g., trocar placement, asking for lights to be turned down). One procedural step, the artery/duct clip step, was overly simplified by the new FSA and was observed to be incorrectly triggered in some cases. In addition, participants produced new sets of utterances not previously observed by researchers nor predicted by SMEs in interviews.

4 DISCUSSION

The overall performance of the speech recognition system demonstrated a general improvement over the earlier version of the VOR. The number of manual transitions necessary in the current version of the system was lower than in the older version. Additionally, the surgeon's ratings indicated that the interactive vocal features of the VOR were the most valuable aspects of the simulation. However, the system also demonstrated some limitations. Unexpected utterances and a general lack of responses to obvious procedural steps accounted for most speech recognition failures. Lag time between participant input and system response was also significant enough to cause participants to doubt that the system had indeed heard them.

Although the FSA modifications did improve the ability of the VOR to automatically transition between procedural steps, further modifications are necessary to establish acceptable system performance. Certain non-critical procedural steps should be reviewed for relevance or automated. Additionally, further surgical observation and SME interviews are necessary to develop a robust glossary of potentially relevant terminology to allow the VOR to detect all reasonable user utterances.

REFERENCES

- Baydogan, E., Belfore, L. A., Scerbo, M.W., & Saurav, M. (2009). Virtual operating room team training via computer-based agents. *International Journal of Intelligent Control and Systems*, 14, 115-122.
- Papelis, Y.E., Croll, M., Garcia, H., Scerbo, M.W., & Kennedy, R. (2014). Behavior authoring and run-time management of computer agents for a virtual operating room training environment. *MODSIM World 2014*, Paper No. 1454, (pp. 1-7). Hampton, VA: MODSIM World.
- Scerbo, M.W., Belfore, L.A., Garcia, H.M., Weireter, L.J., Jackson, M.W., Nalu, A., Baydogan, E., Bliss, J.P., & Seevinck, J. (2007). A Virtual operating room for context-relevant training. *Proceedings of the Human Factors & Ergonomics Society 51st Annual Meeting* (pp. 507-511). Santa Monica, CA: Human Factors & Ergonomics Society.
- W3C 2004, Speech Recognition Grammar specification Version 1.0, Retrieved March 2014 from: www.w3.org/TR/speech-grammar.

THE TEACHLIVE SIMULATOR AS A TOOL FOR LEARNING VERSUS TOOL FOR ASSESSMENT OF TEACHING AND CLASSROOM MANAGEMENT SKILLS

Joshua G Stubbs
Old Dominion University
jgstubbs@odu.edu

Dr. Ginger Watson
Old Dominion University
gswatson@odu.edu

ABSTRACT

TeachLivE, a classroom teaching simulator developed by the University of Central Florida and used at more than forty universities across the United States, was designed to allow pre-service teachers practice teaching and classroom management skills in a practical setting. As more research is undertaken and published regarding this emerging simulator, several studies have used TeachLivE as a tool for formal assessment rather than for learning and practice of skills. Given that TeachLivE was not inherently designed for use in assessment, further investigation into both its effectiveness in this way of use and comparison to its effectiveness as a practice tool is needed.

Author Keywords

TeachLivE; Learning vs. Assessment; Pre-service Teacher Training; Teaching Simulators

INTRODUCTION

The developers of TeachLivE have undertaken multiple studies investigating TeachLivE purely as a training tool; more specifically, a tool in which skills are practiced and developed through immediate feedback and multiple attempts in the simulator (Dieker, Rodriguez, Lignugaris, Hynes, & Hughes, 2013; Dieker, Straub, Hughes, Hynes, & Hardin, 2014).

Copyright to the Work and to any supplemental files integral to the Work which are submitted with it for publication (such as an extended proof, a PowerPoint outline, or appendices that may exceed a printed page limit), including without limitation, the right to publish the Work in whole or in part in any and all forms of media, now or hereafter known, is hereby transferred to the ACM effective as of the date of this agreement, on the understanding that the Work has been or will be accepted for publication

by ACM.

More recently, however, a group of researchers have pushed to consider TeachLivE beyond its use as purely a training tool, arguing that the simulator has substantial promise as an assessment tool. They discuss the use of TeachLivE not only in the assessment of pre-service teachers' skill level, but as a

valid part of in-service evaluations of teachers for the purposes of relicensing and in lieu of in-classroom observations (Barmaki, 2014; Kaufman & Ireland, 2015).

WHAT IS TEACHLIVE?

TeachLivE is a teaching and classroom management simulator developed by the University of Central Florida (UCF). The simulator is physically made up of a projector screen on which the interface is displayed, a first-generation Xbox 360 Kinect motion capture device, speakers placed around the room in which the simulator is housed, and a microphone for the preservice teacher to interact with the interface with. The interface is made up of a classroom, which can be an elementary, middle, or high school classroom depending on the needs of the user. Within the classroom, there are five student avatars for the user to interact with.

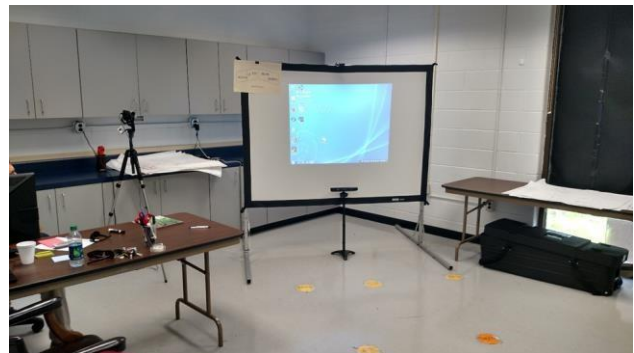


Figure 1. A TeachLivE Lab

Unlike most simulations, the avatars are not entirely A.I.-driven, but are controlled live from a central lab at UCF, which allows for significant customization of classroom management scenarios and provides the sole feedback during the session in the lab.

TEACHLIVE AS A TOOL FOR LEARNING

TeachLivE was originally designed to fill “a gap in teacher education instruction where teacher candidates and struggling teachers can rehearse their skills, improve their skills, and build confidence in their abilities” (Dieker et al., 2013, P.29). Further, it has been considered as a “safe alternative” to putting preservice teachers in a field experience where a loss of valuable student learning time could occur if the novice fails. The simulator allows the teacher-in-training to test and practice specific skills in classroom management and content teaching, with opportunities to pause, reflect, and make adjustments with no real-world consequences (Dieker et al., 2014).

Learning and training practices with the TeachLivE simulator has been found to increase the frequency of higher-order questions and specific feedback towards students by the novice teacher when practicing content teaching (Straub, Dieker, Hynes, & Hughes, 2013). Whitten, Enicks, Wallace, and Morgan (2013) found a growth in effective teaching skills in the field following the use of TeachLivE, but noted that it could not be assumed that this was caused by use in the simulator. Dawson (2016) found that after engaging in regular training sessions, pre-service teachers demonstrated high aptitude in the areas of student error correction and specific feedback. Further, the pre-service teachers showed a significant level of transfer of performance to authentic classroom settings, even over time. It is important to note that these findings were not based upon assessment during use of the simulator, but rather using other assessment tools or even in-field observations.

Hayes, Hardin, & Hughes (2013) noted that there was high amount of skills transfer from training with TeachLivE to the actual classroom. Dawson, Lignugaris & Kraft (2013) noted that teachers demonstrated high rates of transfer on the skills they practiced in TeachLivE in an actual classroom, particularly compared to lower-fidelity simulations, such as role-play exercises.

TEACHLIVE AS A TOOL FOR ASSESSMENT

As early as 2013, research identified the simulation as a potential assessment tool, noting that TeachLivE’s assessment properties may actually be more valuable than its skill practice and learning properties (Sander, 2013). Informal assessment, though immediate feedback, had always been a part of the learning process with TeachLivE, but the findings of Sanders and Kaufman and Ireland (2015) pushed towards a more formal assessment process.

Kaufman and Ireland have heavily pushed for the use of the TeachLivE simulator, noting that the use of similar simulations in other fields are already used to assess learners. However, they fall short of establishing a formal framework for assessment, or even identifying one that may have been established by other research, meaning that while they see the potential of TeachLivE beyond a tool for learning and practice, there is no methodologically and psychometrically sound assessment in place for the use of TeachLivE in this way. Further, they are unclear about even what learning outcomes and/or possessed skills could, in theory, be properly assessed through the use of TeachLivE. They note that it could be used in lieu of in-classroom observation or as professional development, but fail to go into further detail.

Barmarki (2014) does use TeachLivE as an assessment tool with an instrument she developed, but the instrument was not designed to gauge learning outcomes or skill level, rather to gauge body language in certain teaching situations. However, it remains one of the very few empirical uses of TeachLivE as an assessment tool at this point, particularly with a dedicated instrument of assessment developed. Another study that attempted to use TeachLivE as an assessment (Vince Garland, Vasquez III, & Pearl, 2012) did not develop a dedicated instrument at all, but rather used a preexisting instrument meant for in-field observations.

LEARNING TOOL VS. ASSESSMENT TOOL

The majority of the literature that exists regarding TeachLivE discusses the simulator in terms of being a learning tool, rather than an assessment tool. According to the developers themselves, use as a learning tool was the purpose for design in the first place (Dieker et al., 2013). However, as Kaufman and Ireland state, TeachLivE would certainly not be the first learning and practice simulator that also works effectively as a tool

for assessment, as it is common with simulators in the medical and transportation fields (2015).

However, to be effectively used as a tool for assessment, TeachLivE needs to be empirically validated for such use. It must be reliable in measuring skills and practical in terms of cost and logistics. Until that point, the argument for use of TeachLivE as an assessment tool may simply be premature.

BUILDING A RESEARCH FRAMEWORK FOR TEACHLIVE

In order to empirically validate TeachLivE for the purpose of assessment, an entire research framework must be further developed for TeachLivE generally. Very little third party research exists regarding the simulator thus far. More specifically, in order to build the framework, it is important to identify or establish connections between the simulator and currently existing theoretical models and classroom management performance metrics. Vince Garland et al. did attempt this in their 2012 study, using pre-existing observational assessment checklists to evaluate actions during the use of TeachLivE, but it has not been expanded or revisited beyond that study.

Research evaluating performance transfer between tasks, in this case between classroom management scenarios occurring in the use of the simulator and actual in-field, in-classroom classroom management scenarios, has been undertaken in a few cases, but generally with only a particular type of classroom (such as special education classrooms) or in small, exploratory studies (Dawson, et, al., 2013; Harshman, 2012). In both cases, the results are not generalizable. Larger, but specifically focused, research on performance transfer, would help inform the learning side of the learning tool vs. assessment tool discussion.

Another important step in developing a full research framework is the identification of concurrent validity evidence that both good and poor performance can be measured using TeachLive. This step is particularly important because it provides a direct theoretical foundation to the assertion that TeachLivE can be used for assessment.

Finally, basic research on fidelity and immersion regarding TeachLivE, and user perceptions of the same, would provide valuable information on its own, but combined with the rest of the proposed framework,

would help in the development of tasks and scenarios that are a) appropriate transfer-of-training tasks, or b) appropriate assessment tasks. A small study regarding user perceptions of immersion does exist, but more robust research is needed in this area (Stubbs, Baker, and Watson, 2016).

A PATH FORWARD

The establishment of a research framework for TeachLivE would benefit both the validation of the simulator as a learning and/or assessment tool and the development of classroom teaching simulators generally. Until this framework is established, however, it is difficult to determine TeachLivE's effectiveness as either a learning or assessment tool, let alone compare the two to determine best use.

REFERENCES

1. Adcock, A. B., Watson, G. S., Morrison, G. R., & Belfore, L. A. (2011). Effective Knowledge Development in Game-Based Learning Environments: Considering Research in Cognitive Processes and Simulation Design.
2. Barmaki, R. (2014, July). Nonverbal communication and teaching performance. In *Educational Data Mining 2014*.
3. Dawson, Melanie Rees, "From TeachLivE™ to the Classroom: Building Preservice Special Educators' Proficiency with Essential Teaching Skills" (2016). *All Graduate Theses and Dissertations*. Paper 4930.
4. Dawson, M. Lignugaris & Kraft, B.(2013). TLE TeachLivE™ vs. role-play: Comparative effects on special educators' acquisition of basic teaching skills. In A. Hayes, S. Hardin, L. Dieker, C. Hughes, M. Hynes, & C. Straub. *Conference Proceedings for First National TeachLivE Conference. Paper presented at First National TeachLivE Conference: Orlando, FL, University of Central Florida*.
5. Dieker, L. A., Rodriguez, J. A., Lignugaris, B., Hynes, M. C., & Hughes, C. E. (2013). The potential of simulated environments in teacher education: current and future possibilities. *Teacher Education and Special Education: The Journal of the Teacher Education Division of*

- the Council for Exceptional Children*, 37(1), 21-33.
6. Dieker, L. A., Straub, C. L., Hughes, C. E., Hynes, M. C., & Hardin, S. (2014). Learning from virtual students. *Educational Leadership*, 71(8), 54-58.
 7. Harshman, H. E. K. (2012). The influence of TeachLivE on anxiety levels in preservice and in service mathematics teachers. *Proceedings from Ludic Convergence*, 15.
 8. Hayes, A., Hardin, S., & Hughes, C. E. (2013). Perceived presence's role on learning outcomes in a mixed reality classroom of simulated students. Paper presented at the *Human Computer Interaction International*, Las Vegas, Nevada.
 9. Kaufman, D., & Ireland, A. (2015). The potential of simulation for teacher assessment. *The Complexity of Hiring, Supporting, and Retaining New Teachers Across Canada*, 113.
 10. Kaufman, D., & Ireland, A. (2016). Enhancing teacher education with simulations. *TechTrends*, 60(3), 260-267.
 11. Norman, G., Dore, K., & Grierson, L. (2012). The minimal relationship between simulation fidelity and transfer of learning. *Medical Education*, 46(7), 636-647.
 12. Sander, S. Exploring the Impact of Virtual Classroom Technology on Learning to Teach. *Proceedings from Ludic Convergence*, 29.
 13. Straub, C., Dieker, L., Hynes, M., & Hughes, C. TeachLivE National Research Project.
 14. Stubbs, J., Baker, P., & Watson G, (2015). Perceived Fidelity, Workload, and Cognitive Load of Pre-Service Teachers using TeachLivE. Unpublished Manuscript.
 14. Vince Garland, K., Vasquez III, E., & Pearl, C. (2012). Efficacy of individualized clinical coaching in a virtual reality classroom for increasing teachers' fidelity of implementation of discrete trial teaching. *Education and Training in Autism and Developmental Disabilities*, 47(4), 502.
 15. Whitten, E., Enicks, A., Wallace, L., & Morgan, D. (2013). Study of a mixed reality virtual environment used to increase teacher effectiveness in a pre-service preparation program. *Proceedings from TeachLivE Conference 2013*, 38.

A System Dynamic Modeling of Iranian women Movement

Khadijeh Salimi

Old Dominion University

Ksali001@odu.edu

I. Abstract

Iranian women, similar to other women around the globe, took part in several movements during recent centuries to achieve their equal rights. These movements started to coincide with other feminist movements around the world and have continued on to the present day. In each wave of the movements, Iranian women gained some of their rights. These achievements include rights for education, rights to work, and rights to vote, which enable women to participate in politics. But, still, Iranian women suffer from unequal conditions in the 21st century. The objective of this study is to simulate the effect of social medias and advanced technologies on educating Iranian women's through the process of mobilization. To achieve these goals, the study proposes that the main obstacle to have a successful women's movement is the lack of collective gender identity among Iranian women. Thus, to become successful, Iranian women should first educate in the way that they shape gender collective identity, which is particularly possible via social Medias.

II. Author Keywords

Keywords: Women's Movements, Iran, Advanced Technology, Collective gender identity, Social Medias.

III. Introduction

The principal aims of the present study are to explain the important role of technology and social media to educate women, to feel the shared grievance which male dominant structure imposes on them. This study

proposes that Iran's male dominant structures create the ideology of superiority of men and inferiority of women. This leads to citizens, especially women, whose identities are shaped through society and interactions; to perceive themselves in different ways, and do not have a collective gender identity.

A vast number of efforts have been conducted towards Iranian women movement. However, previous studies mostly focus on historical or sociological conditions that led to Iranian women movements, the present research focuses on one of the major obstacles which prohibit a successful women movement, lacking collective gender identity, and offers a solution, educating women, for this existing obstacle via advanced technology. Social medias are crucial tool to educate citizens and help women to have a shared identity, which mostly is based on the feeling of grievance, across a society which is the fundamental element of social movement.

What is collective identity?

Collective identity is different from individual identity. To elaborate, I will define each of these identities. Identity is the thing that all individuals have, seek, and construct through public or social interaction, thus it varies in different contexts. Charles Tilly¹ defines it as people's experiences of different roles in different categories or groups. Categories help people to realize how they and their activities are similar to or different from others. These similarities lead to the shaping of collective identity among individuals. It is important to know how understanding oneself through society and interaction with other people who have the same feelings of shame or

¹ Tilly, Charles. (1996). *Citizenship, Identity, and Social History*. Cambridge: Cambridge University Press.

grievance connects people together and shapes their collective identity.

For the purpose of this paper, I will define women's collective identity as a shared understanding and senses of grievances based on unequal conditions in society, which provide the element of similarity that gathers women together as a group. This collective identity is an identity which is found among the most social movements, causing a group of people to feel an interconnection with or affinity for one another. Members of the group become aware of this affinity when they interact through their social environment. Scholars who study these groups believe that the feeling of grievances, which is mostly ignored, is a leading factor in shaping collective identity and furthering autonomous social movements.

IV. Method and Material

This study covers the Iranian women's movement. Data was collected through library study and needed information was gathered from books, articles, and scientific notes. System dynamic modeling is used to simulate how advanced technology facilitates and fastens the process of educating people to join a shared community. System dynamic is a mathematical modeling which is used to study the behavior of complex systems. It is composed of stocks, feedback loops, and flows.

For the purpose of this paper, in the first stage of the molding, we use the Bass Diffusion Model to analyze the behavior of the society's population who are joined to the educated population. In this step, we suppose that there is no opposition party (Fig.1a). Flows and stocks' diagram is developed as (Fig.1b).

In the flow and stock diagram, there are two level/stock variables which show the potential

population and educated population. The potential population can be educated through social medias such as Facebook, Telegram, WhatsApp, etc. We called This way of educating (factor) as social network process in our model. Second, they can be educated via direct interaction with activists and leaders or via gathering information through different academic sources. This group is called leaders or activists who usually are connected with NGOs.

There are three feedback loops here, two reinforcing and one balancing. The reinforcing loops increase the number of educated people while the balancing loop controls the growth rate of the educated people by controlling the potential population growth.

The model's components are defined as follow:

Potential population= Enter- Growth (person)

Educated population=Growth (person)

Social networking=educated population* social network IF* Networking rate*(potential population/ total population) (person)

Self-education/Activist= Resource IF*(Resource investment/Availability) * potential population

Enter= growth Rate*Potential population (person)

Education growth= Social networking + Self-education/Activist (person)

Resource IF=(0.001/ social political barriers (unit less))* educated population

Opposite barriers IF= 0.1 / (educated population/ total population) (unit less)

Socio-political barriers= Opposition parties' investment* Opposite barriers IF (unit less)

Opposition parties' investment=100

Total population= 80000000 (person)

Growth Rate= 0.02 (d (person)/dt)

Resource investment/availability= 1000 (unit less)

Networking Rate= 1000 (person)

Social network IF=0.015

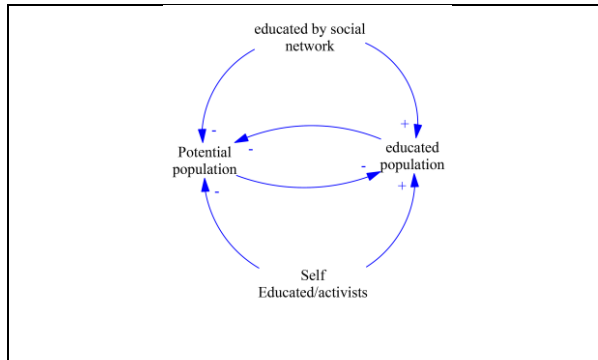


Fig.1a) Causal loop of Education Diffusion- Without Opposition Party

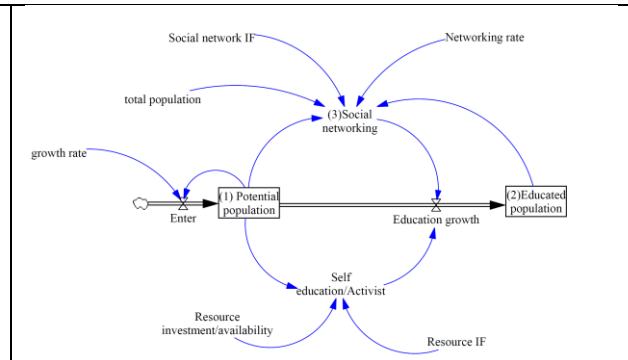


Fig.1b) Stocks and Flow Diagram of Education Diffusion- Without Opposition Party

It should be noted that we have not accessed to empirical data for this study. Because study of Iranian women's situations and movements poses three main problems. First, women are ignored in these kinds of societies, hence accessing enough data and information about their situation is difficult. The second problem, which also originates from the first problem, is lacking enough relevant and primary sources. Third, political reasons make the government manipulate official statistic about women. Despite all these problems, a study on Iranian women's rights is crucial, as they are continually marginalized in the Islamic Republic; indeed, they do not even have any official tribune by which to denounce their suppression.

V. Result and Finding

Iran is a traditional society where many citizens are heavily influenced by the dominant perceptions of traditional masculinity, that often individuals cannot understand their suppression or cannot connect with other who suffer from similar grievances due to their inferior status in society.² This is one of the major problems that Iranian women encounter during their movement. It is important to know how understanding one's self through society and interaction with other people who have the same feelings of shame or grievance brings people together and shapes their collective identity. In other words, while some Iranian women suffer from some inequality in the society, some other not only do not suffer but also admit it as a legitimate law which should be obeyed.

Knowing the problems of Iranian women's movements, lack of collective identity, we use technology and social media to cope this problem. According to Jackson the best way to cope with this problem is to train citizens, and raise their understanding of their status, etc.

Knowing the factors which some of them facilitate and others prohibit the process of educating women, I define three different scenarios to the model:

Scenario 1: In the first scenario, we assume that we already have a population of people who are educated and also the potential population who will be educated through the process of educating. In this way, we can see that all the growth in the educated population is convergent with the change in social networking availability (Fig.2a). If the availability and the Variety of the social networking increase over time, more people will be educated. It means that networking rate and also social network impact factor will increase, then the population of educated people will increase significantly. In other words, the time that our society need to educate people is decreased sharply (Fig.2b). This decreasing in time is very important. As Earl and Kimport³ mention, time is one of the important cost to movement participation.

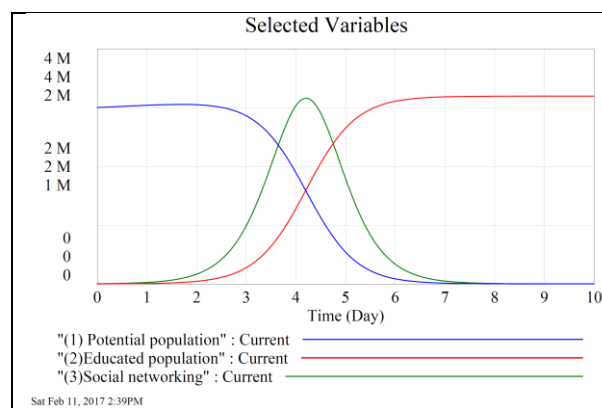


Fig.2a) Result of Basic Model (Social Networking IF= 0.03, Contact rate:1000)

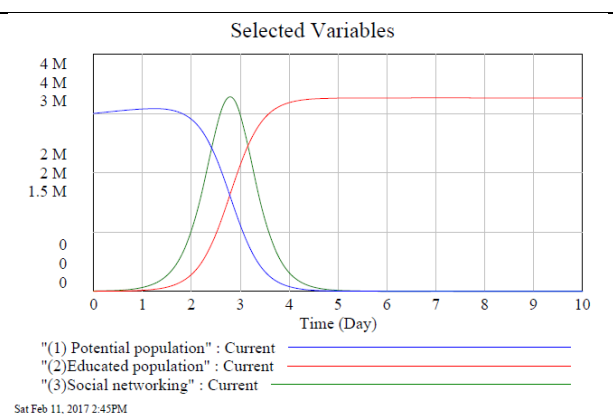


Fig.2b) Result after Change in Parameters (Social Networking IF= 0.05, contact rate:2000)

² Jackson, Peter. (1991). *the Cultural Politics of Masculinity: Towards a Social Geography*. Transactions of the Institute of British Geographers, New Series, 16, No. 2, pp 199-213.

³ Jennifer Earl & Katrina Kimport. (2011) *Digitally Enabled Social Change: Activism in the Internet Age. Taking Action on the Cheap*. The MIT press.

Scenario 2: In the second scenario, the model analyzes the effect of the activist's population growth on educated population growth. Considering the availability of the academic (conferences, seminars, books, publications, etc.) resources increase (Fig.3a).

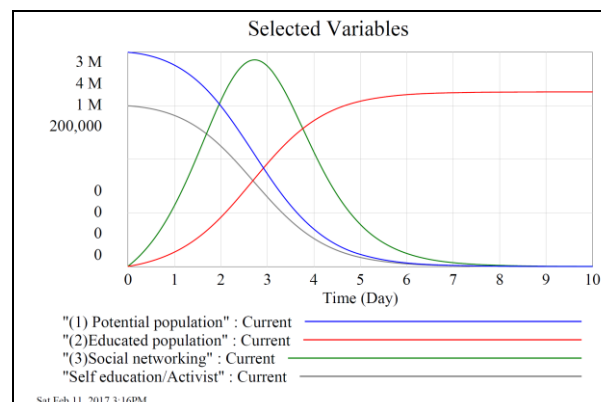


Fig.3a) Result after Change in Parameters (Resource Investment=10, Resource IF=0.005)

If the investment in the available resources for the academic resources and also the impact factor of these resources increases, a more potential population will be educated in the shorter time. (Fig.3b).

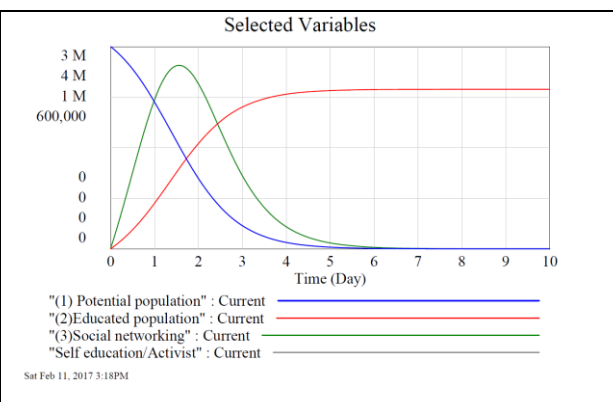


Fig.3b) Result after Change in Parameters (Resource Investment=100, Resource IF=0.01)

However, the case study, Iran, is a country which violated human rights. Many activists such as Narges Mohammadi, Bahare Hedayat, Maryam Bahraman, etc. all are examples of activists who was arrested by the government, and either in prison or do not have a right to public activity. In this atmosphere that traditional activism which Gladwell⁴ called it “high-risk-activism” are almost impossible, the crucial role of advanced technology, social medias, and digital activism are obvious. First of all, they give the opportunity to people to communicate and connect to each other. Second, it facilitates and fastens the process of educating people to a shared idea. Through movies, advertisement, and social web pages we see how feminism attempt to educate people. Social medias help Iranian feminist especially outside Iran and enable them to use advanced technology to cope with this problem. Several NGO Such as Tavana (<https://tavaana.org/en>), Etehad Baraye Iran (<https://www.facebook.com/EtehadBarayeIran>) is offered online free courses to educate Iranian women to understand their right and uncover the falsifying of the dominant ideology.

Scenario 3: in this scenario, we assume that there are two main socio-political barriers. First is Mass Media, which is controlled by an Iranian government and

colonized population minds. This mass media plays an important role to create a dominant culture or ideology. Second, government executive branches such as cops or other forces, who arrest activist and canceled conferences, create a closed atmosphere, etc., (Fig.4a). These factors affect the power and availability of the academic resources for activist and leaders. Thus, these barriers are obstacles and prevent the growth of educated people through academic resources and activists. However, the important role of social medias such as Telegram, Viber, Instagram, etc. will emerge at this time. These networks lead to increase in educated population rate. As a result, the growth in educated population will decrease the effect of socio-political barriers over the time. Shirky⁵ called this as a process which strengthens civil society and the public sphere (Fig.4b). Thus, the balance of power between state and civil society will be shifted in a peaceful way. In addition, Iranian government ability to use executive branches will be weakened and we have a stronger public sphere. Based on the result which is shown in (Fig.4c) we realize that social medias are 10 times more effective than activists in the society such as Iran. While activist educates 3500 persons, social medias are able to educate 3 million and a half individuals. This is really what Boyd⁶ asserts as web and social medias make a process

⁴ Malcolm, Gladwell. (2010). *Small Change*, The New Yorker.

⁵ Clay, Shirky. (2011). *The political power of social Media: Technology, the Public Sphere, and Political*

Change, Published by the council on Foreign Relations.

⁶ Andrew Boyd. (2003). The Nation. The Web Rewires the Movement.

cheaper and quicker. Cheaper here at least in the sense that there is no threat to people lives.

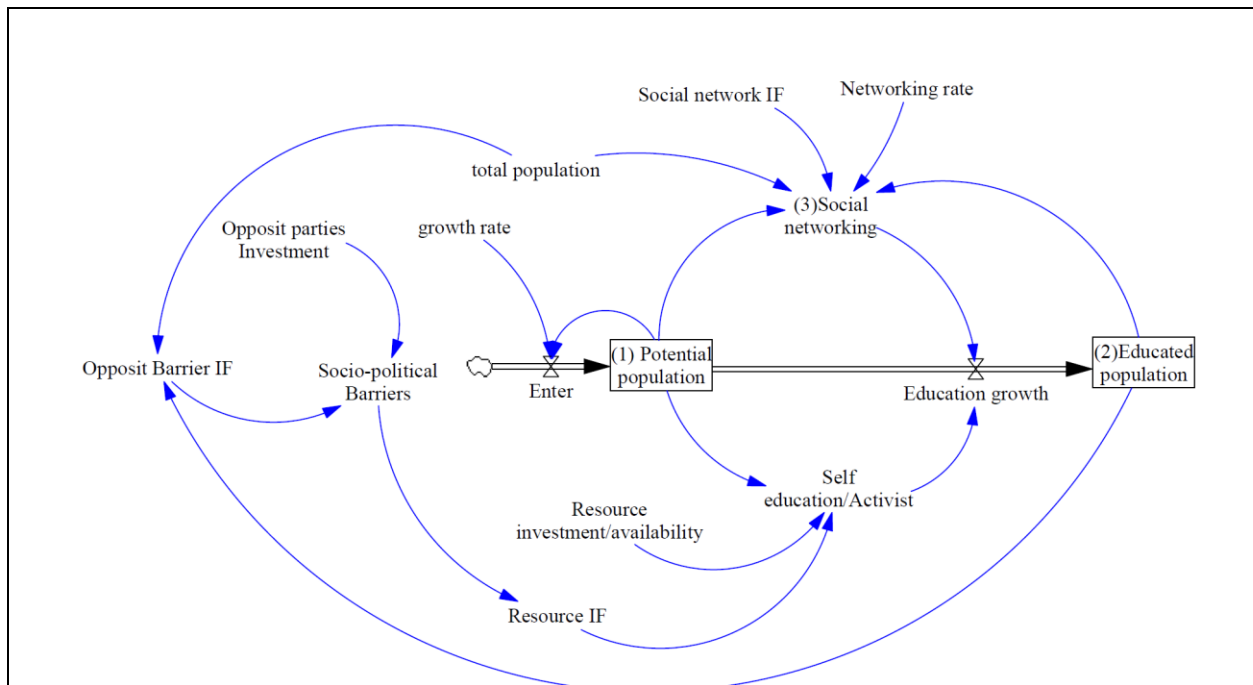
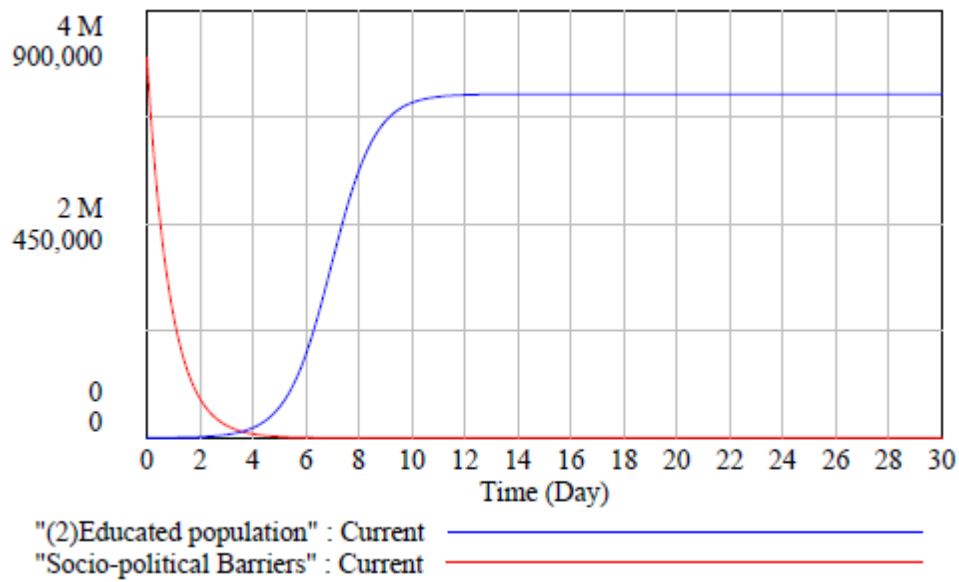


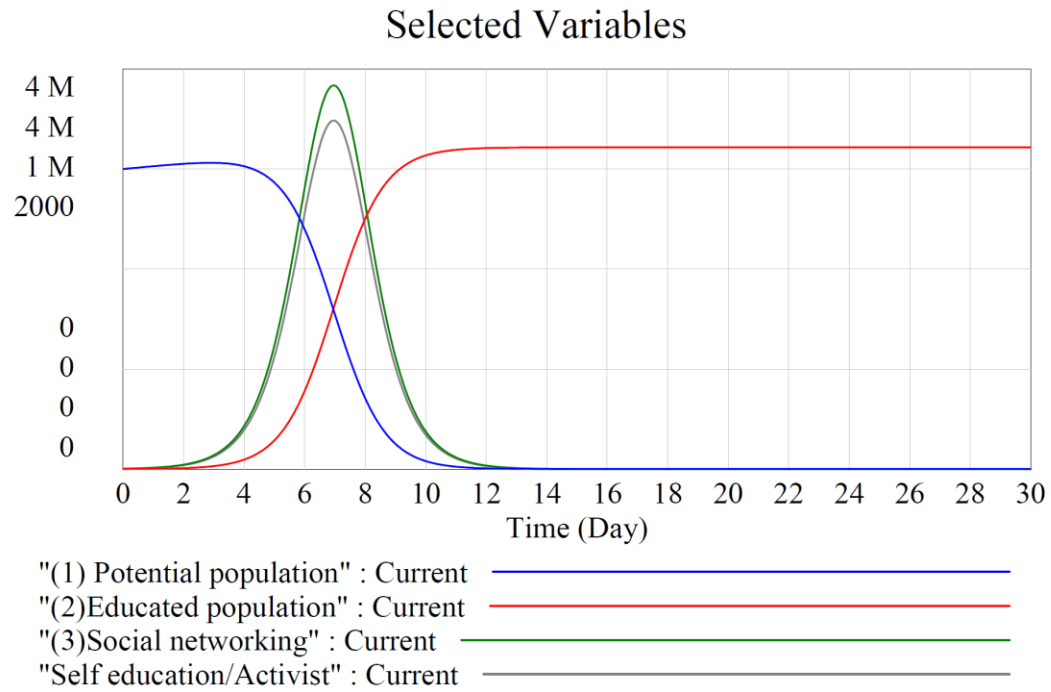
Fig.4a) Model Including the Socio-political Barriers

Selected Variables



Sat Feb 11, 2017 4:26PM

Fig.4b) the changes in the socio-political barriers and educated population



Sat Feb 11, 2017 4:27PM

Fig.4c) the behavior of variables for the scenario 3

VI. Conclusion

This paper considers the lack of gender collective identity among Iranian women, as the major reason for Iranian women movement defeat. And this collective identity is a leading factor in the succession of women's movements. Thus, advanced technologies provide plenty of social networks across societies, which is very hard to control and censure, and enable feminist to educate women in Iran in order to realize their own rights and to come to have a shared identity. Social networks have a significant role in a country like Iran. As simulation results show, in countries where human right activists are under restricting social controls, even a small percentage of increase in the effective rate of social networks have a huge effect on the rate of education of women and increases this rate drastically. Also, model result's shows that under suppression environment, social networking can educate women 10 times more than academic resources. It also shows that social networks fasten the process of educating. This study mainly considers the role of social medias in the process of educating Iranian women, but the role of advanced technologies is not limited to this function; they can also use as an efficient tool in the process of mobilization. Hence, it is recommended to study women movements from these perspectives.

VII. References

1. Boyd, Andrew. (2003). The Nation. The Web Rewires the Movement.
2. Earl, Jennifer & Kimport, Katrina . (2011) Digitally Enabled Social Change: Activism in the Internet Age. *Taking Action on the Cheap*. The MIT press.
3. Gladwell, Malcolm. (2010). *Small Change*, The New Yorker.
4. Jackson, Peter. (1991). *the Cultural Politics of Masculinity: Towards a Social Geography*. Transactions of the Institute of British Geographers, New Series, 16, No. 2, pp 199-213.
5. Shirky, Clay. (2011). *The political power of social Media: Technology, the Public Sphere, and Political Change*, Published by the council on Foreign Relations.
6. Tilly, Charles. (1996). *Citizenship, Identity, and Social History*. Cambridge: Cambridge University Press.

STUDENT CAPSTONE CONFERENCE

2017

BUSINESS & INDUSTRY

- Page 31 Felicia Grey
Old Dominion University
Why do Member Countries Choose not to Participate in the World Trade Organization's Dispute Settlement Body?
- Page 42 Daniel Perez
Old Dominion University
A System Dynamics Approach to Predict the Trend Of Superhero Movies Overtime
- Page 57 Paul Delimarschi, Jonathan Griffith, Mitchel Howard, Sean McBryde and Gene Lesinki
U.S. Military Academy West Point
Simulation and Analysis of the Aircraft Corrosion Control Facility at the Corpus Christi Army Depot

WHY DO MEMBER COUNTRIES CHOOSE NOT TO PARTICIPATE IN THE WORLD TRADE ORGANIZATION'S DISPUTE SETTLEMENT BODY?

ABSTRACT

This study examines why countries would join the World Trade Organization (WTO), but do not use its Dispute Settlement Body (DSB) if a trade dispute arises. To test this expectation, the paper uses an extensive form game with complete and perfect information, costs for litigation (delta), and Prisoner's Dilemma-like payoffs. It finds that with the same litigation, there is pure subgame perfect Nash equilibrium where both states will protect and avoid filing. Free trade and no DSB is also an equilibrium solution. If one state has a higher delta, its trading partner will protect as its dominant strategy. The affected state is then forced to maintain free trade and avoid the DSB, or respond with protectionism and then acquiesce since it cannot afford the full litigation process. These findings, however, do not capture third party litigation; neither do they consider what may happen if a country is unable to retaliate after WTO approved sanctions.

Keywords: World Trade Organization, delta, litigation costs.

1. INTRODUCTION

The World Trade Organization (WTO) is the main international framework for regulating trade among countries. Its main function is to ensure that "trade flows as smoothly, predictably and freely as possible." (WTO 2017). As a result, it administers WTO trade agreements, provides a forum for trade negotiations, handles trade disputes, monitors national trade assistance and training for developing countries, and cooperates with other international organizations. There are three main trade remedies that are available to all members of the WTO as recourse for disputes over trade in goods. These include countervailing duties, safeguards and antidumping, of which, antidumping has been the most frequently used remedy. An aggrieved party may therefore use the Dispute Settlement Body (DSB) which evokes a process of consultation, adjudication, and implementation to get redress for trade violations. (Michalopoulos 2002).

In an ideal world, all trading partners avoid protectionism and engage in free trade. In reality however, this is not always true. Countries frequently flout the principles and provisions to which they have agreed to be bound. Some affected parties are able to unilaterally retaliate, others find recourse through bilateral and regional arrangements, while some find reprieve through case settlement at the WTO. Since the WTO seeks to promote and facilitate free trade, it is useful to explore if the presence of a Dispute Settlement Body inhibits states' inclination to cheat.

Although the Dispute Settlement Body (DSB) is available to all WTO members, utilization is very uneven. Many states have never utilized the body. This failure to use the DSB could mean any of several things. It could reflect the effective functioning of this institution – perhaps the mere threat of the DSB is enough to deter most violations of WTO rules. On the other hand, it could reflect important

inequalities in access to this institution. It could also mean that wealthier and more powerful nations may be much able to take advantage of this institution. Understanding the puzzle of institutional membership but not participation, may therefore shed light of the deliberations that countries make whenever a trade dispute emerges.

1.1 Research Question

The general question that paper intends to model is: Why do member countries choose not to participate in the WTO's Dispute Settlement Body?

Some secondary questions that will be answered include:

1. What effect does participation or nonparticipation have on states' trading relations?
2. Does the DSB create opportunities for trading partners to exploit members?
3. Does the DSB mitigate defection between trading partners with asymmetric interests?

1.2 Relevance

There is much debate about institutions and their role in international affairs. Robert Keohane, for example, explicates the assumptions of neoliberal institutionalism. Here, the positive-sum logic of neoliberalism is advanced. Multilateral institutions arguably cause voluntary cooperation, which in turn effectuates utility gains for each cooperating state or government. (Keohane 1984). Realism proponents like Waltz, Grieco, Mastaduno and Mearsheimer however, attack these tenets. In their estimation, "cooperation under anarchy" is problematic because decentralized enforcement, national interests, and relative gains impede the efficiency of institutions. (Gruber 2000, p. 18-32). Multilateralism supporters point to the general membership and success of the WTO as evidence for their theory. This optimism however, has been countered by the seemingly disparities in how developed and developing countries use the DSB for trade recourse. The paucity of cases from developing countries suggests that the system may be inherently biased against them and so they are to some extent disenfranchised. Scholars who explore the extent to which the DSB functions in satisfying the needs of its developing country Members highlight power asymmetries, initiation and retaliation costs, start-up expenses and low domestic, institutional capacity as possible impeding factors. (Guzman & Simmons 2005; Bown & Hoekman 2005; Davis & Bermeo 2009).

Every conflict within the DSB is fundamentally a dyadic / relational grievance. Highlighting solely the variance in usage between developing and developed countries is therefore intellectually myopic. If participation in this mechanism is taken as the dependent variable and power asymmetries an independent variable, then there is an implicit assumption that trade violations follow only a unidirectional path. This reasoning takes it for granted that only large states are violators, that they exploit weaker states, and that weaker states do not contravene WTO provisions. How then would one account for trade disputes between developing countries and also those between developed ones? Moreover, if economic and institutional capabilities are directly related to a state's tendency to file a dispute, what explains the fact that not all wealthy countries litigate although they may have the ability to do so? Moreover, some affluent nations are more frequent users of the DSB than others. What explains this?

General participation in the DSB is taken as an indication that its provisions are accessible to all its Members. A state's usage as a complainant or a defendant therefore indicates its ability to at least file or respond to a dispute. Participation by itself however, does not account for the calculated opportunity cost of participation versus nonparticipation. In essence, several factors outside of those mentioned may precipitate participation and conversely, nonparticipation, even if the state is able to do so. If the world trading system is an international chess board upon which moves and countermoves are weighed based on preferences, perceived options and payoffs, then participation in the DSB needs to be revisited. Under conventional views of the DSB, nonparticipation could indicate that a state:

- a. Has no trading rights that are being violated.

- b. Has been violated but is unable to file a dispute proceeding.
- c. Has been violated but is fearful that litigation may make it worse off *ex ante*.

Less examined are the possibilities that a violated state may:

- a. Choose not to file although it is able to do so.
- b. Choose not to file because it is fearful of retaliation
- c. File outside of the WTO.
- d. Retaliate.

Many Members of the WTO are simultaneously bound in bilateral and regional arrangements. What therefore explains their choice to proceed with the formal dispute settlement arrangements within the WTO versus informal means, or even the selection of multilateral over bilateral and regional mechanisms and vice versa? Examining this phenomenon may add value to the debate about the (in) efficacy of the WTO generally, and the Dispute Settlement System specifically, since states have other options at their disposal and may therefore choose the one that gives the best payoff at the moment in question. As one tries to make a conclusion about the usefulness of institutions in facilitating cooperation, the example of the DSB which is embedded in an interwoven international trading system may also help to explain state behaviour in other dispute resolution mechanisms and especially their ability to opt in and out at will.

2. STATE OF THE ART

2.1 Current Approaches

Trade is an area of study that interests economic and political scientists. Many types of approaches have therefore been used to understand the dynamics of interstate and international trade. Some of the models that have been used include transmission models, general equilibrium, supply / demand / price adjustment models and discrete decision event algorithm. (Pollins 1982, p. 504-533). An example of a transmission model is Project Link that was spearheaded by the World Bank. It uses a sophisticated trade algorithm to model supply and demand of goods across OECD countries, planned economies, and less developed countries. Where the general equilibrium model is concerned, “complex systems are represented as sets of simultaneous equations that are driven by an objective function such as minimizing prices or maximizing employment.” (Pollins 1982, p. 511). In order to utilize this method, “the values of system variables are derived analytically given the objective function and subject to specified constraints.” (Pollins 1982, p. 511). For the other two types, the supply / demand / price adjustment model uses a market clearing logic that has a trade algorithm, while the discrete decision event algorithm aims to represent dyadic trade flows and incorporates political factors in how these flows are calculated. (Pollins 1982, p. 520).

While several models have been used to represent trade generally, others have been used to examine trade litigation specifically. These include probabilistic models, static, simultaneous games such as Prisoner’s Dilemma and The Battle of the Sexes, and sequential models like the Divide the Dollar game. All of these involve some knowledge or assumption about what the other player might do, the payoffs for those strategies, and the consequent equilibria that can be formed.

2.2 Shortfalls of Current Approaches

All of the current approaches have their usefulness. Their potential shortcomings however, give rise to the use of alternative methods, based on the question that the researcher is trying to answer, as well as the tools and data that are available to him or her. One drawback to the transmission models for example, is that they do not adequately represent dyadic trade because global supply and demand are

pooled. In order to compensate for this, the LINK Project used a dyadic interaction matrix. This, however, does not account for the variation in exports and imports in particular countries over time. Conversely, the general equilibrium models offer a lot of flexibility in their application to issues in international politics. One disadvantage though, is that they are not designed for point-predictive forecasting. Additionally, they have limited utility for state-to-system linkages, and cannot track exchange relations from one state to the next. The supply / demand / price adjustment model on the other hand, is most applicable to primary commodity markets, which restricts its general usage. Discrete event algorithm models also have their setbacks. For instance, although they are widely used, but specific types like TRADER may be underutilized due to the highly technical computations that are needed. (Pollins, 1982, p. 517-525).

The static games also have their shortcomings. These misrepresent the fact that trade is sequential, and that it is also cyclic. This allows for “learning” and updating of beliefs, which might not be captured in a single iteration. Where the sequential ones are concerned, one drawback is that they are sometimes too parsimonious, and do not adequately represent the dynamics of the global trading system.

3. PROPOSED APPROACH

3.1 Method

In order to answer the research question, the paper uses instantiations of an extensive form game with complete and perfect information, penalties for litigation (δ) and Prisoner’s Dilemma-like payoffs in which global welfare is maximized by cooperation at free trade, but individual states have incentives to cheat. The conditions affect the sequential moves that players make and the consequent equilibria that are formed. This game builds on the common model of trade as a Prisoner’s Dilemma and adds the complex interplay of trade at the WTO / DSB. Using an extensive form game, there are two main branches which outline all the possibilities and payoffs that player B must contemplate as it interacts with either a free trading or protectionist State

The structure of the game is outlined in Figure 1-1. The game begins with State A and then State B selecting strategies of Free Trade or Protectionism. The underlying payoffs for this stage match those of the typical prisoner’s dilemma model of trade. Both receive positive payoffs with mutual cooperation at free trade. The lowest payoffs are associated with enacting a free trade policy in the face of protectionism on the part of a trading partner while enacting a protectionist policy when the trading partner has selected free trade yields the highest. Mutual protectionism provides lower payoffs than mutual free trade, but absent the DSB mechanism, this outcome is the unique Nash Equilibrium of the game.

The DSB potentially provides a route out of this trade dilemma. In the model, State A, whose turn it is to move, chooses between filing a case in the Dispute Settlement Body and foregoing that choice. If State A proceeds with the DSB alternative, State B, selects between acquiescence and litigation. If however, State A does not pursue a case, State B may opt to file or avoid doing so. If the DSB is chosen, State A now selects between acquiescing and litigating.

There are twenty possible payoffs, which all fall under the mutual cooperation, mutual defection and “sucker” situations well known in the literature on the Prisoner’s dilemma. There are two additional dynamics, however, to this game. Upon successful litigation, the WTO allows the complainant to impose countermeasures. (Martin 1992, p. 765-792). This is accommodated in the game by making it synonymous to “cheat, cheat.” The thinking is that if a country is protectionist, then the free trader now has the WTO’s permission to also protect its trading entities. There is a cost, however, for all filing and litigation. As a result, a penalty, δ , is used to capture what states have to pay financially as well as audience costs for going to the DSB. It should be noted however, that δ may not be the same for both parties. Some scholars also contend that it is the initial cost for filing that is cumbersome, but this often becomes negligible once a culture of litigation has been formed. (Davis & Bermeo, 2009). It is therefore possible for States A and B to be regular DSB users and this may make δ less of an inhibiting factor.

Exploring delta could therefore illuminate how the DSB functions. In essence this model probes how the different costs of filing regardless of what the other state pays affect the tendency to use the DSB versus engaging in tit-for-tat strategies. Variations of this game can answer the broader question of whether (and when) the WTO's Dispute Settlement Body tempers defection between trading partners with asymmetric interests.

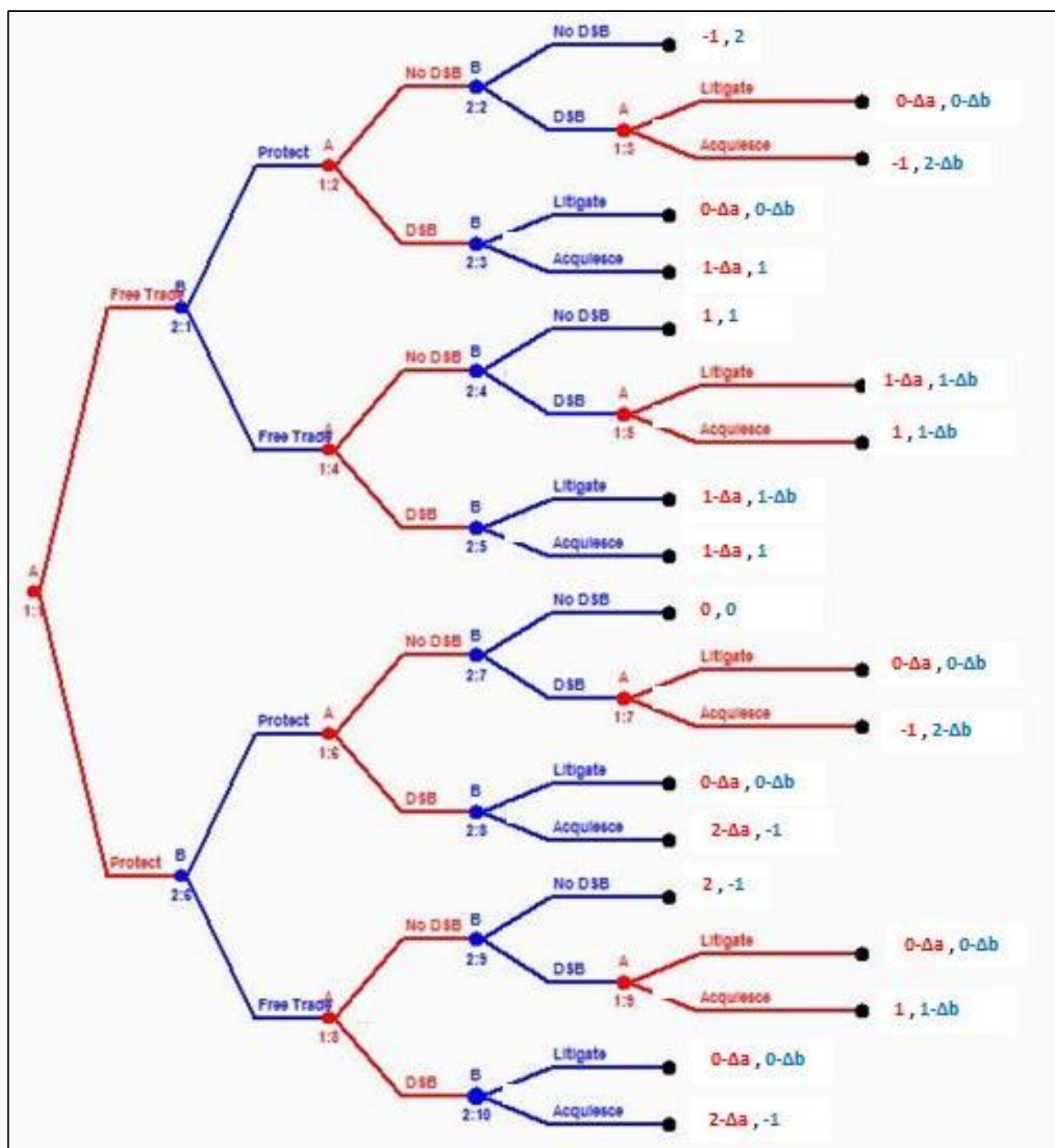


Figure 1.1: A Model of Trade and Dispute Resolution

3.2 Justification of the Approach

Having Prisoner's Dilemma as the premise for this paper is important because it helps to explain players, (two player game), their strategic choices and payoffs, as well as the Nash Equilibrium. Since however, trade is sequential and not simultaneous, an extensive form game better explicates the strategic moves that trading partners make as they interact with each other. In essence, they show that states oscillate between playing free trade and protectionism, evoking the DSB or avoiding it and litigating or acquiescing.

The inclusion of delta, the cost for litigation, is also an important contribution to the model. This is due to the fact that different states pay varied costs when they have to litigate. This in turn, influences whether they will take advantage of institutional reprieve whenever a trade dispute arises, or if they will avoid the DSB and engage in protectionism. Some scholars also contend that it is the initial cost for filing that is cumbersome, but this often becomes negligible once a culture of litigation has been formed. (Davis & Bermeo, 2009). It is therefore possible for States A and B to be regular DSB users and this may make delta less of an inhibiting factor. Exploring delta can therefore illuminate how the DSB functions. In essence this model probes how the different costs of filing regardless of what the other state pays, affect the tendency to use the DSB versus engaging in tit-for-tat strategies. Variations of this model can therefore answer the broader question of non-participatory membership in the WTO's Dispute Settlement Body.

4. MODEL RESULTS

Proposition 1: In a game with symmetric information and the same delta (δ), the presence of the DSB will mitigate the tendency to defect if $1-\delta$, but the state which moves first has an advantage.

Proof:

In an extensive form game, the sequence of operations can be garnered through backward induction. There are two subgame perfect Nash equilibria in pure strategies, one in which the DSB is never used, and another in which the DSB is used regularly by State B.

"Protect, No DSB, Acquiesce; Free Trade, DSB, Acquiesce" is a possible equilibrium. (State A's strategies are listed first and State B's follow). This would give a payoff of $1, 1-\delta$ (In all payoffs, State A's are listed first). What this equilibrium shows is that State A is more inclined to be protectionist. Remarkably, in response to a trade violation by State A, State B's strategy is to still be a free trader. It does, however, seek reprieve at the DSB. Once it reaches that stage, State A acquiesces and both end up better off even with the delta that State B pays for initiating the litigation.

The DSB institution eliminates the mutual-protectionism equilibrium. "Protect, No DSB; Protect, No DSB" is not an equilibrium if $\delta \leq 1$. The payoff for choosing this route is $0, 0$. However, after State A has selected "protect," State B can do better by choosing "free trade" and then bringing a DSB case. As discussed above, the payoff for State B from that strategy is $1-\delta$. Therefore, if delta is less than 1, the mutual protection equilibrium is no longer viable.

The second equilibrium enables State A to be a free trader and for State B to reciprocate accordingly: "Free Trade, No DSB; Free Trade, No DSB." The payoffs for this strategy are $1, 1$. This is a stable equilibrium because neither has an incentive to deviate from this strategy. If State B for example deduces that State A is inclined to be a free trader and avoid using the DSB, then State B could change its strategy from free trade to protectionism. State A however, would take State B to the DSB where B would acquiesce because the payoffs for acquiescence are greater than those for litigation. The resulting outcome would be "Free Trade, DSB, Acquiesce; Protect, DSB, Acquiesce." This would provide payoffs

of **.95**, **1** for States A and B respectively, so B would be no better off selecting Protect, and hence has no incentive to deviate. Moreover, if State A perceives that State B will be protectionist, then it will also play protectionism in the first instance. This would force State B to choose between free trade and protectionism, both of which provide outcomes that are less than the **1**, **1** payoffs.

This, in principle, is how the DSB and all other dispute resolution mechanisms aim to function. The DSB's presence constrains defection if litigation costs are small. Regardless of what delta is, provided it is less than 1 for both parties, an aggrieved party will seek reprieve at the DSB and that the guilty party will surrender. Proposition 1 showed that with $\delta \leq 1$ violators are clearly identified and recourse at the DSB is sought accordingly. Culpable states in turn acquiesce when brought to trial. The equilibrium at mutual protectionism is eliminated, and the remaining equilibria provide both states with payoffs associated with relatively free global trade.

In the context of trade therefore and with Prisoner's Dilemma-like situations very possible and probable, the DSB is a sufficient arrangement for inducing participation and mitigating unfair trading practices. This however, occurs only when costs are not prohibitively expensive. It is therefore prudent to examine how the strategies, payoffs and outcomes change when the conditions under which trade must ensue include an expensive dispute settlement process.

Proposition 2: If dispute settlement costs become too expensive ($\delta \geq 1$), countries will avoid the DSB and engage in protectionism outside the institution.

Proof:

Suppose that A has selected Protect. In this context, the best response by B is to select Protect as well. If B selects protect, the payoff for both players is 0. If B selects Free Trade and then files with the DSB, A will acquiesce, but since $\delta \geq 1$, the payoff for B in this scenario of $1-\delta$ is less than the payoff (0) from simply selecting Protect. Hence, B will not use the DSB. A will also never select Free Trade with $\delta \geq 1$. If A has chosen free trade, and B has then chosen protect, A could choose DSB which would induce B to acquiesce, but the cost of bringing the case to the DSB render this option unappealing. Because $\delta \geq 1$, the payoff $1-\delta$ that A receives is worse than the payoff of 0 associated with selecting Protect initially.

The only subgame perfect Nash Equilibrium is "Protect, No DSB; Protect, No DSB." This gives a payoff of **0**, **0**. A consideration is for the states to play "Free Trade, No DSB; Free Trade, No DSB." This would give both countries a payoff of **1**, **1**. This however, is not a stable equilibrium. This is because if State B knows that State A will choose free trade, as the second mover, it can quickly opt for protectionism. This would result in the "sucker" situation whereby State A would get **-1** and State B would get the much larger payoff of **2**. To avoid this possibility, both players will avoid the institution and use tit-for-tat strategies. This case of a prohibitively expensive DSB therefore shows that when dispute settlement costs are too high for both states, the presence of the WTO is irrelevant. This is because neither state is willing to bear the institutional costs and is willing to simply engage in protectionism.

Proposition 3: In cases where Player A has a significantly lower cost ($\delta_A \leq 1$) than Player B ($\delta_B \geq 2$) to use the DSB, Player A will use protectionism to simultaneously force concessions from Player B and make it worse off than it would be in a world where the DSB does not exist.

Proof:

There are two equilibria, both of which place B at a substantial disadvantage. One equilibrium is "Protect, No DSB; Free Trade, No DSB." In the first scenario where State B responds to State A's

protectionism with free trade, the equilibrium path shows that both states will avoid the DSB. This would lead to an equilibrium of “Protect, No DSB; Free Trade, No DSB.” The payoffs for this strategy are for **2** for State A and for **-1**. State B. This equilibrium is noteworthy because though State A is the culpable party, State B, because of its cumbersome litigation costs, avoids the institution and ends up worse than State A. This is a case of double exploitation by State A in that it firstly trades unfairly with State B. Secondly, since State B finds the DSB process too expensive, it ends up with a payoff that is far worse than State A’s. B will not defect from this equilibrium by filing with the DSB because even though A will acquiesce, yielding B a payoff of $1-\delta$, since $\delta \geq 2$, this payoff is worse than the sucker payoff of -1.

A second equilibrium has state B retaliate with protectionism, only to have that retaliation curtailed by the DSB. On the equilibrium path the moves selected by each state are “Protect, DSB; Protect, Acquiesce.” State B’s retaliation with protectionism ultimately comes to nothing. What we see here is that State A, as a frequent user of the DSB, or a country with greater resources, would simply take State B before the DSB because it can afford to do so. State B is now forced to choose between litigation which gives it a payoff of **- δ** and acquiesce, which yields **-1**. Since $\delta > 1$ for state B, litigation provides worse utility than acquiescing and receiving a payoff of -1. This equilibrium highlights some of the shenanigans that State A can use to exploit the trading relationship that it has with State B. Notice for example that, because State B cannot access the DSB but State A can, A can exploit that institution to deny B the recourse to a protect-protect equilibrium. What this means is that in the case where B has an exorbitant cost to go to the DSB, the inclusion of the institution in the trading dynamics makes B worse off than B would be in the absence of the WTO/DSB.

5. DISCUSSION

This paper has utilized several iterations of a stage game to probe how and when countries choose the use or avoid the DSB whenever a trade dispute arises. This was done under conditions of certainty. The following are some empirical implications that have held constant across the models.

5.1 Countries consider the cost of litigation when determining whether to utilize the DSB mechanism.

All countries that face a trade dispute have several alternatives. A vulnerably interdependent state for example, may choose to ignore the grievance because it is either unable to retaliate, or seeking recourse may make it worse off *ex ante*. A state with a relatively comparable economy, however, has two possible options. One is to engage in tit-for-tat strategies outside the DSB, or to pursue its case formally. Since there is a cost attached to filing, countries will weigh the costs and benefits of litigation versus unilateral retaliation and will act accordingly.

5.2 Sufficiently high costs can lock a country out of enjoying the benefits of the international institution.

The World Trade Organization is a multilateral institution. Most countries accede to it because of the perceived benefit. If however, there are costs to litigation, some of which being prohibitively expensive, then Members will be unable to access the very provisions that they hope to evoke in need of need. If costs are too high as demonstrated by delta, then limited participation when catalyzed by procedural costs would mean that not all countries can freely access the institution. In this way, sufficiently high costs can impede some countries from enjoying the benefits of the institution.

5.3 Lower costs make the threat to utilize the DSB more credible. Hence countries with lower costs will be more likely to successfully resolve disputes they initiated prior to full DSB consideration.

Since countries are more likely to pursue the DSB if the filing costs are moderate, then it also follows that lower costs make the threat to use the institution more credible. Consequently, a culpable state that is threatened with DSB litigation is more likely to settle with the aggrieved party. This is because with no foreseen barrier to the DSB in terms of costs, filing would be sufficient to signal the affected country's intent to get recourse. Since based on the models employed the guilty party is exposed and counter protectionism sanctioned, lower costs make the threat of DSB usage credible and could propel settlement prior to full engagement in the dispute settlement process.

5.4 Lower costs make litigation more available. Hence, countries with high costs will be more likely to settle disputes initiated by trading partners prior to full DSB consideration.

Whenever the DSB is evoked, the defendant has the opportunity to litigate or acquiesce. With lower costs, a respondent, especially a free trader, would be more likely to pursue the case to the fullest extent and get redress. If, however, the litigation costs are too high, then continuing would become more expensive than “protect, protect” and would leave both states worse off. Additionally, litigation would mean that both countries have to pay some cost and not just the party that initiates the process. With these considerations, only lower costs would make litigation attractive and countries knowing that they face great expenses to engage the DSB, could be induced to settle and not go through the whole process.

5.5 With certainty, the presence of the institution prevents full-scale litigation because the guilty party acquiesces.

The iterations of the model show that the presence of the DSB is sufficient to curtail unfair trading practices. This happens regardless of the costs associated with dispute settlement. These instantiations highlight how transparency affects state behaviour. Since it is obvious who is a free trader and who is not, the guilty party acknowledges culpability and does not pursue any litigation. Both states therefore have all the information they need to converge around equilibrium paths of joint free trade or joint protectionism.

5.6 Under certainty, the presence of the institution is sufficient to generate an equilibrium with both sides playing free trade, even though the full process is never utilized.

Prisoner's Dilemma is the game from which this model is generated. Based on its premises, moves are simultaneous and cheating is the dominant strategy. In an extensive form game version, however, moves are sequential and the players have complete and perfect information about each other's strategies and payoffs. Since concurrently playing free trade both a greater payoff than joint protectionism, we see that the presence of the institution serves as a credible monitor of state behavior. Here, while the process is never fully utilized, both states being fully cognizant of the mechanism and their choices still form an equilibrium with free trade and avoid full litigation.

6. CONCLUSION

In the final analysis, this paper sought to test whether the WTO's Dispute Settlement Body is effective in translating membership into participation. This question forms part of the larger debate in international politics about the efficacy of institutions. In an attempt to answer the research question, several instantiations of a formal model are used, each of which has the players moving sequentially, both with certainty about each other's choices and strategies. Each version has Prisoner's Dilemma-like payoffs and the inclusion of delta, the cost for litigation, which varies across the iterations. This affects the different Nash Equilibria that are formed.

The paper examines how information symmetry and sequential moves affect the strategies and outcomes of the game. Delta is varied so that there are instances when it is the same, or either state

alternately has a higher cost. Implicit in these iterations is also the fact that the states deliberate between playing free trade and protectionism, using the DSB or avoiding it, and litigating or acquiescing. These cases are important because they demonstrate that dispute settlement costs as evidenced by delta, affect the decision to use the DSB. This is even more pronounced for cases where the affected country pays more. With a comparable delta however, trade relations are straight forward. Each player is able to see where it is in the game, where the other player is, and also the consequences of each move. This makes it easy to determine if one party is in breach of a WTO provision. The guilty state in turn acknowledges its guilt and accepts punitive countermeasures.

Extremely high costs, however, can make the institution irrelevant, or susceptible to manipulation from countries that pay less. If for example, both players find the cost to go before the DSB too high, they will opt to be protectionists and avoid the institution. If however, one player has an extremely high cost in relation to the other player, the player with the lesser cost is incentivized to play protectionism. The other player, because of the expensive filing costs, will continue to choose free trade and avoid the DSB. If, however, the State with the higher DSB costs retaliates with protectionism, the other player will take it to the DSB, with the only feasible option being to not pursue full scale litigation. This shows that beyond a certain limit, high costs will lock countries out of protecting themselves through protectionism, or make them suffer at the hands of more capable states if the institution is evoked.

There are, however, limitations to these findings. The focus of the paper, for example, is on a model that has complete and perfect information. This scenario does not obtain in the real world because no state ever has full information about what the other is doing. A next step to this project might therefore be to incorporate asymmetric information to see how the equilibria changes. This could be done with some element of signaling and the updating of beliefs based on the probability that countries are engaging in free trade or protectionism. Another dimension to the dispute settlement process that these models do not consider is what may happen to a WTO Member that retaliates after successful litigation. In some cases, a developing country that wins a dispute may be unable to initiate WTO sanctioned countermeasures because of inequitable volumes of trade flow. The outcome would be that resources are spent filing the dispute, but the complainant loses in the end if it is unable to truly singly punish the defector. These possibilities make the decision to litigate a serious matter, which many countries pursue, only if they believe the odds of winning and punishing are very high. Those that are unable to secure these outcomes may choose to respond outside the DSB, or do not act at all.

In general, however, the employment of extensive form games are useful in analyzing the strategies, outcomes and payoffs that are available to countries as they alternate between states of free trade and protectionism. If the variations hold across cases, they therefore can feature in the debate about non participatory membership; that is, if and under what conditions do institutions matter and especially, what makes states use them versus finding recourse elsewhere.

REFERENCES

- Bown, Chad P., and Bernard M. Hoekman. 2005. "WTO Dispute Settlement and the Missing Developing Country Cases: Engaging the Private Sector." *Journal of International Economic Law*. Volume 8, Number 4: 861-890.
- Davis, Christina L., and Sarah Blodgett Bermeo. 2009. "Who Files? Developing Country Participation in GATT/ WTO Adjudication." *The Journal of Politics*. Volume 71, Number 3: 1033 – 1049.
- Gruber, Lloyd. 2000. *Ruling the World: Power Politics and the Rise of Supranational Institutions*. New Jersey: Princeton University Press, 2000.
- Guzman, Andrew T., and Beth A. Simmons. 2005. "Power Plays and Capacity Constraints: the Selection of Defendants in World Trade Organization Disputes." *Journal of Legal Studies*. Volume 34, Number 2: 557-598.
- Hoekman, Bernard, Aaditya Mattoo and Phillip English, eds. 2002. *Development, Trade, and the WTO: A Handbook*. Washington, D.C.: The World Bank. http://www-wds.worldbank.org/servlet/WDSCContentServer/WDSP/IB/2004/08/19/000160016_20040819140633/Rendered/PDF/297990018213149971x.pdf Retrieved February 16, 2016.
- Keohane, Robert O. 1984. *After Hegemony: Cooperation and Discord in the World Political Economy*. New Jersey: Princeton University Press.
- Martin, Lisa L. 1992. "Interests, Power, and Multilateralism." *International Organization*. Volume 46, Number 4. (Autumn) pp. pp. 765 – 792. <http://www.jstor.org/stable/2706874>
- Matsushita, Mitsuo and Thomas Schoenbaum. 2003. *The World Trade Organization: Law, Practice and Policy*. New York: Oxford University Press.
- Morrow, James D. 1994. *Game Theory for Political Scientists*. New Jersey: Princeton University Press.
- Pollins, Brian M. 1982. "Modeling International Trade Flows: A Survey and Comparison of Simulation Approaches." *International Political Science Review / Revue internationale de science politique*, Volume 3, Number 4, Mathematical Approaches to International Relations, pp. 504 – 533.
- Snidal, Duncan. 1985. "Coordination Versus Prisoners' Dilemma: Implications for International Cooperation and Regimes." *The American Political Science Review*, Volume 79, Number 4 (December), pp. 923 – 942.
- "What is the WTO?" https://www.wto.org/english/thewto_e/whatis_e/whatis_e.htm Accessed January 8, 2016.

A SYSTEM DYNAMICS APPROACH TO PREDICT THE TREND OF SUPERHERO MOVIES OVER TIME

Daniel Perez Ibanez

Department of Modeling, Simulation & Visualization Engineering
Old Dominion University
Norfolk, VA
dpere013@odu.edu

ABSTRACT

This research studies the trend of the movie industry over the years. Specifically, the goal of the project is to predict the number of movies from a specific genre based on the results of similar movies during the previous years. Due to their recent popularity, data from superhero movies was collected to perform the modeling and simulation. However, the system was designed so that it could be applied to different types of movies.

The development of the system consisted of collecting necessary data from different movie databases and then applying it to a System Dynamics model developed with Vensim. A Stock and Flow diagram was generated and simulated to perform the prediction during the period of time from 1977 to 2020. The results obtained demonstrate that the system achieves prediction with an average error of 14.19%, which makes the system a useful tool for individuals or companies interested on the trend of the movie industry.

Keywords: System Dynamics, Vensim, Stock and Flow Diagram, Movie Industry

1 INTRODUCTION

The movie industry is in constant change. The trends in blockbusters change over the years, and the popularity of a certain type of movie can vary significantly with time. For instance, Western movies were very popular for many years, but in the last decade, the number of westerns has been constantly decreasing (Agresta 2013). During the last years, superhero movies have been one of the most popular genres of blockbusters. Both Marvel and DC have increased their movie catalog and they already announced most of the movies that they plan to release in the upcoming years (Wheeler 2014).

This project covers the modeling and simulation of a system that is able to predict the trend of a certain type of movie over the years. The research will be focused on predicting the number of movies that will be released at a certain time based on the results obtained by similar movies during the previous years. Ideally, this model could be applied to different types of movies, but it will be focused on superhero movies due to their recent popularity. This model could be useful to individuals or companies who are interested on the future popularity of a type (or different types) of movies. For instance, merchandise companies could use the prediction to determine in what area they want to focus during the next years to improve their sales.

The modeling and the simulation of this system were realized using the principles of System Dynamics: an approach to studying complex systems that relies on the interaction of one system variable with another (Sokolowski and Banks 2009). The software utilized to model and simulate the system was Vensim (Ventana Systems 2015). The remainder of the paper is organized as follows: related work is presented on Section 2. Then, the development of the model is covered in Section 3. Section 4 analyzes the results of the simulation and finally, conclusions are drawn in Section 5.

2 RELATED WORK

There have been several attempts to predict the behavior of the movie industry. However, the majority of them focused on predicting the success of movies in terms of box-office revenue and viewer ratings.

A common practice in the literature is to use data extracted from social media in order to obtain the predictions. Asur and Huberman (2010) built a linear regression model that used twitter data to forecast box-office revenue for movies in advance of their release. Krauss *et al.* (2008) used social network analysis and web data mining to run a model for forecasting movie success. By analyzing the data from online forum discussions on the Internet Movie Database (IMDb), they successfully predicted 9 Oscar nominations and they found an existing correlation between the discussions and the box office success of 20 top grossing movies of 2006.

Doshi *et al.* (2010) also used data from social media, but in this case, they focused on the prediction of movie prices. Specifically, they tried to predict prices on the Hollywood Stock Exchange (HSX) and the ratio of gross income to production budget. In order to accomplish that, they used three types of metrics: (1) movie ratings from IMDb and Rotten Tomatoes, as well as Box Office performance data from Box Office Mojo, (2) general discussions from the web and from bloggers and (3) posts from IMDb forums. Using a multi-linear and non-linear regression model, they achieved an accuracy of 80% on predicting whether a movie will be a blockbuster.

Zhang and Skiena (2009) followed a slightly different approach on the prediction of movie trends. They used a system for large-scale news analysis called *Lydia* in conjunction with two different models (regression and *k*-nearest neighbor) to predict movie gross. They achieved an overall accuracy of 92.1%, which was increased to 96.81% when the news data was combined with quantitative data from IMDb. Their experiments prove that both qualitative and quantitative data can be used in the prediction of movie gross.

A different practice on the prediction of movie success consists of using quantitative data that is publicly available. When following this approach, the variables used have a very important role on the performance of the prediction model. On an attempt to predict movie gross, Simonoff and Sparrow (2000) used the following predictor variables: (1) the genre of the film, (2) the Motion Picture Association of America (MPAA), (3) the origin country of the movie, (4) the popularity of the actors and actresses appearing in the movie, (5) the production budget of the film, (6) whether or not the movie was a sequel to an earlier movie, (7) the season when the movie was released, (8) the number of screenings of the movie during its first weekend of general release, (9) the gross revenue during its first weekend of general release, (10) the rating of the movie by the critic Roger Ebert, and (11) the number of Academy Awards nominations and wins for the movie. Using all these variables, they created a model that was capable of predict grosses of newly-released movies. Specifically, their model performed well in 21 of their 24 prediction intervals. A similar study was performed by Chang and Ki (2005), who determined that sequel, actors, budget, genre, MPAA rating, release periods and number of first-week screenings were significantly related to total box office performance.

System Dynamics is not very common on this area. This modeling technique was used by Ramsden (2009) in order to predict box office revenue. His model used the purchase of a movie as the stock variable and the

potential and served customers as the flows. Different variables such as *word-of-mouth* and *advertisement recruitment* were used, and a mean absolute error of 7.5% was obtained.

It can be appreciated that there have been several attempts to predict movie gross and ratings. However, to the best of the author's knowledge, a model that predicts the number of movies that will be released during a certain time has not been studied before. Besides, the use of System Dynamics in this area can be further explored. Based on these assumptions, this paper offers an innovative method to predict the trend of movies over time.

3 DEVELOPMENT

The development of the system consisted on two main phases: the collection of the data needed for the simulation and the design of the model. During the first phase, data from superhero movies was collected in a table and plotted in different graphs so that it could be analyzed. The second phase used this analysis of the data to design a System Dynamics model in Vensim. The following subsections cover the details of the development process.

3.1 Data Collection

The first stage of the project consisted of collecting information about superhero movies during the past years. Based on the literature reviewed, it was determined that the most important variables that affect the success of a movie are the money that it makes and the viewer ratings that it receives. With this assumption, a table with all the superhero movies from 1966 to 2016 was constructed (Appendix A). Each entry of the table contains the following information: release year of the movie, name of the movie, worldwide gross (how much money it made from ticket sales), budget (money spent on the movie) and ratings from two independent websites.

The information was extracted from different public movie databases. The list of movies was obtained from IMDb. Box Office Mojo was used to get the worldwide gross of each movie, while the budget was collected from The Numbers website. The ratings were collected from two different sites: IMDb and Rotten Tomatoes. IMDb includes ratings from anyone with an account on their site, while Rotten Tomatoes' ratings are based on the published opinions of movie critics.

Once all the data was obtained, a new table was computed with the average results of superhero movies by year (Appendix B). This table contains the number of superhero movies for each year from 1966 to 2016, as well their corresponding average gross, budget and ratings.

The initial idea behind this project was to predict how many movies would come out in a specific year based on the results of similar movies during the previous year. However, looking at the table from Appendix B, it can be seen that there are some years in which no superhero movies were released. Since there is no data during those years, a prediction for the next year is not possible in those cases. To solve this problem, the years were grouped in seasons. In order to determine how many years should be in a season, it was asserted that a low number should be chosen to make the number of outliers as large as possible. Sets of 2 years were not possible because there are points on time in which no superhero movies are released during 2 years. To guarantee that each period has at least two years of released movies, the number of years in each season was set to 4. Besides, since there is just one superhero movie from 1966 to 1977, it was decided to discard that period of time. With these premises, a final table (Table 1) with the data necessary to perform the simulation was obtained.

Table 1: Average Results of Superhero Movies by Season (set of 4 years)

Season	Number of Superhero Movies	Average Gross	Average Budget	Av. IMDb Rating	Av. Rotten Tomatoes Rating
1977-1980	3	\$181,139,120	\$54,500,000	6.90	82%
1981-1984	3	\$39,102,155	\$28,368,750	5.08	32%
1985-1988	2	\$26,821,897	\$23,500,000	4.10	14%
1989-1992	10	\$131,041,320	\$34,716,667	6.00	51%
1993-1996	13	\$76,590,164	\$33,333,333	5.18	35%
1997-2000	11	\$107,120,547	\$54,062,500	5.76	50%
2001-2004	11	\$323,665,722	\$95,183,333	6.35	57%
2005-2008	25	\$289,632,037	\$117,143,750	5.95	45%
2009-2012	19	\$316,803,980	\$110,041,667	6.75	56%
2013-2016	21	\$671,858,408	\$171,265,000	6.94	59%

Then, 3 variables were plotted (Figure 1) to observe how the results of previous years affected the number of movies in the next season. These variables are: the number of movies, the average benefits (difference of gross and budget), and the average ratings (on a scale from 0 to 10).

Looking at Figure 1, it can be assessed that both the average benefits and the ratings influence the number of movies that will be released in the following season. It was discovered that generally, if the difference of the results from one season to another was positive, then the number of movies in the next season would increase, and vice versa. For instance, both the ratings and the benefits of the 1989-1992 season were better than the ones in the previous season (1985-1988) and consequentially, the number of movies during the next season (1993-1996) increased.

3.2 Design of the Model

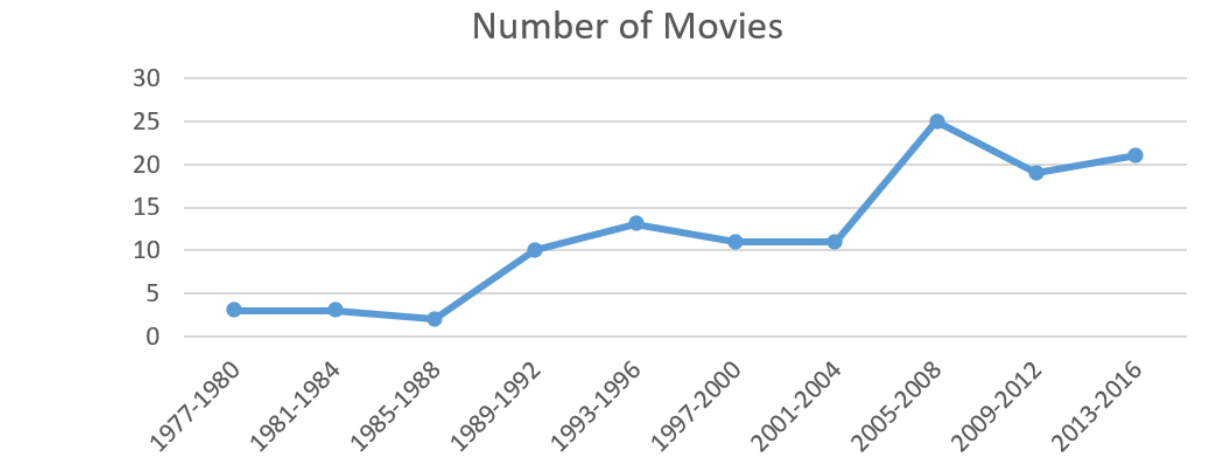
Once all the data was collected and analyzed, a stock-flow diagram was created with Vensim. The model is shown in Figure 2.

The variables *IMDb Ratings*, *Rotten Tomatoes Ratings*, *Average Budget* and *Average Gross* are lookup tables that contain the data displayed in Table 1. An auxiliary variable called *Season* was created so that it could be used as the input in the lookup table functions. This variable is increased by one in each time step, simulating the passing of the time.

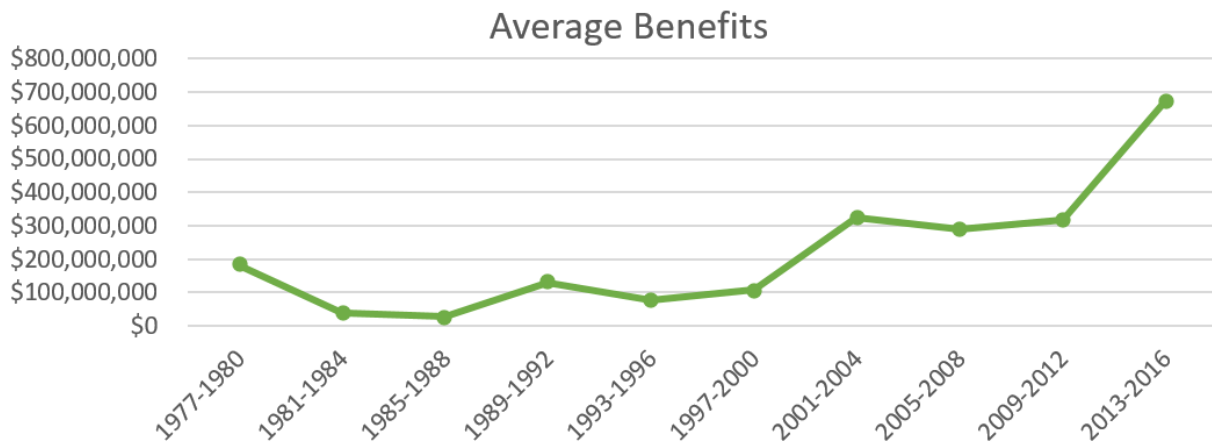
The lookup table variables are summarized into two variables: *Benefits (B)* and *Overall Ratings (R)*. These variables are the ones shown in Figure 1b and 1c. The benefits are the difference between the gross and the budget, and the overall ratings are the average of the scores from IMDb and Rotten Tomatoes. Since the Rotten Tomatoes ratings are percentages, they had to be multiplied by ten so that they matched the IMDb scores. Equations 1 and 2 show the computation of these variables.

$$B = AvGross - AvBudget \quad (1)$$

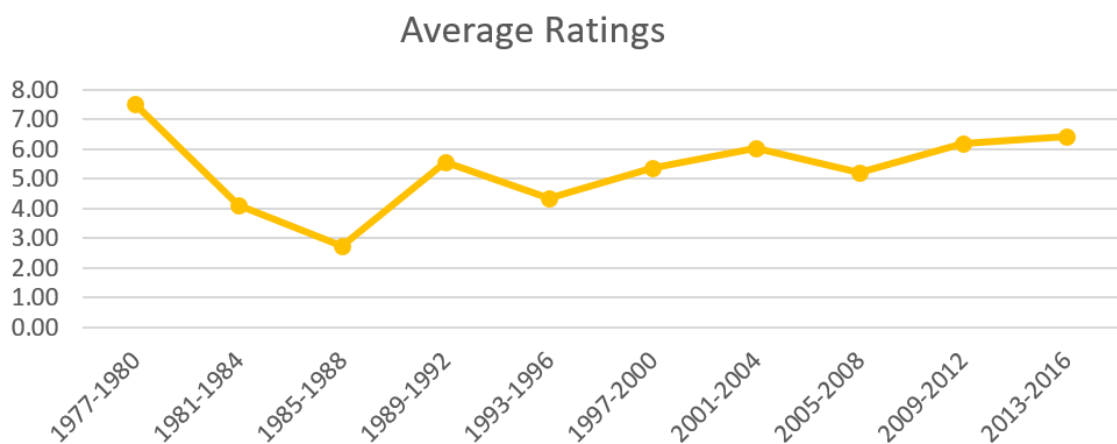
$$R = \frac{R_{IMDb} + R_{RT} * 10}{2} \quad (2)$$



(a) Number of superhero movies over time



(b) Average benefits of superhero movies over time



(c) Average ratings of superhero movies over time

Figure 1: Charts showing the collected data after being processed

The ratings and benefits are the two main variables that will determine the result of the simulation. After analyzing the data in Figure 1, it was asserted that what really influences the number of movies on a season is the change of the variables with respect to the previous season. If this change is positive, the number of movies in the next season is likely to be increased, otherwise, the number of movies tends to decrease. With this in mind, the fraction change of both ratings and benefits was computed. Before calculating the change, the ratings and benefits from the previous years had to be entered. To do that, two variables called *Previous Ratings* and *Previous Benefits* were created. These variables store the values of the ratings and benefits (respectively) from the previous period using a delay fixed function. The delay fixed function takes a variable and delays it for an amount of time. In this case, the amount of time was 1 unit, and the variables to be delayed were the ratings and the benefits. Equations 3 and 4 show the calculations of the changes for a certain time t .

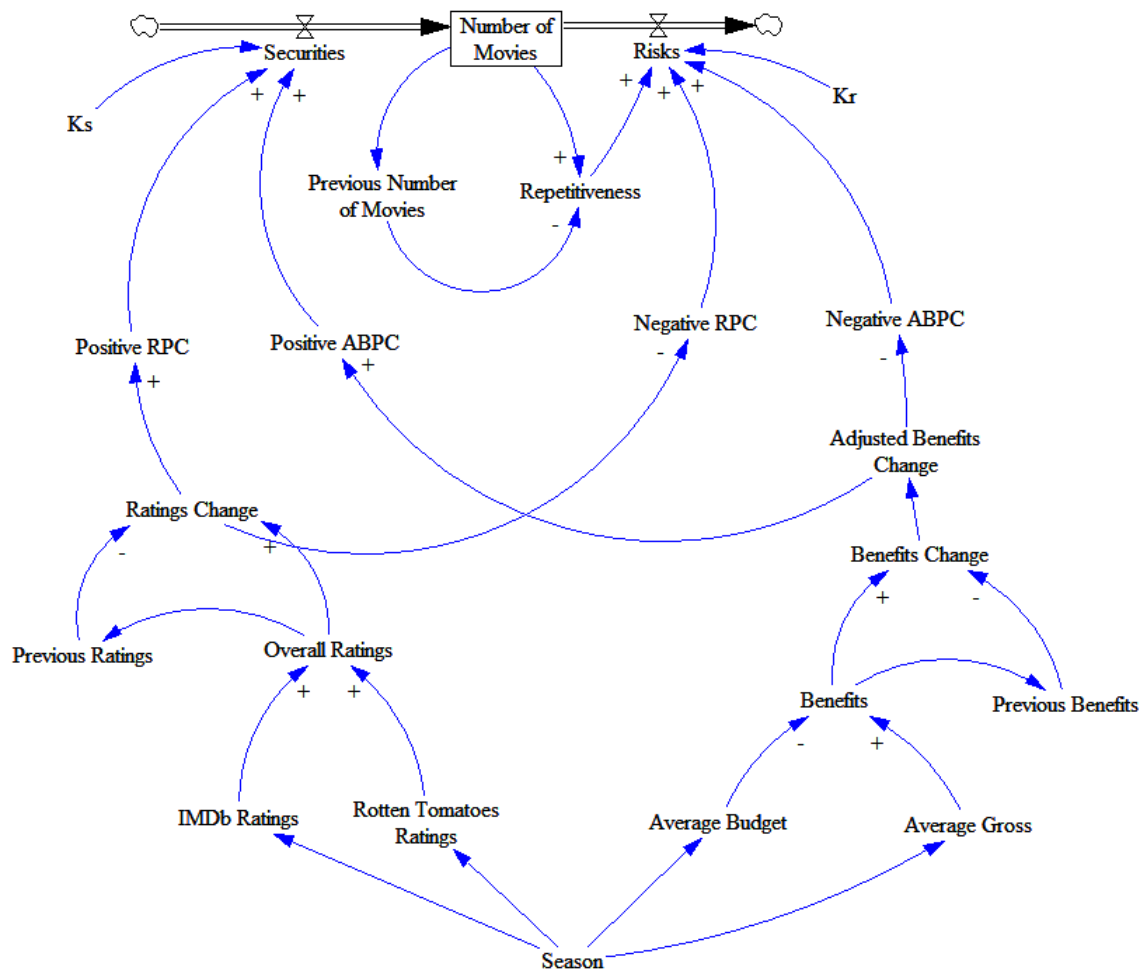


Figure 2: Stock and flow diagram of the model

$$RC_t = \frac{R_t - R_{t-1}}{R_{t-1}} \quad (3)$$

$$BC_t = \frac{B_t - B_{t-1}}{B_{t-1}} \quad (4)$$

In the case of the *Benefits Change* (BC), there was one issue that had to be addressed. There is an increase of 2,825% between the benefits in 1985-1988 and the ones in 1989-1992. This substantial increase happens because the benefits obtained in the first season were really low compared to the ones in the next season. Looking at the real data, it can be seen that this increase affects the number of movies during 1993-1996 as expected. However, this boost in the simulation resulted in an immense number of movies during the 1993-1996 period, which made the model impractical after that season. To make the simulation more realistic, the BC had to be clamped. After performing different simulations, it was determined that the best threshold to be applied was 500% (or 5.0 in fraction notation). In order to do that, the variable *Adjusted Benefits Change* (ABC) was created. The value of this variable is calculated with an IF-ELSE statement as shown in Equation 5.

$$ABC = \begin{cases} BC, & \text{if } BC < 5 \\ 5, & \text{otherwise} \end{cases} \quad (5)$$

Once the *Ratings Change* (RC) and *Adjusted Benefits Change* have been defined, each of them is divided into two new variables. These variables filter the RC and ABC so that only the positive or negative values are selected respectively. To obtain this behavior in the model, a new set of IF-ELSE statements was defined. In the case of the positive variables, the output is the RC or ABC if they are greater than zero, otherwise, they are zero. The negative case is similar, but returning the absolute values of the input variables. The equations for the positive and negative RC are included below. The corresponding equations for the benefits are similar, but using the ABC instead of the RC .

$$PRC = \begin{cases} RC, & \text{if } RC \geq 0 \\ 0, & \text{otherwise} \end{cases} \quad NRC = \begin{cases} |RC|, & \text{if } RC < 0 \\ 0, & \text{otherwise} \end{cases} \quad (6)$$

In an effort to make the system realistic to the dynamic change of people tastes', a new variable called *Repetitiveness* (Rp) was introduced as an hypothesis. This variable takes into account the fact that viewers tend to lose the interest in a movie genre when it gets too repetitive. To simulate this behavior, the repetitiveness was defined as the change between the number of movies of a season and the number of movies in the previous season. Equation 7 shows the computation of Rp for a certain time-point t , where N_t is the number of movies at time t .

$$Rp_t = \frac{N_t - N_{t-1}}{N_{t-1}} \quad (7)$$

Once those variables were declared, the stock and flows of the model were defined. The flows of the model are the *Securities* (S) and *Risks* (Rs) of making superhero movies, while the stock is the *Number of Movies* (N) that will be released during that season. The stock N is defined as the integral of the flows S and Rs during time, as shown in Equation 8.

$$N = \int (S - Rs)dt \quad (8)$$

The *Securities* variable is influenced by the Positive Ratings Change (*PRC*) and the Positive Average Benefits Change (*PABC*), while the *Risks* variable depends on the negative variables and the repetitiveness. The computation of this variables is shown in Equations 9 and 10. The weights for each variable were defined after several trial-error experiments.

$$S = Ks * (PABC * 0.75 + PRC * 0.25) \quad (9)$$

$$Rs = Kr * (NABC * 0.5 + NRC * 0.25 + Rp * 0.25) \quad (10)$$

It can be appreciated that the risks and securities depend on the variables *Ks* and *Kr* respectively. These variables are constants used to adjust the result of the simulation. The first approach to adjust these variables relied on the fact that the securities to make a movie would affect more the outcome of the simulation than the risks. Following this hypothesis, the securities' constant (*Ks*) was set as twice the value of *Kr*. The first simulations generated a result with a trend slightly similar to reality: the number of movies increased and decreased in a similar way, but the actual numbers differed considerably. Thus, although the hypothesis was correct, the constants needed to be adjusted to get results closer to the real numbers. After performing different simulations, it was determined that the optimum values for the constants were:

$$Ks = 2.9 \quad Kr = 1.5$$

4 RESULTS

Once the model was designed, the simulation was performed for 11 units of time, corresponding to the 11 sets of 4 years between 1977 and 2020. Table 2 and Figure 3 show the results obtained for the *Number of Movies* variable and compare them to the ground truth from Figure 1a. The values obtained were rounded to the nearest integer. It can be appreciated that the simulation coincides with the real data at most of the points. The only transition in which the simulation does not match the real data is the one between 1985-1988 and 1989-1992. Looking at Figure 1b and 1c, it can be seen that the changes during the previous seasons were negative, but the number of movies significantly increased. Thus, although the simulation does not match reality at that point, it is consistent with the assumptions made in the model.

The period 2017-2020 was added in Table 2 and Figure 3. Since the model predicts the number of movies in a specific season based on the results of the previous seasons, the model is able to predict the results during the next years. According to the simulation, the number of superhero movies released in the next four years will increase with respect to the last season.

The percentage error in the predictions is shown in the last column of Table 2. For each entry in the table, the error was defined as:

$$Error = \frac{|RealData - Prediction|}{RealData} * 100\% \quad (11)$$

It can be seen that the greatest error was 80% in the season 1989-1992 (which was previously discussed). However, in half of the cases, the error is 0%, and the average is 14.19%. These low values prove that the model has the ability to predict the trend of superhero movies over the years.

Table 2: Results obtained after running the simulation

Season	Real Data	Prediction	Error
1977-1980	3	3	0.00%
1981-1984	3	3	0.00%
1985-1988	2	2	0.00%
1989-1992	10	2	80.00%
1993-1996	13	13	0.00%
1997-2000	11	10	9.09%
2001-2004	11	11	0.00%
2005-2008	25	18	28.00%
2009-2012	19	17	10.53%
2013-2016	21	18	14.29%
2017-2020	–	21	–

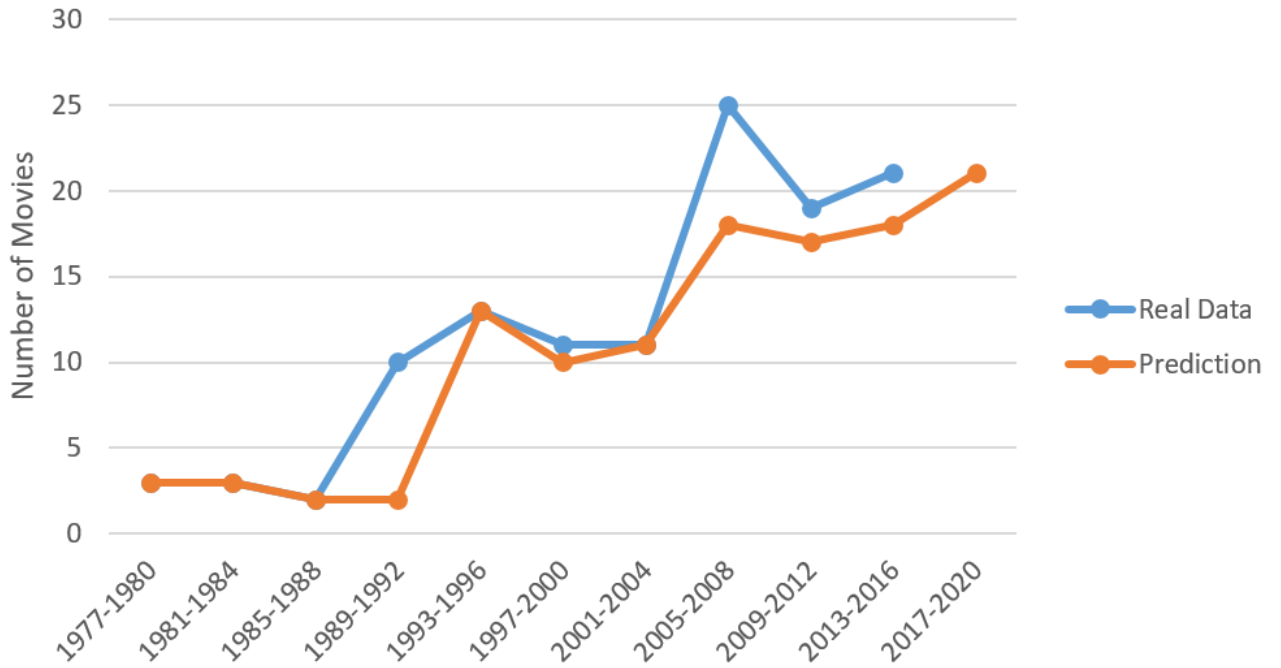


Figure 3: Results plotted and compared to the real data

5 CONCLUSIONS

In conclusion, this model successfully predicts the number of superhero movies based on results obtained during the previous years by the same genre. The system is able to determine the number of superhero movies during the last four decades by analyzing the previous benefits and ratings obtained by similar movies, as well as to predict the trend of the genre during the next four years. The results obtained demonstrate that the simulation is close to reality and thus, the prediction obtained for the next years can be considered reliable, which makes this model a potentially useful tool for people or companies interested on predicting the trends of the movie industry. Although the model is focused on the trend of superhero movies,

this simulation could be applied to different genres by updating the lookup tables with the corresponding data and adjusting the values of the securities' and risks' weights and constants.

With this model, I have attempted to demonstrate an ability to predict trends of movies by looking at the results obtained in the previous years by similar films. The results obtained show that System Dynamics is a good approach to generate these predictions. However, in order to prove its full potential, the system needs to be tested with more genres. This simulation demonstrate the feasibility of the model and marks a promising start on the prediction of the number of movies released during a certain period of time. For future work, I plan to feed the simulation with more data to improve its reliability and efficiency.

ACKNOWLEDGMENTS

The author would like to thank Dr. John Sokolowski and Dr. Catherine M. Banks for their help and assistance during the development of this project, as this paper is a result from their class MSIM 872 'Modeling Global Events'.

A TABLE OF SUPERHERO MOVIES FROM 1966 TO 2016

Year	Movie	Worldwide Gross	Budget	IMDb Rating	RT Rating
1966	Batman	\$1,700,000		6.8	80%
1978	Superman	\$300,218,018	\$55,000,000	7.3	93%
1980	Hero at Large	\$15,934,737		6.2	52%
1980	Superman II	\$108,185,706	\$54,000,000	6.8	89%
1983	Superman III	\$70,656,090	\$39,000,000	4.9	26%
1984	The Toxic Avenger	\$800,000	\$475,000	6.2	68%
1984	Supergirl	\$14,296,438	\$35,000,000	4.3	7%
1986	Howard the Duck	\$37,962,774	\$30,000,000	4.6	15%
1987	Superman IV: The Quest for Peace	\$15,681,020	\$17,000,000	3.6	12%
1989	The Toxic Avenger Part II	\$792,966	\$2,300,000	5.1	0%
1989	The Return of Swamp Thing	\$192,816		4.4	33%
1989	Batman	\$411,348,924	\$35,000,000	7.6	72%
1989	The Toxic Avenger Part III: The Last Temptation of Toxie	\$363,561	\$2,300,000	4.3	N/A
1990	Teenage Mutant Ninja Turtles	\$201,965,915	\$13,500,000	6.7	40%
1990	Darkman	\$48,878,502	\$16,000,000	6.4	83%
1990	Captain America	\$675,437	\$10,000,000	3.3	8%
1991	Teenage Mutant Ninja Turtles II: The Secret of the Ooze	\$78,656,813	\$25,000,000	6	32%
1991	The Rocketeer	\$62,000,000	\$40,000,000	6.4	62%
1992	Batman Returns	\$266,822,354	\$80,000,000	7	80%
1993	Teenage Mutant Ninja Turtles III	\$42,273,609	\$21,000,000	4.8	21%
1993	The Meteor Man	\$8,023,147	\$30,000,000	5	29%
1994	The Crow	\$144,693,129	\$15,000,000	7.6	82%
1994	The Shadow	\$48,063,435	\$40,000,000	6	35%
1994	The Mask	\$351,583,407	\$18,000,000	6.9	77%
1994	Blankman	\$7,941,977		4.8	13%
1995	Tank Girl	\$4,064,495	\$25,000,000	5.2	38%
1995	Batman Forever	\$336,529,844	\$100,000,000	5.4	40%
1995	Judge Dredd	\$113,493,481	\$90,000,000	5.5	18%
1995	Mighty Morphin Power Rangers: The Movie	\$66,433,194	\$15,000,000	5.1	50%
1996	Barb Wire	\$3,794,000	\$23,000,000	3.2	28%
1996	The Phantom	\$17,323,326	\$42,000,000	4.9	41%
1996	The Crow: City of Angels	\$17,917,287	\$13,000,000	4.5	12%
1997	Turbo: A Power Rangers Movie	\$8,363,899	\$8,000,000	3.4	17%
1997	Batman & Robin	\$238,207,122	\$125,000,000	3.7	11%
1997	Spawn	\$87,840,042	\$40,000,000	5.2	19%
1997	Steel	\$1,710,972	\$16,000,000	2.8	12%
1998	Star Kid	\$7,029,025	\$12,000,000	5.4	38%
1998	The Mask of Zorro	\$250,288,523	\$95,000,000	6.7	83%
1998	Blade	\$131,183,530	\$45,000,000	7.1	54%

1999	Mystery Men	\$33,461,011	\$68,000,000	6	60%
2000	X-Men	\$296,339,527	\$75,000,000	7.4	81%
2000	The Specials	\$13,276	\$1,000,000	6	47%
2000	Unbreakable	\$248,118,121	\$75,000,000	7.2	68%
2002	Blade II	\$155,010,032	\$55,000,000	6.7	57%
2002	Spider-Man	\$821,708,551	\$139,000,000	7.3	89%
2003	Daredevil	\$179,179,718	\$75,000,000	5.3	44%
2003	X2: X-Men United	\$407,711,549	\$110,000,000	7.5	86%
2003	Hulk	\$245,360,480	\$120,000,000	5.7	61%
2003	The League of Extraordinary Gentlemen	\$179,265,204	\$78,000,000	5.8	17%
2004	Hellboy	\$99,318,987	\$66,000,000	6.8	81%
2004	The Punisher	\$54,700,105	\$33,000,000	6.5	29%
2004	Spider-Man 2	\$783,766,341	\$200,000,000	7.3	93%
2004	Catwoman	\$82,102,379	\$100,000,000	3.3	9%
2004	Blade: Trinity	\$128,905,366	\$65,000,000	5.9	25%
2005	Elektra	\$56,681,566	\$43,000,000	4.8	10%
2005	Constantine	\$230,884,728	\$100,000,000	6.9	46%
2005	Son of the Mask	\$57,552,641	\$84,000,000	2.2	6%
2005	The Adventures of Sharkboy and Lavagirl in 3-D	\$69,425,967	\$350,000,000	3.5	20%
2005	Batman Begins	\$374,218,673	\$150,000,000	8.3	84%
2005	Fantastic Four	\$330,579,719	\$100,000,000	5.7	27%
2005	Sky High	\$86,369,815	\$35,000,000	6.2	73%
2005	The Legend of Zorro	\$142,400,065	\$80,000,000	5.9	26%
2006	V for Vendetta	\$132,511,035	\$54,000,000	8.2	73%
2006	X-Men: The Last Stand	\$459,359,555	\$168,000,000	6.8	58%
2006	Superman Returns	\$391,081,192	\$204,000,000	6.1	76%
2006	My Super Ex-Girlfriend	\$60,984,606	\$65,000,000	5.1	40%
2006	Zoom	\$12,506,188	\$35,000,000	4.2	3%
2007	Ghost Rider	\$228,738,393	\$120,000,000	5.2	26%
2007	Fantastic Four: Rise of the Silver Surfer	\$289,047,763	\$130,000,000	5.6	37%
2007	Spider-Man 3	\$890,871,626	\$250,000,000	6.2	63%
2007	Underdog	\$65,270,477	\$25,000,000	4.8	14%
2008	Jumper	\$222,231,186	\$85,000,000	6.1	16%
2008	Iron Man	\$585,174,222	\$140,000,000	7.9	94%
2008	The Incredible Hulk	\$263,427,551	\$150,000,000	6.8	67%
2008	Hancock	\$624,386,746	\$150,000,000	6.4	41%
2008	Hellboy II: The Golden Army	\$160,388,063	\$85,000,000	7	85%
2008	The Dark Knight	\$1,004,558,444	\$180,000,000	9	94%
2008	Punisher: War Zone	\$10,100,036	\$35,000,000	6	27%
2008	The Spirit	\$39,031,337	\$60,000,000	4.8	14%
2009	Push	\$48,808,215	\$38,000,000	6.1	23%
2009	Watchmen	\$185,258,983	\$120,000,000	7.6	65%
2009	X-Men Origins: Wolverine	\$373,062,864	\$150,000,000	6.7	38%

2009	Defendor	\$44,462	\$3,000,000	6.8	68%
2010	Kick-Ass	\$96,188,903	\$30,000,000	7.7	76%
2010	Iron Man 2	\$623,933,331	\$200,000,000	7	72%
2010	Jonah Hex	\$10,903,312	\$47,000,000	4.7	12%
2010	Super	\$327,716	\$2,000,000	6.8	49%
2011	The Green Hornet	\$227,817,248	\$120,000,000	5.8	43%
2011	Thor	\$449,326,618	\$150,000,000	7.1	66%
2011	X-Men: First Class	\$353,624,124	\$160,000,000	7.8	65%
2011	Green Lantern	\$219,851,172	\$200,000,000	5.6	26%
2011	Captain America: The First Avenger	\$370,569,774	\$140,000,000	6.9	80%
2012	Chronicle	\$126,636,097	\$12,000,000	7.1	69%
2012	Ghost Rider: Spirit of Vengeance	\$132,563,930	\$75,000,000	4.3	17%
2012	The Avengers	\$1,518,594,910	\$220,000,000	8.1	92%
2012	The Amazing Spider-Man	\$752,216,557	\$230,000,000	7	72%
2012	The Dark Knight Rises	\$1,084,439,099	\$250,000,000	8.5	87%
2012	Dredd	\$35,626,525	\$45,000,000	7.1	78%
2013	Iron Man 3	\$1,215,439,994	\$200,000,000	7.2	79%
2013	Man of Steel	\$668,045,518	\$225,000,000	7.1	55%
2013	The Wolverine	\$414,828,246	\$120,000,000	6.7	69%
2013	Kick-Ass 2	\$60,795,985	\$28,000,000	6.6	31%
2013	Thor: The Dark World	\$644,783,140	\$170,000,000	7.1	66%
2014	Captain America: The Winter Soldier	\$714,766,572	\$170,000,000	7.8	89%
2014	The Amazing Spider-Man 2	\$708,982,323	\$250,000,000	6.7	52%
2014	X-Men: Days of Future Past	\$748,121,534	\$200,000,000	8	91%
2014	Guardians of the Galaxy	\$774,176,600	\$232,300,000	8.1	91%
2014	Teenage Mutant Ninja Turtles	\$485,004,754	\$125,000,000	5.9	22%
2015	Avengers: Age of Ultron	\$1,405,413,868	\$250,000,000	7.5	75%
2015	Ant-Man	\$519,445,163	\$130,000,000	7.4	81%
2015	Fantastic Four	\$167,977,596	\$120,000,000	4.3	9%
2016	Deadpool	\$782,370,866	\$58,000,000	8.1	84%
2016	Batman v Superman: Dawn of Justice	\$872,662,631	\$250,000,000	6.8	27%
2016	Captain America: Civil War	\$1,151,314,869	\$250,000,000	8	90%
2016	X-Men: Apocalypse	\$534,348,852	\$178,000,000	7.2	48%
2016	Teenage Mutant Ninja Turtles: Out of the Shadows	\$235,180,501	\$135,000,000	6.1	38%
2016	Suicide Squad	\$641,117,209	\$175,000,000	6.6	26%

B B. TABLE OF AVERAGE RESULTS OF SUPERHERO MOVIES BY YEAR

Year	Number of Movies	Average Gross	Average Budget	Av. IMDb Rating	Av. RT Rating
1966	1	\$1,700,000	N/A	6.80	80%
1967	0				
1968	0				
1969	0				

1970	0				
1971	0				
1972	0				
1973	0				
1974	0				
1975	0				
1976	0				
1977	0				
1978	1	\$300,218,018	\$55,000,000	7.30	93%
1979	0				
1980	2	\$62,060,222	\$54,000,000	6.50	71%
1981	0				
1982	0				
1983	1	\$70,656,090	\$39,000,000	4.90	26%
1984	2	\$7,548,219	\$17,737,500	5.25	38%
1985	0				
1986	1	\$37,962,774	\$30,000,000	4.60	15%
1987	1	\$15,681,020	\$17,000,000	3.60	12%
1988	0				
1989	4	\$103,174,567	\$13,200,000	5.35	35%
1990	3	\$83,839,951	\$13,166,667	5.47	44%
1991	2	\$70,328,407	\$32,500,000	6.20	47%
1992	1	\$266,822,354	\$80,000,000	7.00	80%
1993	2	\$25,148,378	\$25,500,000	4.90	25%
1994	4	\$138,070,487	\$24,333,333	6.33	52%
1995	4	\$130,130,254	\$57,500,000	5.30	37%
1996	3	\$13,011,538	\$26,000,000	4.20	27%
1997	4	\$84,030,509	\$47,250,000	3.78	15%
1998	3	\$129,500,359	\$50,666,667	6.40	58%
1999	1	\$33,461,011	\$68,000,000	6.00	60%
2000	3	\$181,490,308	\$50,333,333	6.87	65%
2001	0				
2002	2	\$488,359,292	\$97,000,000	7.00	73%
2003	4	\$252,879,238	\$95,750,000	6.08	52%
2004	5	\$229,758,636	\$92,800,000	5.96	47%
2005	8	\$168,514,147	\$117,750,000	5.44	37%
2006	5	\$211,288,515	\$105,200,000	6.08	50%
2007	4	\$415,063,289	\$135,000,000	5.53	38%
2008	8	\$363,662,198	\$110,625,000	6.75	55%
2009	4	\$151,793,631	\$77,750,000	6.80	49%
2010	4	\$182,838,316	\$69,750,000	6.55	52%
2011	5	\$324,237,787	\$154,000,000	6.64	56%
2012	6	\$608,346,186	\$138,666,667	7.02	69%
2013	5	\$600,778,577	\$148,600,000	6.94	60%

2014	5	\$686,210,357	\$195,460,000	7.30	69%
2015	3	\$697,612,209	\$166,666,667	6.40	55%
2016	8	\$702,832,488	\$174,333,333	7.13	52%

REFERENCES

- Agresta, M. 2013. “How the Western Was Lost (and Why It Matters)”. <http://www.theatlantic.com/entertainment/archive/2013/07/how-the-western-was-lost-and-why-it-matters/278057/>. Accessed March 14, 2017.
- Asur, S., and B. A. Huberman. 2010. “Predicting the future with social media”. In *Web Intelligence and Intelligent Agent Technology (WI-IAT), 2010 IEEE/WIC/ACM International Conference on*, Volume 1, pp. 492–499. IEEE.
- Chang, B.-H., and E.-J. Ki. 2005. “Devising a practical model for predicting theatrical movie success: Focusing on the experience good property”. *Journal of Media Economics* vol. 18 (4), pp. 247–269.
- Doshi, L., J. Krauss, S. Nann, and P. Gloor. 2010. “Predicting movie prices through dynamic social network analysis”. *Procedia-Social and Behavioral Sciences* vol. 2 (4), pp. 6423–6433.
- Krauss, J., S. Nann, D. Simon, P. A. Gloor, and K. Fischbach. 2008. “Predicting Movie Success and Academy Awards through Sentiment and Social Network Analysis.”. In *ECIS*, pp. 2026–2037.
- Ramsden, E. 2009. “Box office revenue forecasting: a system dynamics approach”. *White Paper, Oregon*.
- Simonoff, J. S., and I. R. Sparrow. 2000. “Predicting movie grosses: Winners and losers, blockbusters and sleepers”. *Chance* vol. 13 (3), pp. 15–24.
- Sokolowski, J. A., and C. M. Banks. 2009. *Modeling and simulation for analyzing global events*. John Wiley & Sons.
- Ventana Systems, Inc. 2015. “Vensim”. <http://vensim.com/>. Accessed March 14, 2017.
- Wheeler, A. 2014. “Infographic: New Superhero Movies Between Now And 2020”. <http://comicsalliance.com/your-supermovie-timeline-updated-with-marvel-studios-phase-three-releases-infographic/>. Accessed March 14, 2017.
- Zhang, W., and S. Skiena. 2009. “Improving movie gross prediction through news analysis”. In *Proceedings of the 2009 IEEE/WIC/ACM International Joint Conference on Web Intelligence and Intelligent Agent Technology-Volume 01*, pp. 301–304. IEEE Computer Society.

SIMULATION AND ANALYSIS OF THE AIRCRAFT CORROSION CONTROL FACILITY AT THE CORPUS CHRISTI ARMY DEPOT

Paul Delimarschi

Department of Systems Engineering
United States Military Academy
West Point, NY, MOLDOVA
Paul.Delimarschi@usma.edu

Jonathan Griffith

Department of Systems Engineering
United States Military Academy
West Point, NY, USA
Jonathan.Griffith@usma.edu

Mitchell Howard

Department of Systems Engineering
United States Military Academy
West Point, NY, USA
Mitchell.Howard@usma.edu

Sean McBryde

Department of Systems Engineering
United States Military Academy
West Point, NY, USA
Sean.McBryde@usma.edu

Gene Lesinski

Department of Systems Engineering
United States Military Academy
West Point, NY, USA
Eugene.Lesinski@usma.edu

Abstract: Corpus Christi Army Depot (CCAD) is the principle agent for Army rotary wing aircraft depot maintenance and is constructing a new \$32.4M paint facility to support their mission. In this research, a discrete-event simulation with custom Excel user interface is developed to model production capacity and support production “what if” analysis for this new facility. The user interface allows users, unfamiliar with the simulation modeling language, to change major production parameters of interest and conduct analysis to examine the impact of production factors on key metrics. A multi-criteria quantitative value model is developed and integrated with the discrete-event simulation to quantify the value of each production scenario to support production decisions. A full factorial design of experiments, conducted across five production factors, identifies fourteen Pareto efficient production scenarios. These Pareto efficient solutions provide CCAD several production options that can be used to conduct cost-value trade space analysis.

Keywords: Multi-Criteria Analysis, Cost-Value Analysis, Pareto-Efficiency

1 INTRODUCTION

1.1 Background

The Corpus Christi Army Depot (CCAD) leads the Department of Defense in helicopter maintenance, repair, recapitalization, and overhaul capability. CCAD aims to restore rotary wing aircraft and components and return them to the Department of Defense and other government organizations with uncompromising

quality, at the lowest possible cost, in the shortest amount of time possible. Though CCAD is not the only depot of its kind, it is the principal agent for rotary wing depot repair. Today CCAD operates on the Naval Air Station located in Corpus Christi, TX. The facilities and equipment are valued at over \$746 million. As one of south Texas' largest industrial employers, CCAD employs more than 5500 personnel and contractors providing an overall economic impact of more than \$1.14 Billion to the local community (Best and Neil 2010). Due to shrinking Defense budgets and decreased new aircraft purchases, CCAD has seen an increase in their workload. The continual growth in the demand for upgraded and refurbished aircraft such as the UH-60, its Air Force variant, and the AH-64 has created several production challenges for CCAD. These challenges include bottlenecks, backlogs, equipment maintenance issues, worker shortages, and space availability. In December of 2013, CCAD broke ground on a new modernized \$32.4M Aircraft Corrosion Control Facility (ACCF). The purpose of the recently built ACCF is to increase throughput and reduce aircraft time in system, continue to meet the state's air permit requirements by consolidating paint operations under one roof, and minimize costly downtime associated with aircraft depot refurbishment. The ACCF features six paint booths, a cleaning bay, and an alodine treatment bay - all large enough to fit a CH-47 Chinook helicopter (Rox 2017).

1.2 Model Motivation and Problem Statement

The CCAD motivation for developing a model of their ACCF facility revolves around mitigating inefficiencies while increasing revenue. Making their new paint facility as efficient as possible will enable CCAD to meet their larger mission of increasing the number of aircraft and rotor blades they recapitalize each quarter and thus increasing their revenue for each fiscal year. Additionally, increased throughput and reduced time in system will improve the operational availability of the rotary wing aircraft fleet. Previous CCAD modeling and analysis efforts focused on rotor blade refurbishment and aircraft recapitalization (Green et al. 2015, Harvey et al. 2016). **The purpose of this research effort is to develop a discrete event simulation model of CCAD's new ACCF and an accompanying multi-criteria value model in order to analyze production scenarios and support production decisions.**

2 METHODOLOGY

Figure 1 highlights our methodology and general approach. The initial phase focused on identifying the scope of the effort and problem definition. During this stage research and stakeholder analysis is conducted to thoroughly understand the CCAD paint process, capture desired model features, identify and gather necessary supporting data, and identify the key output metrics of interest.



Figure 1. Research Methodology

In the next phase, a discrete event simulation model and supporting multi-objective decision support tool are developed and integrated. Model development was preceded by formulation of an architecture that integrated the user interface, discrete event simulation model, decision support tool, and output module. Model verification and validation are then conducted to ensure the “model is right” and the “right model” was built. In the final phase a full factorial design is conducted to enumerate the possible production scenarios. Each production scenario was run utilizing the developed model and analyzed to identify key production factor insights.

3 PROBLEM DEFINITION

3.1 Paint Process Flow

Figure 2 highlights the process flow for the CCAD paint process. The paint facility processes aircraft shells, rotor blades, and major aircraft component containers. The routing, processing steps, required workers, and process times are different for each of the three processed entities. Aircraft shells are routed through the paint facility twice. The first pass consists of 13 tasks with an estimated average total processing time of 14 days. The second pass consists of only 6 tasks with an estimated average total processing time of 7 days. Between the first and second pass, the aircraft shells are processed in another facility before returning. Aircraft shells must pass through the cleaning and alodine treatment booths prior to entering the paint booths.

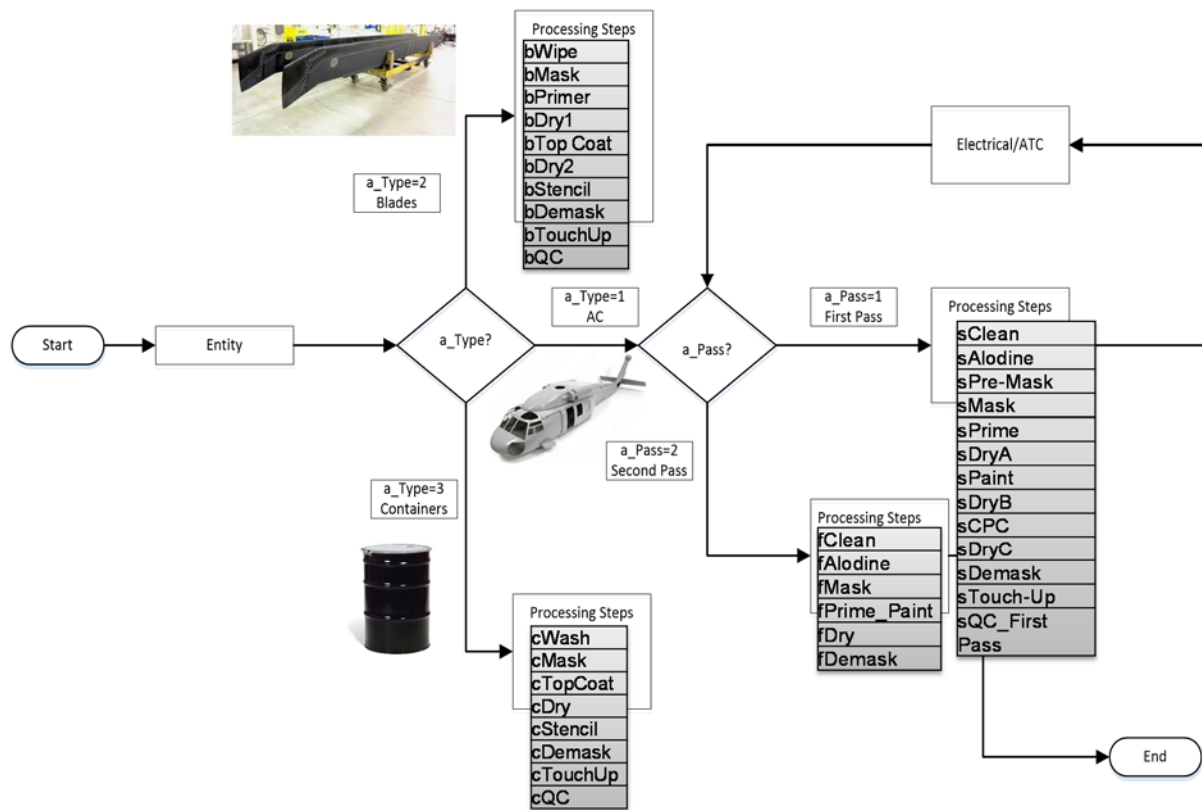


Figure 2. Paint Facility Process Flow

Rotor blade processing consists of 10 tasks with an estimated average processing time of 28.5 hours for a batch of 8 blades. Container processing consists of 8 tasks with an estimated average processing time of 2.7 hours for a batch of 6 containers. Containers and blades are only processed once in the paint facility and do not require routing through the cleaning and alodine treatment booths. Aircraft shells have priority for processing over blades and containers. Each workload task requires an established number of workers. The total number of workers per shift are limited in quantity - resulting in key resource constraints. Additionally, there are a maximum six available paint booths that can only process one aircraft or one batch of blades or containers at a time - resulting in processing capacity constraints.

3.2 Key Performance Metrics

In *Optimizing Factory Performance*, Ignizio identifies a large collection of production metrics commonly used to improve manufacturing efficiency. A few of these metrics include: throughput, work in progress, utilization rates, resource down times, product wait times, and product time in system (Ignizio 2009).

During the initial stages of the project, CCAD staff members identified several paint facility key performance metrics of interest. The most important stakeholder performance metric is annual throughput. Throughput is the number of aircraft, blades and containers that exit the ACCF per year. The goal is to maximize annual throughput. The second most important stakeholder performance metric was paint bay utilization. In this research data is gathered on individual paint bay utilization as well as average paint bay utilization. The goal is to maximize paint bay utilization. The third most important stakeholder performance metric is worker utilization. Data on shift worker utilization and overall worker utilization is gathered as key model output. Table 1 below highlights the key stakeholder metrics and their priority.

Table 1. Key Performance Metrics

Priority	Metric	Data Collected	Goal
1	Throughput	Annual Aircraft Exits	Maximize Throughput
		Annual Blade Exits	
		Annual Container Exits	
2	Paint Bay Utilization	Individual Bay Utilization	Maximize Bay Utilization
		Average Bay Utilization	
		Workload-Bay Utilization	
3	Worker Utilization	Average Worker Utilization	Maximize Worker Utilization
		Shift-Worker Utilization	

3.3 Desired Model Features

Constantine and Lockwood note that a well-designed user interface allows people who understand the problem domain to use the application or tool without having to read manuals or receive extensive training. Additionally, the better the user interface, the higher the probability that the model or application will be used (Constantine and Lockwood 1999). During the course of stakeholder analysis, desired model and user interface features are elicited to create a model that would be useful and used. CCAD desired a model in which production factors and production scenarios could be modified and run by users with no ProModel software programming knowledge. The specific model requirements are that it shall:

- Allow adjustment of workload arrival rates
- Allow adjustment of the number of open paint bays
- Allow assignment of workload types to paint bays
- Allow adjustment of the number of shifts and workers/shift
- Allow adjustment of blade and container batch sizes
- Allow adjustment of workload type priorities
- Capture key performance metrics: throughput, worker utilization, paint bay utilization
- Be capable of running 10 scenarios at a time
- Allow a user with no ProModel programming experience to adjust production factors and run the model
- Receive input from and export data to Excel to assist ProModel novice users
- Incorporate a multi-criteria value model that quantifies the “value” of each scenario

4 MODEL

4.1 Model Architecture

Effective, documented software architecture can increase user and developer understanding of the data, components, and interactions of the software model. Additionally, it facilitates reuse, future model extension, and model management (Garlan and Perry 1995). In this research, considerable prior planning of a model architecture (See Figure 3) was conducted to facilitate understanding, incremental development, and verification of the modeling effort. The three major elements of model architecture include an input

module, simulation module, and output module. The module elements highlighted in red will be discussed in further detail in this section. The Input module consists of the user interface worksheet, a shift and worker assignment macro written in VBA, the ProModel shift calendar files, and the processing data worksheet. To satisfy the model requirement that the model shall allow a user with no ProModel programming experience to adjust production factors and run the model, required development of a user interface. The user interface receives production factor changes from the user and pushes those changes to the ProModel base model, ProModel bay assignment subroutine, and ProModel scenario manager.

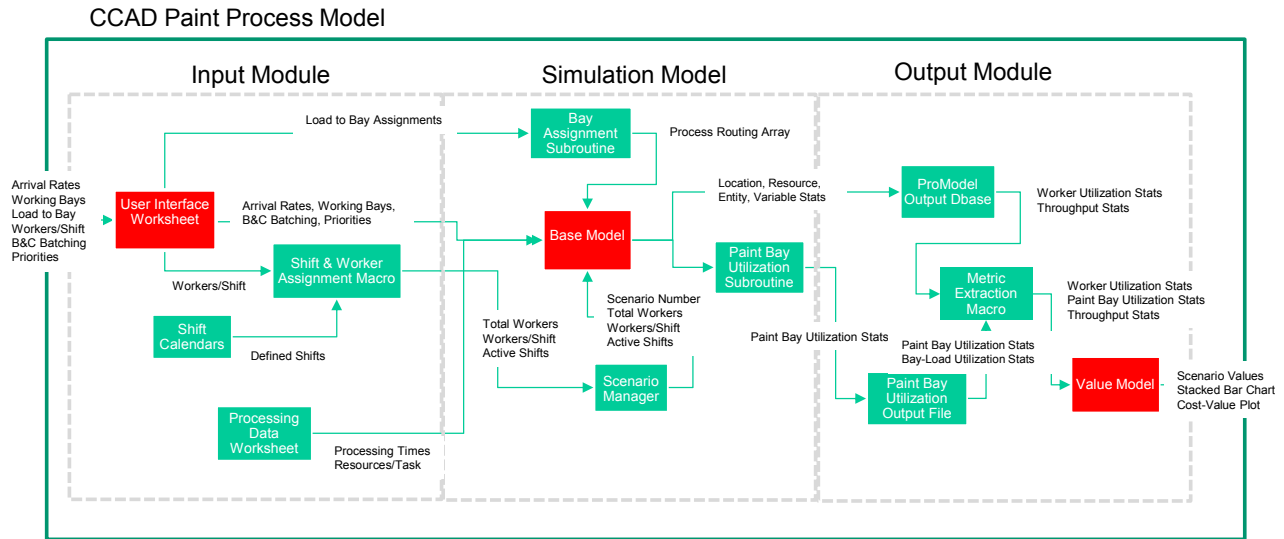


Figure 3. Model Architecture

The simulation module consists of a base ProModel discrete-event model, a ProModel bay assignment subroutine, the ProModel scenario manager, and a ProModel paint bay utilization data subroutine. The base model is discussed in further detail in Section 4.3. The Output module consists of the ProModel output database, a paint bay utilization data file, a metric extraction macro, and a value model. When ProModel runs, all output is saved in an external ProModel output database. The metric extraction macro pulls the desired metrics from the ProModel database and populates a “results” Excel worksheet. The value model uses the key scenario metrics to calculate a multi-criteria value score for each scenario which allows CCAD to compare the relative merit of competing scenarios.

4.2 User Interface

Figure 4 highlights the model user interface. The user interface is an Excel worksheet within a larger workbook. Rows 2-4 allow adjustment of the annual arrival rate of aircraft, blades, and containers. Row 6 allows the user to designate the number of paint bays that are open for operation. Rows 15-20 allow the user to specify bay specific workload. Rows 23-25 in the user interface allow the user to specify the number of shifts and workers per shift. Rows 35-36 specify the batch size for blades and containers. Finally cells E39-E41 allow the user to adjust the relative priority of aircraft, blade, and container processing. The “Update Model” command button pushes this information to ProModel through VBA and ActiveX.

	A	B	C	D	E	F	G	H	I	J	K	L
1			Scenario 1	Scenario 2	Scenario 3	Scenario 4	Scenario 5	Scenario 6	Scenario 7	Scenario 8	Scenario 9	Scenario 10
2	Annual AC	63	1301 125% 100%	76	135% 130% 105%	45	103% 65% 60%	45	105% 100% 100%	60	110% 105% 100%	45
3	Annual Container	4027	135% 120% 115%	5034	135% 130% 120%	2618	103% 65% 60%	2618	105% 100% 100%	3826	110% 105% 100%	2618
4	Annual Blades	1367	135% 120% 115%	1709	135% 130% 120%	821	103% 65% 60%	821	105% 100% 100%	1299	110% 105% 100%	821
5												
6	Open Booths	6	5	3	5	3	4	4	4	4	5	4
7	Paint1	Open	Closed	Closed	Closed	Closed	Closed	Closed	Closed	Closed	Closed	Closed
8	Paint2	Open	Open	Closed	Closed	Closed	Closed	Closed	Closed	Closed	Open	Closed
9	Paint3	Open	Open	Closed	Open	Closed	Open	Open	Open	Open	Open	Open
10	Paint4	Open	Open	Open	Open	Open	Open	Open	Open	Open	Open	Open
11	Paint5	Open	Open	Open	Open	Open	Open	Open	Open	Open	Open	Open
12	Paint6	Open	Open	Open	Open	Open	Open	Open	Open	Open	Open	Open
13												
14	Bay/Assignments	Assign Load to Bay	Assign Load to Bay	Assign Load to Bay	Assign Load to Bay	Assign Load to Bay	Assign Load to Bay	Assign Load to Bay	Assign Load to Bay	Assign Load to Bay	Assign Load to Bay	Assign Load to Bay
15	Paint1	Containers	Containers		Containers			Blades Container	Blades Container		Containers	
16	Paint2	Blades	Containers		Containers			Blades Container	Blades Container		Containers	
17	Paint3	Everything	Blades		Containers			Blades Container	Blades Container		Containers	
18	Paint4	AC	AC	Everything	AC	Everything	AC	AC	AC	Everything	AC	Blades
19	Paint5	AC	AC	Everything	AC	Everything	AC	AC	AC	Everything	AC	AC
20	Paint6	AC	AC	Everything	AC	Everything	AC	AC	AC	Everything	AC	AC
21												
22	Workers	90	60	26	32	30	30	48	48	48	48	48
23	Shift 1 Workers	30	20	10	16	10	10	16	16	16	16	24
24	Shift 2 Workers	30	20	8	16	10	10	16	16	16	16	24
25	Shift 3 Workers	30	20	8	16	10	10	16	16	16	16	24
26												
27												
28												
29	UPDATE MODEL	1-30	1-20	1-10	1-16	1-10	1-10	1-16	1-16	1-16	1-16	1-24
30		31-60	21-40	11-18	17-32	11-20	11-20	17-32	17-32	17-32	17-32	25-48
31		61-90	41-60	19-26	33	21-30	21-30	33-48	33-48	33-48	33-48	49
32												
33												
34	Batching											
35	Blades	8	8	8	8	8	8	8	8	8	8	8
36	Containers	6	6	6	6	6	6	6	6	6	6	6
37												
38	SHIFT FILE NAMES											
39	ccadshifts1.pmcsl				AC Priority	99						
40	ccadshifts2.pmcsl				Blade Priority	50						
41	ccadshifts3.pmcsl				Container Priority	50						
42	ccadNotWorking.pmcsl											
43												

Figure 4. Model User Interface

4.3 Discrete Event Simulation

Discrete event and system dynamics are two common and distinct approaches to simulation modeling. System dynamics modeling is typically focused on strategic level problems where state changes are continuous. While discrete event simulation state changes occur at discrete points in time and the focus is on tactical-level, entity specific information (Brailsford and Hilton 2001). In this research, a discrete event simulation is used in order to capture tactical level, time dependent, information about resources, workloads, and processes. Figure 5 below highlights the ProModel discrete event simulation. The discrete event simulation model includes a blueprint of the ACCF as the model background. The key locations of interest are annotated on the production layout graphic: cleaning booth, alodine booth, paint booths 1-6, and holding areas for workload types. Runtime model graphics include a depiction of task specific workers, process task indicators, and graphics for each workload type. The model also includes a scoreboard that tracks the quantity of workload types in the system at any time and workload types processed by the system. The scoreboard also displays the average time in system for each workload type.

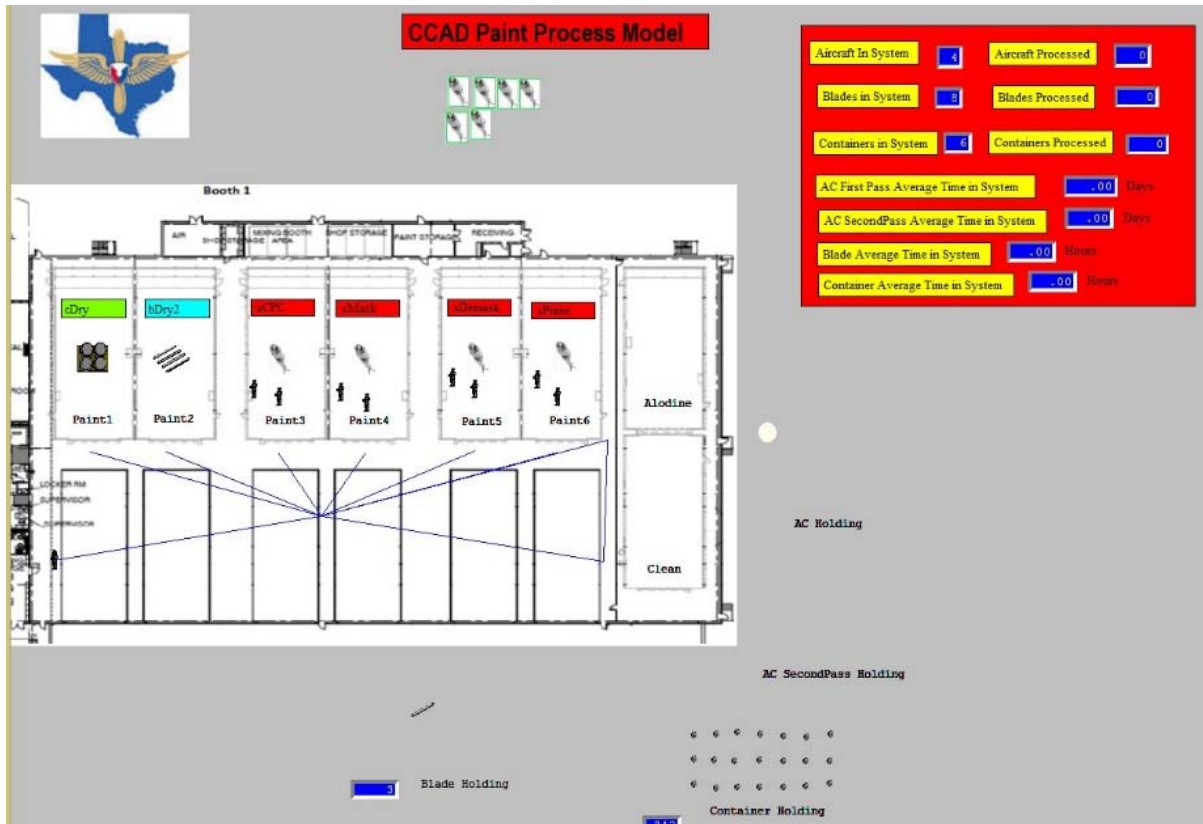


Figure 5. Discrete Event Simulation of CCAD Paint Process

4.4 Multi-Criteria Value Model

Value focused thinking is used to generate, compare and analyze solution alternatives with an emphasis on what is of value to the stakeholder. Typically, multiple problem or solution criteria are of interest (Keeney & Keeney, 2009). The user interface in this research is designed so the user can enter and run up to 10 production scenarios at a time while capturing the key outputs of throughput, bay utilization, and worker utilization. A quantitative value model is developed and integrated with the discrete-event simulation to assist CCAD in comparing the relative value of these competing production scenarios. The key performance output metrics are combined to formulate a production scenario value score.

$$\text{Scenario Value} = w_1 * \text{Throughput} + w_2 * \text{Bay Utilization} + w_3 * \text{Worker Utilization} \quad (1)$$

The production scenario value score (see Equation 1) is a weighted combination of throughput, bay utilization, and worker utilization. The weights $w_1 - w_3$ are adjustable by the user to modify the relative importance of throughput, bay utilization, and worker utilization.

5 VERIFICATION AND VALIDATION

Throughout the development of the simulation model several techniques were used to verify and validate the model with the goal of producing a product that is accurate and credible. The simulation itself is an approximate imitation of the real-world Aircraft Corrosion Control Facility (ACCF). Through verification and validation developers ensure the “Model is right” and “Right model” is developed (Sargent 2013). To verify the model a combination of graphics, counters, and debug statements are incorporated within the

CCAD paint facility model. Several counters were placed in the model at key locations to monitor aircraft, blade and container arrivals, quantities in the system, number at key holding areas, and number processed by the system. Incorporation of animation, graphics, and counters assisted in verifying correct entity routing as well as resource usage. Use of entity graphics also allowed verification of operational bays and workload assigned to bays. A series of labels are incorporated that indicate the particular task being performed on each entity while they were processed in the paint booths which allowed verification that the correct sequence of tasks were performed on each entity. Debug and display statements were incorporated throughout the ProModel code to ensure the entity task processing times were correct. Lastly, the model was observed during runtime to verify that resources were not working outside of scheduled shift times (i.e. weekends and break times). To validate the model, historical data and stakeholder feedback were utilized to assess if the model is an accurate representation of “reality.” CCAD analysts repeatedly ran the model using historical data from fiscal years 2012 through 2015 but since the new facility is markedly different than the previous facility in size and layout, the analysts provided a subjective assessment of the model validity.

6 MODEL RESULTS AND ANALYSIS

6.1 Full Factorial Design of Production Scenarios

A full factorial design is a planned set of tests on the response variables with one or more factors with all combinations of levels. Through a full factorial design the factor main effects can be identified as well as set the stage for analysis of variance (ANOVA) and regression analysis. ANOVA is used to show which factors are significant while regression derives appropriate coefficients for predictive models (Goupy and Creighton 2007). Through stakeholder analysis five key production factors were identified as well as the typical level settings for each factor. Major production factors and associated levels include: arrivals (low, med, high), open bays (3-6), bay-workload assignment strategies (aircraft heavy, balanced), shifts (1-3), and workers per shift (10, 16, 20, 24, 30). Table 2 highlights these five factors as well as their typical levels.

Table 2. Production Scenario Factors and Levels

Input Factor	Type	Number of Levels	Level Description
Arrivals	Categorical	3	Low, Medium, High
Open Bays	Integer	4	3,4,5,6
Bay-Workload Assignment	Categorical	2	Aircraft Heavy, Balanced
Shifts	Integer	3	1,2,3
Workers per Shift	Integer	5	10,16,20,24,30

A full factorial design of these factors and levels results in 360 potential production scenarios which are processed by the model while capturing the key performance metrics for each.

6.2 Results and Analysis

Regression analysis of the full factorial, simulation model outputs resulted in predictive models for production scenario value based upon number of operational bays, workload-bay assignment strategy, number of shifts, and workers per shift. Equations (2) through (4) are the predictive models for production scenario value for high, medium, and low workload arrival rates.

$$V_{High} = 19.83 + 4.90x_{Bays} - 1.46x_{Work-Assign} + 2.07x_{Shifts} - .78x_{Workers/Shift} \quad (2)$$

$$V_{Medium} = 20.26 + 4.63x_{Bays} - 1.32x_{Work-Assign} + 1.78x_{Shifts} - .77x_{Workers/Shift} \quad (3)$$

$$V_{Low} = 18.55 + 4.26x_{Bays} - 1.14x_{Work-Assign} + 3.14x_{Shifts} - .76x_{Workers/Shift} \quad (4)$$

V_i represents the production scenario value for workload arrival rate i . As presented and discussed earlier, production scenario value is a weighted combination of throughput, bay utilization and worker utilization. Statistical analysis indicates that number of bays and workers per shift are the most significant factors in predicting production scenario value. Table 3 displays the correlation coefficients (R^2) and analysis of variance (ANOVA) results for the three linear regression models.

Table 3. Model Key Statistics

Equation	R^2	F Statistic	P-Value
V_{High}	.61	45.59	<.0001
V_{Medium}	.64	42.90	<.0001
V_{Low}	.66	56.94	<.0001

A production cost equation was developed to quantify the annual costs associated with each production scenario and to facilitate subsequent cost-value analysis. The cost equation (5) assumes a salary differential between shifts, 2400 annual shift work hours, and a \$1M annual operations and maintenance cost per operational bay.

$$Scenario_{Cost} = 2400(\$30Workers_{1st\ Shift} + \$40Workers_{2nd\ Shift} + \$50Workers_{3rd\ Shift} + \$1Mx_{Bays}) \quad (5)$$

After calculating production scenario costs, a Cost-Value analysis was conducted to highlight Pareto efficient scenarios within each workload arrival category (High, Medium, Low). Figure 6 highlights five production scenarios for low workload arrival rates that are Pareto efficient. All other solutions are dominated since they provide less value at an increased annual cost. A Pareto approach provides CCAD several options to choose from and they can conduct trade space analysis to decide if the additional value is worth the additional cost. Table 4 highlights the production scenario factors, value, and cost of these Pareto efficient solutions.

Table 4. Pareto Efficient Production Solutions for Low Workload Arrival

Scenario	Production Scenario Factors	Value	Cost
205	6 Bays: Balanced: 3 Shifts: 10 Workers/Shift	56.61	\$8,880,000
229	5 Bays: AC Heavy: 3 Shifts: 10 Workers/Shift	52.30	\$7,880,000
303	5 Bays: AC Heavy: 2 Shift: 10 Workers/Shift	49.72	\$6,680,000
133	4 Bays: Balanced: 1 Shifts: 10 Workers/Shift	45.62	\$4,720,000
167	3 Bays: Balanced: 1 Shifts: 10 Workers/Shift	28.83	\$3,720,000

Table 5 highlights the Pareto efficient production scenarios for High, Medium, and Low workload arrival rates. Note that all Pareto efficient production solutions utilize 10 workers per shift. Additionally, if 3 or 4 paint bays are available, 1 shift of 10 workers results in Pareto efficient production solutions.

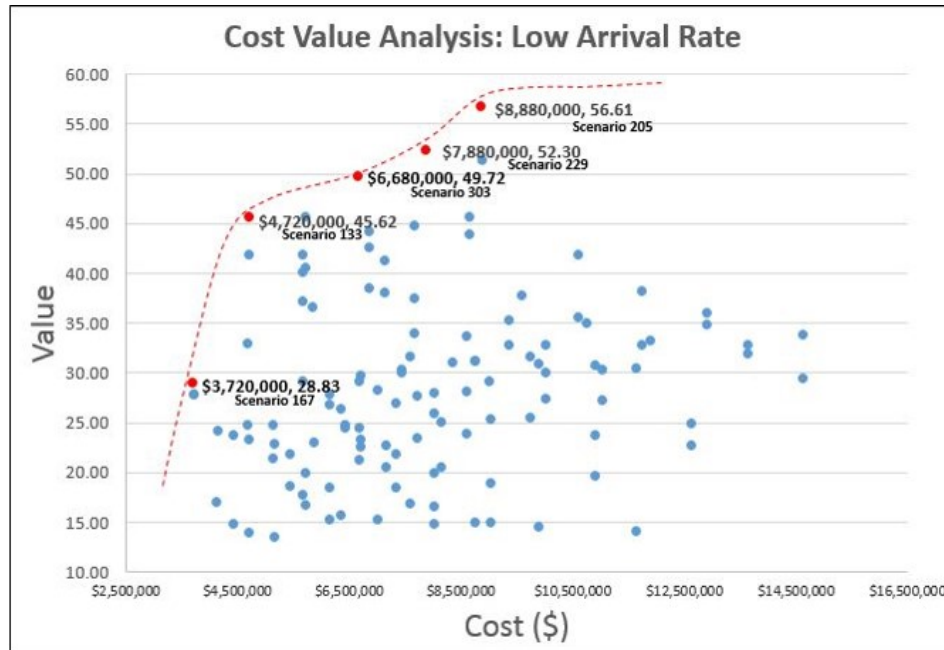


Figure 6. Cost versus Value Plot for Low Workload Arrival Production Scenarios

Table 5. Pareto Efficient Production Scenarios (High, Medium, Low Arrival Rates)

Arrival Rate	Production Scenario Factors	Value	Cost
High	5 Bays: AC Heavy: 3 Shifts: 10 Workers/Shift	55.18	\$7,880,000
	5 Bays: Balanced: 1 Shift: 10 Workers/Shift	52.83	\$5,720,000
	4 Bays: AC Heavy: 1 Shifts: 10 Workers/Shift	39.48	\$4,720,000
	3 Bays: AC Heavy: 1 Shifts: 10 Workers/Shift	33.01	\$3,720,000
Medium	5 Bays: AC Heavy: 3 Shifts: 10 Workers/Shift	51.90	\$7,880,000
	5 Bays: AC Heavy: 2 Shifts: 10 Workers/Shift	50.52	\$6,680,000
	5 Bays: Balanced: 1 Shifts: 10 Workers/Shift	50.46	\$5,720,000
	4 Bays: AC Heavy: 1 Shifts: 10 Workers/Shift	37.15	\$4,720,000
	3 Bays: AC Heavy: 1 Shifts: 10 Workers/Shift	30.53	\$3,720,000
Low	6 Bays: Balanced: 3 Shifts: 10 Workers/Shift	56.61	\$8,880,000
	5 Bays: AC Heavy: 3 Shifts: 10 Workers/Shift	52.30	\$7,880,000
	5 Bays: AC Heavy: 2 Shift: 10 Workers/Shift	49.72	\$6,680,000
	4 Bays: Balanced: 1 Shifts: 10 Workers/Shift	45.62	\$4,720,000
	3 Bays: Balanced: 1 Shifts: 10 Workers/Shift	28.83	\$3,720,000

7 CONCLUSIONS

In this research, a discrete-event simulation with custom Excel user interface is developed and integrated with a multi-criteria value model to support paint production facility “what if” analysis. The Excel user interface allows users, unfamiliar with the simulation modeling language, to change major parameters of interest and conduct “what if” analysis to examine the impact of production factors on key metrics of interest. A multi-criteria value model is integrated within the model architecture to allow users to quantitatively compare the value of competing production scenarios to support production decisions. The multi-objective qualitative value model combines the three production metrics of interest (throughput, bay utilization, and worker utilization) into a production scenario value score. A full factorial design of experiments conducted across five production factors, identifies fourteen Pareto efficient production scenarios. Analysis also reveals that the number of bays and workers per shift are the most significant

factors in predicting production scenario value. All identified Pareto efficient production solutions utilize 10 workers per shift. If 3 or 4 paint bays are available, 1 shift of 10 workers results in Pareto efficient solutions. The identified Pareto efficient solutions provide CCAD several production options that can be used to conduct cost-value trade space analysis.

REFERENCES

- Best, Z. and T. Neil. April 2010. "UH-60 Recapitalization (Recap) - Black Hawk's Cornerstone for Fleet Sustainment". *Army AL&T* pp. 30-33.
- Brailsford, S. C. and N.A. Hilton. 2001. "A comparison of discrete event simulation and system dynamics for modelling health care systems".
- Constantine, L. L., and L.A. Lockwood. 1999. *Software for use: a practical guide to the models and methods of usage-centered design*. Pearson Education.
- Garlan, D. and D.E. Perry. 1995. "Introduction to the special issue on software architecture". *IEEE Trans. Software Eng.* vol. 21(4), pp. 269-274.
- Goupy, J. and L. Creighton. 2007. *Introduction to design of experiments with JMP examples*. SAS Publishing.
- Green, N., D. Jaye, S. Kerns, and G. Lesinski. 2015. "Modeling and Analysis of the Rotor Blade Refurbishment Process at the Corpus Christi Army Depot". *Industrial and Systems Engineering Review* vol. 3(2), pp. 124-130.
- Harvey, J., H. McCormick, K. O'Brien, and T. Ritchie. 2016. "Modeling and Analysis of the UH-60 Refurbishment Process at the Corpus Christi Army Depot (CCAD)". *Proceedings of the Annual General Donald R. Keith Memorial Conference*.
- Ignizio, J. 2009. *Optimizing factory performance: cost-effective ways to achieve significant and sustainable improvement*. McGraw Hill Professional.
- Keeney, R. L. and R.L. Keeney. 2009. *Value-focused thinking: A path to creative decision making*. Harvard University Press.
- Rox, B. 2017. "CCAD Breaks Ground for New Helicopter Painting Facility". *The United States Army*. Corpus Christi Army Depot. CCAD Public Affairs.
- Sargent, R. G. 2013. "Verification and validation of simulation models". *Journal of Simulation* vol. 7(1), pp. 12-24.

STUDENT CAPSTONE CONFERENCE

2017

AGENT BASED MODELING

- Page 69 Julie Zhou and Chenyun Zhang
Governor's School for Science & Technology
A Simulation On The Effect Of A Major World War On The Population Of The World
- Page 79 Wessam Elhefnawy
Old Dominion University
Make my Neighborhood Safe Again : An Agent Based Model for Burglary Crime Prediction and Capture its Patterns
- Page 91 Saturnina Nisperos, Sonali Kakde and Frederic McKenzie
Old Dominion University
An Agent-based Simulation of the Impact of Yogic Breathing Adoption in Hampton Roads

A SIMULATION ON THE EFFECT OF A MAJOR WORLD WAR ON THE POPULATION OF THE WORLD

Julie Zhou
Governor's School for Science & Technology
Hampton, VA
julie.zhou@nhgs.tec.va.us

Chenyun Zhang
Governor's School for Science & Technology
Hampton, VA
chenyun.zhang@nhgs.tec.va.us

ABSTRACT

The growth in nuclear weapon testing, political instability, and increase in terrorist activity make war and casualties a global concern. The purpose of this study was to simulate changes in the world population under a variety of war scenarios. The current population and growth rate of each country from the CIA World Factbook and historical war casualty data were put into a map shapefile. The model applies the growth rate and, if a battle took place, the coefficients derived from Lanchester's Theory of Warfare to the population, updating the map. The user can specify the initial soldier amount for both sides and the weapon efficiency increase since the historical wars. This process is repeated for each tick, which is count of days passed in the simulation. While the population will inevitably decrease in the case of a war, the different coefficients of war have varying impacts on the population.

Keywords: Lanchester's Theory of Warfare, War, Population, NetLogo

1 INTRODUCTION

The idea of a "World War III" has been an active area of speculation and interest for many people, and the speculation is for a good reason. In recent years, North Korea announced their nuclear missile testing and willingness to use these weapons to counter foreign threats, resulting in heightening tensions with multiple countries, such as the United States (Gale and Lee 2016). Growth in nuclear missile testing is not the only problem that the world is facing, but terrorist acts have also increased within the last century and have quickly developed into a global problem. Especially after the September 2001 attacks, people around the world have recognized the threat of terrorism on civilization (Vertigans 2010). Furthermore, improvements in weaponry enables weapons to affect a wider area and cause heavier damage to a population. These developments have led to the possibility of a modern world war taking place, which could lead to mass military and civilian casualties.

However, in the case of an outbreak of a major war, war simulations can help gauge the effects of a conflict. They can also test the outcome of different scenarios. In a recent study, two researchers incorporated the Lanchester's Theory of Warfare in their simulation of the American Civil War battle of Pickett's Charge in order to see whether or not the Confederate army could have won in three different situations (Armstrong and Sodergren 2015). Although there are many factors involved in warfare that are hard to trace, war simulations are still used to analyze past conflicts; this analysis can provide more foresight into future conflicts (Helmbold 1961). Simulations that use the Lanchester's theory do not aim to model war combat because of its complexity and, as mentioned before, difficulty in measuring war factors, such as movement of troops, but rather they aim to measure combat attrition, which takes in the casualties of one side and compares it to the size of its forces (Perry 2011). Similarly, this study implements the Lancaster Theory of Warfare in order to replicate battles of World War I and World War

II and also calculate casualties based on user-defined efficiencies and soldier amount. In order to observe changes in population of different countries during a war, this study incorporates past world war data and user-input war factors to produce different war casualties.

2 OBJECTIVES AND OVERVIEW

A simulation was created to predict the condition of the world population during and after the outbreak of a modern world war. Patterns in growth or decline in the different populations, and which countries would change the most dramatically, were observed in the results. These changes led to results that could be possible consequences of a major world war.

The model has two different options: one based on World War I and another based on World War II. A majority of the major war battles were implemented into the model, and the casualties in each battle were used. The user is able to define what day the war will start on, what percentage of a country would become soldiers, the percentage that the weapon efficiency increased from the war that data was gathered from, and which country's population to track. After running the simulation, the map projection displays various color changes for the countries, which represent shifts in population.

3 RESEARCH

For the map projection of the model, a free country Geographic Information System (GIS) shapefile from the Natural Earth vector and raster map database was used to visualize the changes in the population of each country. The current population and the yearly growth rate for each country were taken from the Central Intelligence Agency (CIA) World Factbook and were joined into the country shapefile through the open-source QGIS software. The population growth rate uses the annual average percent change, which is the difference of the current population and past population all divided by the past population and the outcome of a surplus or deficit in births in comparison to deaths. The population growth rate also uses the annual average percent change to the net migration of the country, which is the difference between the immigrants and emigrants of a country. These data served as the starting attributes of each country in the model and provided the basis for how the countries would develop without war (Figure 1). Casualty data, which includes the people killed and wounded, from the battles in World War I and World War II were taken from online databases. These data were also added into the shapefile under columns that were titled by the date of the battle, such as F08231914; F08231914 stands for August 23, 1914. Rough sizes of the military of each country involved in World War I and World War II were added into the file for the replication of past battles and as comparison to the user-input war factors. If the country was not in World War I or World War II, the program defaulted the size of the military to 1,000 in order to avoid any cases of divide by zero.

name_long	Pop	Pop_grow
Australia	22751014	1.0700000...
Austria	8665550	0.5500000...
Azerbaijan	9780780	0.9600000...
Burundi	10742276	3.2700000...
Belgium	11323973	0.7600000...
Benin	10448647	2.7800000...
Burkina Faso	18931686	3.0300000...
Bangladesh	168957745	1.6000000...
Bulgaria	7186893	-0.5800000...
Bahamas	324597	0.8400000...
Bosnia and...	3867055	-0.1300000...
Belarus	9589689	-0.2000000...
Belize	347369	1.8700000...
Bolivia	10800882	1.5600000...
Brazil	204259812	0.7700000...

Figure 1: QGIS Shapefile Attributes

4 NETLOGO

NetLogo is an agent-based programming software and relies on patches, which are the pixels that the model is made up of and turtles, which are the agents that move on the grid. While the simulation did not utilize any turtles, it was heavily dependent on patches. Within NetLogo, patches can have their own variables; the variables are used to characterize each patch. In the simulation, each patch was identified as a part of a country, which allowed it to contain the name of the country, the population of the country, and the population growth rate of each country. The data that are contained within the variables that the patch owns are extracted from the shapefile, which are implemented through GIS extension from NetLogo.

NetLogo also provides different extensions, such as array, matrix, and GIS. The GIS extension within NetLogo allowed for the implementation of the shapefile and projection file. While agent-based modelling is often more suited for large scale population modeling, the ability to load vector dataset and raster dataset through shapefiles and ASCII grid files (GIS Extension) allowed for the implementation of past data. Thus, the program used the extension to load the projection file, which created a map, and the datasets within the shapefile. Due to the apply-coverage method, each column of the shapefile can be designated to a variable, which is owned by every patch. These variables were used to display and change the population.

5 DISPLAY

In order to project the populations of a country on the map, the patch-owned variables were used. When the population data of a patch falls within the range for a specific color, the patch will turn to that specific color, which is described by the key that is next to the map (Figure 2). The number next to the color is the upper bound, and the lower bound is found through the color above. The colors fluctuate between purple and pink. The colors ascend from a lighter shade to a darker shade; however, if the two colors are on the same scale, such as violet + 1 and pink + 1, the pink color represents a larger population. The key does not follow a specific trend because some of the population cluster around a specific color. Furthermore, China and India skews the key to larger numbers. This causes the rest of the map to have the same color

because the range that China and India skews is very large. To allow for flexibility, users can define different variables, such as the percentage of the population that will participate in the war that will impact the outcomes of the simulation.

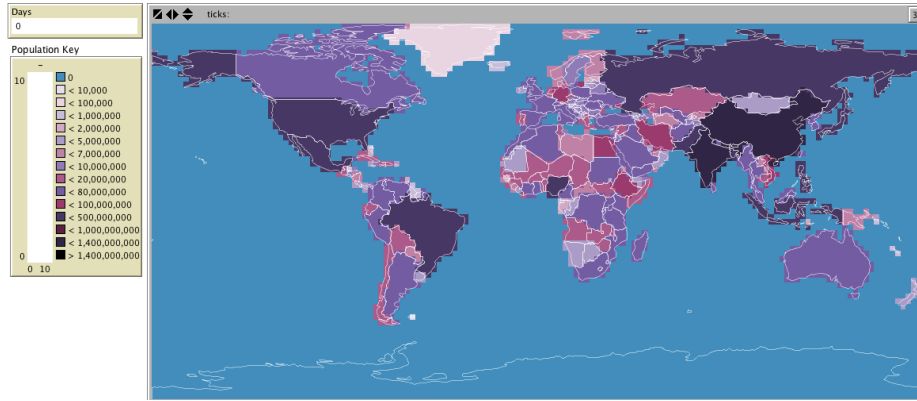


Figure 2: The map projection with the key

6 SIMULATION

Since there are many scenarios that may happen when a war breaks out, the simulation should be flexible enough to cover the different situations, such as when there are not many soldiers participating in the war. The users can define when the war starts by inputting which day will begin the war. The user can also choose which country to monitor. Through a drop-down bar, the user can select a country and the population of the country will be displayed below it, which allows the user to carefully track one specific country during the war. In order to monitor the population, the program checks to see if the name variable of the patch is equal to the name of the country that the user would like to monitor. If it is, the program assigns the variable that is displayed with the population data that the patch contains. The user can also choose whether the World War I data or the World War II data will be used. These two wars were chosen because the wars were the two large-scale wars that are closest to modern day. Information on the battles and the deaths, therefore, were readily available. After the user selects a war, the program will find the first battle of the war through its property name: F08051914, which is for World War I, and F03171939, which is for World War II.

The program finds the name and assigns the variable, `start_index`, as the column number that the property name was found. The deaths that occurred during each battle can also be adjusted by the user; the user can input numbers for the percentage of the citizens that are soldiers for each country and for the percentage that the weapon efficiency increased since the original battle. These two variables are key components of the Lanchester's Weapon Theory. The theory states that the rate of casualties for the two forces are the effectiveness of the two forces times the numerical strength of the two forces (Watkins n.d.). Thus, the percentage that numerical strength of the forces increased is the percentage of the citizens that are soldiers divided by the original number of soldiers, which is multiplied by the weapon efficiency divided by 100. Since the simulation divides the new number of soldiers by the old number of soldiers, the model decreases the effect that the soldier amount during the past war. This is due to the fact that if the new soldier amount is less than the old soldier amount and if the weapon efficiency stays the same, the death count actually decreases. This number is multiplied by the original death count and added the death count to create a new variable, `new_death`. Through the Lanchester's Weapon Theory, the users are able to manipulate the strength of each country. After these variables are set, the simulation ends the set-up stage.

Once the simulation starts, the program takes the growth rate of each country and multiply it by the current population. The resulting number will be rounded and added to the current population. If the war has started, the program will determine if the day of which the model is on is equal to the day that a battle happened during one of the major wars. The days of the battles are provided in a list of numbers. For example, if a battle happened on the first day of the war, the number zero will be the first number in the list. If the two numbers are equal, the model will obtain the data on the battle by finding the column in the shapefile. The program finds the battle in the shapefile by taking the start_index and adding the array count to it, which is initially zero. Once the death data are obtained, the program calculate the new_death variable and subtract it from the current population. If the resulting number is greater than the current population, the population will be set to zero because a country cannot have a negative population size. After the population is manipulated through the growth rate and the battle, the new population will be implemented. The day count will increase by one, and the population will be manipulated again.

7 RESULTS

In the simulation, the soldier percentage was set to 10%. This was chosen because during World War I, around 4% of people in America went to war and around 19% of people in United Kingdom went to war. In World War II, around 9% of Americans went to war, while the percentage in United Kingdom was very similar. Other countries, such as Germany, that were active in the war had between 0% to 20% percent of the population in the war. Thus, the soldier percent was set to 10%, while different efficiencies were tested.

During the initial stage of the simulation, which was when the day is equal to zero, North America, South America, Asia, and Europe mainly consisted of darker colors; Africa on the other hand, was consisted of some lighter colors (Figure 3).

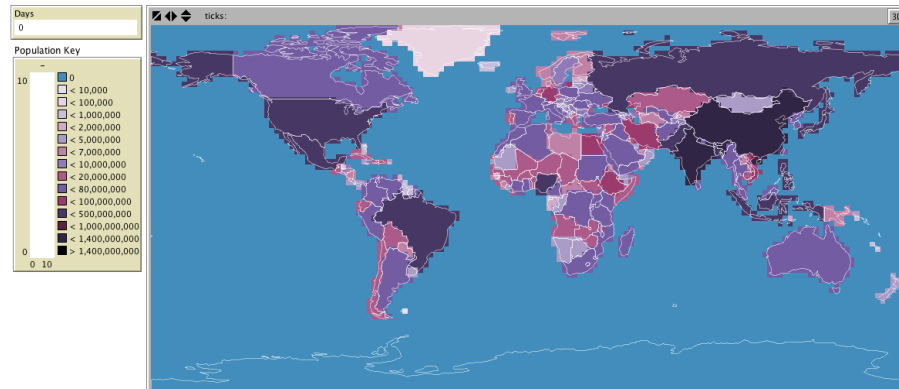


Figure 3: Initial Phase

When the weapon efficiency was set to 50%, and the selected option was World War I, there was not a drastic visual change in the colors of the map when it was stopped at around 1521 days, which was when the war ended (Figure 4). Sweden and some African countries, such as the Democratic Republic of Congo, were a few that got noticeably darker at the end of the simulation. However, there were still some countries, such as Libya, that got lighter in color, which shows that the population in that country decreased. However, the change in the major countries, such as Russia, America, and China, was not displayed through the colors, since the change did not exceed the limit for the color to change.

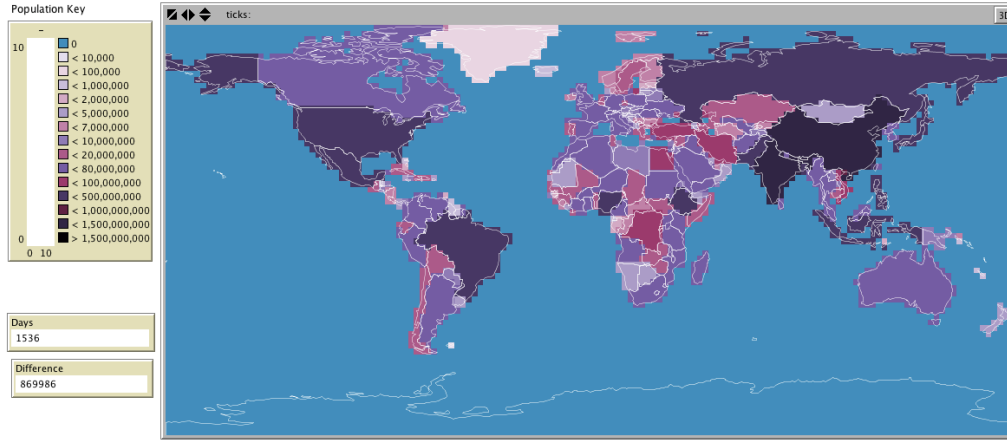


Figure 4: End of simulation with the weapon efficiency of 50% and with World War I data

In the World War II setting, most of the visual change occurred in the Northern African and Western European regions when it was stopped at around 2320 days (Figure 5). The Northern African countries, such as Mali, turned darker, which showed that their population increased. The Western European regions, except for Germany also became darker.

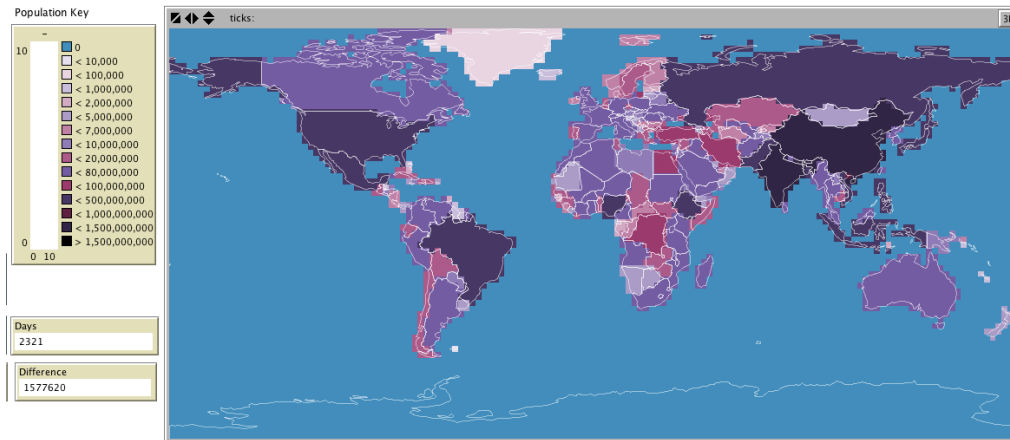


Figure 5: End of simulation with the weapon efficiency of 50% and with World War II data

When the efficiency rate was set to 100% with the World War I dataset, the colors of the map did not dramatically change from the ones with the weapon efficiency of 50%. Due to the fact that the changes in color represent a dramatic change in the population, this stable color trend also occurred within the World War II dataset. Even though the color did not change dramatically, the population of each country were still impacted through the change of weapon efficiency (Figure 6). The changes between the different options do not seem like much but for Japan, the population went down 868,880 from the original 126,919,659 in the World War I option and 1,584,757 in the World War 2 option. Germany had losses of 9,446,042 in World War I scenario and 1,001,349 in the World War II scenario.

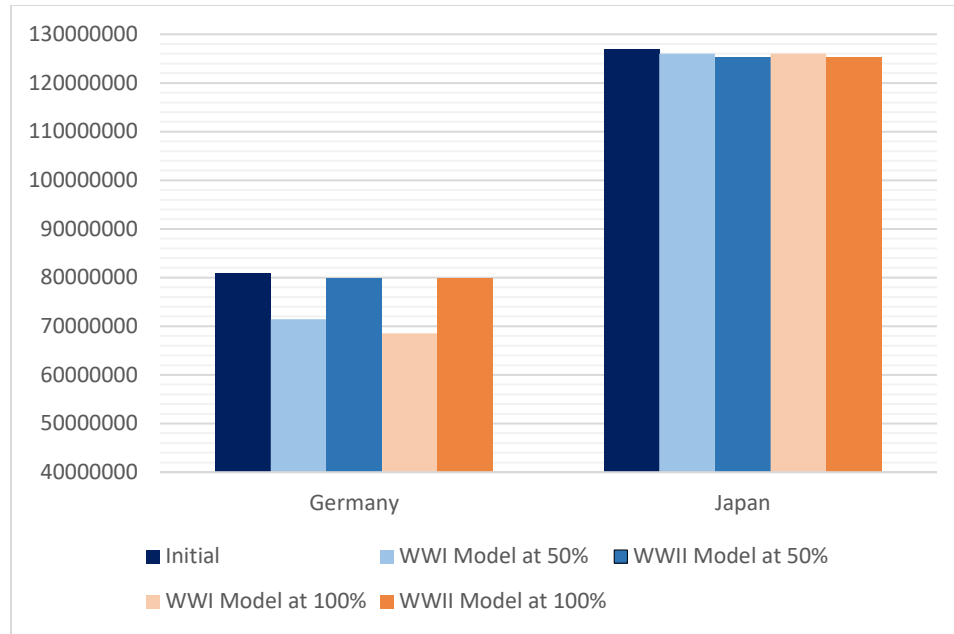


Figure 6: Populations of Japan and Germany with the World War I and World War II options

8 DISCUSSION

This model is a user-friendly, basic simulation that shows the consequences of war on population through different weapon efficiencies and soldier count. It could serve as an aid for the military and first responders to assess if an area is capable of surviving a battle. They can perform what-if analysis with the model by asking questions such as “what if the weapon efficiency increased by 100% since World War I” or “what if all the countries deploy 10% of the population during the war.” The possible consequences shown can also help determine what still needs to be done to further prepare the area in case of a major world war. History teachers could also use this model to help students get a better understanding of the effects of war. In most cases, the number of casualties seem very minute compared to the current population of the world, which skews the perception that students may have on the impact of the two major wars during that time period. However, since the model is based on past world war data, it is unable to display the effects of countries that were not included in those battles. Because not all of the battles of World War II were included and not all the factors of war were taken into account, the model does not provide a simulation of warfare. This may also be the cause of why there was a greater decrease in population in the World War I option, especially for Germany, even though World War II is known for being one of the deadliest conflicts in history. However, Japan suffered more losses in the World War II model. This may be due to the fact that the bombings of Nagasaki and Hiroshima were included in the data set.

Since the colors in the model did not dramatically change with a 50% weapon efficiency and a 100% weapon efficiency, it shows that the population of the world will not completely change. While the colors did not change, there is still a significant impact to the population. Furthermore, since nuclear weaponry were integrated with common weaponry, such as guns, a larger impact can be expected when nuclear weaponry is specifically simulated.

9 FUTURE WORK

In order to develop a more accurate simulation of war, more battles of each world war would need to be added into the model. The current strengths of the militaries of each country would also need to be taken into account, as well as data on migration during war and economic status. Typically, there would be a higher rate of emigrants from a war-torn country, so the migration rate in the model would need to be manipulated to fit a war situation.

Furthermore, the patch issue with the incorrect data must be fixed. Since the patches within NetLogo would cover more than one countries, the data within the patch would be skewed. While the data within most of the patches were fixed by hardcoding the data in, there were still some patches remaining that contained incorrect data. However, this did not significantly impact the map due to the fact that the correct data will quickly override the incorrect data. While the map was not significantly impacted, the ability to graph the data was significantly impacted. Since there were incorrect data points, the graph will suddenly dip when there were no battles. Thus, it provides a false sense that many people died when, in reality, no one did. This can either be fixed by using a different strategy, such as by finding the center of the country, which usually indicates that there are no other countries within the patch, or using another map that fits the shape of the patch.

For the interface, users could specify which countries to go to war with each other. By implementing war theories to determine how they will fight each other, the simulation could then run a battle and display the consequences of it on the user-selected countries. This allows for more possibilities to be shown, which creates a more realistic simulation of what a major world war may be like. Furthermore, through the addition of population centers, the simulation will become more realistic due to how it follows more closely with the real world.

Although it was part of the original plan of this study, due to time constraints, the economic portion was not included. By including the gross domestic product (GDP) per capita, the GDP growth rate of each country, and the trends of war economies, the effects of war on the economy of the world could be simulated. This allows for the factor of how financial needs will impact the outcome of the war. During a war, it is possible that a country's economy is fall apart, which causes their military to become weaker through the lack of weaponry and necessities. This can lead to a larger count of death, which can also end the war earlier.

ACKNOWLEDGMENTS

We would like to thank Dr. Andreas Tolk of MITRE for taking the time to give us excellent advice and guidance on war simulation. We would also like to thank Dr. Bedir and Dr. Woo for taking their time to allow us to explain our project and give their suggestions, which included different approaches, to this study. Lastly, we would like to acknowledge and thank our Environmental Science and Mentor teacher, Dr. Margaret Mulvey, for helping us with the revisions on the paper, allowing us to use her room to work on the model, and coordinating our meetings with Dr. Tolk with our mentorship coordinator, Mrs. Laura Vobrak.

REFERENCES

- Armstrong, Michael J., and Steven E. Sodergren. 2015. Refighting Pickett's Charge: Mathematical Modeling of the Civil War Battlefield. *Social Science Quarterly* 96(4): 1153-1168. <http://onlinelibrary.wiley.com/doi/10.1111/ssqu.12178/full> (accessed October 8, 2016).
- Gale, Alastair, and Carol E. Lee. 2016. "US Agreed to North Korea Peace Talks Before Latest Nuclear Test." *Wall Street Journal*, February 21, 2016. <http://www.wsj.com/articles/u-s-agreed-to-north-korea-peace-talks-1456076019> (accessed June 5, 2016).
- GIS Extension. "NetLogo User Manual version 5.3.1." Center for Connected Learning and Computer-Based Modeling, Northwestern University, Evanston, IL. <http://ccl.northwestern.edu/netlogo/docs/gis.html> (accessed June 28, 2016).
- Helmbold, Robert L. 1961. *Historical Data and Lanchester's Theory of Combat. Part 1*. No. Corg-sp-128. Technical Operations Inc Fort Belvoir VA Operations Research Group. <http://oai.dtic.mil/oai/oai?verb=getRecord&metadataPrefix=html&identifier=AD048095> (accessed October 16, 2016).
- NetLogo Code. <http://pasted.co/cce259dd>
- NetLogo. "Download NetLogo." Center for Connected Learning and Computer-Based Modeling, Northwestern University, Evanston, IL. <https://ccl.northwestern.edu/netlogo/download.shtml> (accessed June 28, 2016).
- Perry, Nigel. 2011. *Applications of Historical Analyses in Combat Modelling*. No. Dsto-Tr-2643. Defence Science and Technology Organization (Australia) Joint Operation Division. <http://oai.dtic.mil/oai/oai?verb=getRecord&metadataPrefix=html&identifier=ADA557493> (accessed October 16, 2016).
- Vertigans, Stephen. 2010. British Muslims and the UK government's 'war on terror' within: evidence of a clash of civilizations or emergent de-civilizing processes?. *The British Journal of Sociology* 61(1): 26-44. <http://onlinelibrary.wiley.com/doi/10.1111/j.1468-4446.2009.01300.x/abstract> (accessed October 8, 2016).
- Watkins, Thayer. n.d. "Lanchester's Theory of Warfare." San Jose State University Economics Department. <http://www.sjsu.edu/faculty/watkins/war.htm> (accessed September 20, 2016).

AUTHOR BIOGRAPHIES

JULIE ZHOU is a senior at the Governor's for Science & Technology in Hampton, Virginia. She is taking rigorous courses in scientific programming and inquiry physics, which are in the Scientific Programming strand. Furthermore, she is taking calculus, research application/ mentorship, and environmental science in the Governor's School. She participates in the York High School's Girls' Varsity Tennis Team, as well as other activities such as National Honor Society and the Green Club. Her email address is julie.zhou@nhgs.tec.va.us.

CHENYUN ZHANG is a senior at the Governor's for Science & Technology in Hampton, Virginia. She is completing coursework in the Governor's School's Scientific Programming strand, which contains classes in scientific programming and inquiry physics. Along with those classes, she is taking multivariable calculus/linear algebra, research application/ mentorship, and environmental science. At

Grafton High School, Chenyun participates in the Girls' Varsity Tennis Team and in multiple honor societies. Her email address is chenyun.zhang@nhgs.tec.va.us.

MAKE MY NEIGHBORHOOD SAFE AGAIN : AN AGENT BASED MODEL FOR BURGLARY CRIME PREDICTION AND CAPTURE ITS PATTERNS

Wessam Elhefnawy
Department of Computer Science
Old Dominion University
Norfolk, VA, USA
welhefna@cs.odu.edu

ABSTRACT

Establish effective methods for crimes prevention have a long history in modern societies. For this purpose, researcher have employed various techniques to analyze and model crime data, from simple regression to data mining methods. However, these methods cannot model actions and behavior of individuals. Agent based modelling is promising method that has ability to model the occurrence of crime and the motivations behind it. This paper presents the construction of an agent-based model (ABM) framework for simulating occurrences of burglary crime. It presents a framework that allows environmental elements to be simulated. A realistic urban environment, based on the real historical crime records, demography census data, and geographical information system (GIS) were constructed, and experiments were conducted to explore the potential of the model to realistically simulate the main processes and drivers within this system. The results are highly promising, demonstrating the potential of this approach for both understanding processes behind burglary crime and developing effective crime prevention strategies.

Keywords: Burglary Crime Modeling, Agent Based Model, Automatic Calibration

1 INTRODUCTION

Burglary crime reduction and prevention is a critical challenge all over the world, with various negative impacts ranging from economy losses to social life wellbeing. The purpose of this study is to provide a realistic prediction system which can help in burglary prevention and reduction. In particular, Agent based modeling (ABM) is employed to find burglary crimes patterns and understand the surrounded environment attributes that impact to the burglary crimes. The central research questions that guided this study are:

1. What motivates burglars to engage in burglary?
2. What factors are considered by burglars during target selection?
3. What deters burglars from burglarizing specific targets?
4. What techniques do burglars use when engaging in burglary?

In order to contribute to the burglary crime prevention using ABM, in this study we employed criminology theorems, demography census data, and GIS information to develop a realistic ABM to simulate burglary crime and capture its patterns.

1.1 Burglary And Criminology Theorems

Criminology theorems focus specifically on the spatial-temporal behavior of the individual involved in crime events and the details of the immediate surrounding physical environment. In this section we will explain such theories:

Theorem 1. *Routine activity theory* (Cohen, L.E. and Felson, M., 1979.) explores the interactions between victims, offenders and other people who might influence an individual crime event. The theory stipulates

that for a crime to occur an offender must meet a victim at a time and place with an absence of others who might prevent the crime. This convergence depends on the routine activities of the people involved.

Theorem 2. *Geometric theory of crime* (Bottoms, A.E. and Wiles, P., 1997) shares many similarities with routine activities theory, but focuses more explicitly on the interdependencies between a person's knowledge of the environment. The theory considers how the routes used to travel around a city influence a person's awareness space and hence the spatial-temporal locations in which offenders are likely to commit a crime. Burglars do not search for targets at random; instead they are likely to search near important places such as schools, work places.

Theorem 3. *Rational choice perspective* (Clarke, R.V. and Cornish, D.B., 1985) criminal's decision to commit crime is a cost-benefit analysis weighing up potential rewards of a successful crime with the risks of being apprehended. Thus a crime will only be committed if it is perceived as profitable.

Although they describe different elements of the crime system, the theorems largely agree on the mechanisms that lead to the spatial-temporal patterns of crime. A factor that is particularly relevant to crime modelling is that in each theorem the emphasis is on the individual-level nature of crime occurrences. The crime system is driven by the behavior and interactions of individual people situated in a highly detailed local environment which incorporate lower-level dynamics that ultimately explain why crime takes place.

1.2 AMB and Crime Prediction

ABM has been introduced by researcher to simulate crimes with the computational criminology technique to better understand of crime. In (Birks, D.J., Donkin, S. and Wellsmith, M., 2008), ABM used to test and validate criminological theorem, two models are used, one for general offending, and one for burglary, to arguing the potential of ABM method to simulate the crime dynamics. (Melo, A., Belchior, M. and Furtado, V., 2005) used an ABM to study strategies of preventive policing by simulating criminals and police patrols. (Hayslett-McCall, et. al. 2008) proposes an ABM combined with Cellular Automata based on social disorganization and *Routine activity theory*.

Researchers have attempted to simulate crime in many different ways, starting with simple models, aiming to represent the environment with real data (Winoto, P., 2002), who created a basic multi-agent model of crime, based on the *Rational choice perspective* with a goal to examine the equilibrium of crime and punishment. These models showed that dynamics of crime can be studied using simulation techniques. (Brantingham, et. al., 2009) present a framework for modeling crime, integrating theorems of crime analysis and prediction, by means of multi-agent modeling combined with state machines. This model includes a more realistic urban environment with a road network, which the agents can navigate, and a basic temporal component, for a hypothetical analysis of crime.

The most complete framework introduced by (Malleon, N., Heppenstall, A. and See, L., 2010) which attempts to simulate burglary in the city of Leeds by means of ABM. The agents can navigate the virtual environment and decide whether to burglar a property or not, depending on inner drives and believes, as well as on the characteristics of the property. In addition, real crime data has been used to manually calibrate the model. However, the proposed framework is most realistic in environment representation and offender architecture, this work suffer of non-realistic implementation of agents, and lack of automated calibration of the model. Therefore, automating the calibration process of the simulation parameters could lead to a self-learning model which would adapt to the new types of criminal behavior.

2 FARMWORK

2.1 Model Architecture

As the burglary crimes represents one of the wide majority of crimes committed in any region. In the proposed framework, our aim is to building an ABM to predict the burglary crime. Figure 1 show our

proposed framework architecture. In the first step of the process, historical real crime records, environment information, and neighborhood census data will be gathered. In second step, the simulation parameters are estimated from the gathered data and employed in ABM for burglary crime prediction. The prediction results from the running model will be compared to real burglary crime historical records for the same region of interest (ROI). finally, the generated simulation results from the system are compared to the gathered information to determine whether the initial simulation parameters should be adjusted, and additional data should be added to improve the predictions accuracy. This automated comparison is employed iteratively before more complexity added to it to verify the model prediction accuracy, and its logical consistent.

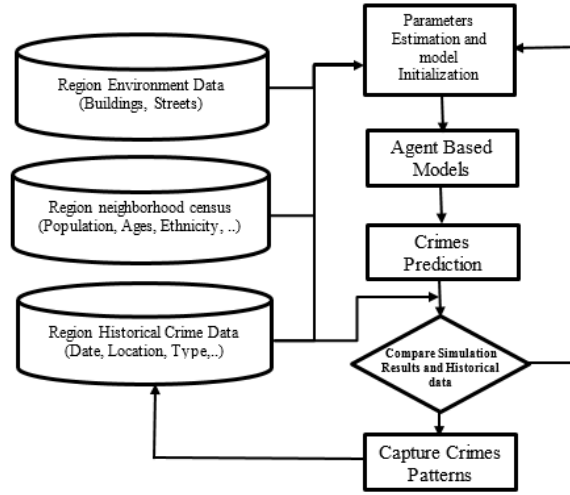


Figure 1 Proposed Framework Architecture

2.2 Data Gathering and Environment Reconstruction

Our model can be applied to any ROI defined by area zip code with available appropriate crime records, demography census information, and GIS information. The collected data is used as input to construct the burglar agents, and the realistic urban environment by estimating and initializing the model simulation parameters. The crimes information consist of all crimes recorded collected by Police Department (PD) for couple of years. The crimes information obtained from PD web services for the ROI, it includes the places of arrested drug dealers and provides the number of expected burglaries for burglary model output validation during simulation. Besides, it includes the following attributes for each crime: date, place, and type. All these data are processed to estimate burglary hotspots, living place, and extract geocoded information (such as latitude, longitude and place id) of each crime record.

The environment data, in which burglary crimes took place, plays important rule in model accuracy for crime simulation. Burglaries, victims, and navigation of the real environment have a great impact on the criminal chances. We build a web crawler to extract environment geocoded data from OpenStreetMaps (OSM) (<https://www.openstreetmap.org/>). The collected environment records contain road networks, public transportation, building data (such as building identity: commercial, residential, parking ...etc.), building address, and geocoded corners of building boundary. We also collected the social data and statistics (such as operating hours, population, and occupation density during time of day ... etc.) of the local businesses places which attract burglar during daily routine to socialize (such as restaurant, schools, parking ... etc.). A web crawler based on Google places APIs (<https://developers.google.com/places/>) is used to collect places statistics using OSM extracted information for buildings.

Finally, we collected neighborhood demographic census information of ROI to model burglaries in a data driven manner, and provide more representation of agents flow in the environment. Neighborhood demographic census statistics defined by attributes such as total population, density per square mile, age,

genders, ethnicity, marital status, occupation, family size, annual household income, highest degree attained, Housing Inventory, and property occupation to provide synthetic list of entire population of ROI.

The importance of synthetic population for constructing a realistic urban environment is to represent individual community (transportations, and houses) with the absence of household population data. The individual community used to estimate occupancy and attractiveness for every individual household in the simulation of the ROI rather than assuming all households are similar. Also, it impact the global vision and sense of the community, which defined as neighborhood cohesion. Table 1, shows community factors and their descriptions.

Table 1 Synthetic population community factors and their description

Community factors	Description
Attractiveness	Weight of abundance of item within houses in the ROI based on annual household income.
Occupancy	Estimate whether household is occupied at specific time based on employment data of people who live in ROI. i.e. Family houses are more occupied during the day and evening time.
Neighborhood cohesion	A measure of how well interrelation among the community households. The community defined by neighborhood census data, such as ethnicity, age, and language.

2.3 Agent Architecture

The main component of our framework is a self-governing action entity called burglar agent, which is defined to model the burglary behavior in a realistic urban environment. The burglar agent designed to perform day to day behaviors with aim to capture the patterns of physical or social environment behaviors in order to construct accurate model of burglary crimes. These behaviors characterizes the motive of the burglar and include the needs to generate wealth, sleep, and socialize. Burglary is commonly linked to a drug addiction [reference] so addiction must form a part of the burglar's agent model behaviors. Table 2, shows the burglar behaviors and their descriptions.

Table 2 Burglar agent behaviors and their descriptions

Behavior	Description
Wealth	Measure used to control burglar agents behaviors that require money to perform the behavior such as socialize, or drug addiction.
Sleep	Measure of how much burglar agent needs to sleep. Sleeping behavior helps to enforce realistic temporal behavior on the agent, which is increase over night time and force burglar agent to travel to home, if it required.
Socialize	Measure of amount of socializing performed by burglar agent. Burglar can spend time socializing in specific places such as bars, schools, or friends' houses. This behavior decreases overtime and costs money.
Addiction	Weight of level of substance in the burglar agent which is increases after drug use and decreases overtime. Agent can purchase drugs from drug dealers, if required.

The burglar agents is structured into several modules as shown in Figure 2. The burglar agent model decomposed into logical component following the urban environment complements extant theories of the environment in the behavioral science (Akers, R.L., 2011).

- **The Space Evaluation Module** plays an important rule on the space activity and space perception of the burglar agent. It uses a navigation algorithm to move the agent from an origin to a destination considering the agents preferences such as chose the shorter routes, avoid traffic, use familiar roads ... etc.
- **The Target Selection Module** is responsible for monitoring all possible targets on the routes taken by the burglar agent, and select the attractive targets based on selection criteria. This leads to creation burglary space of the agent.

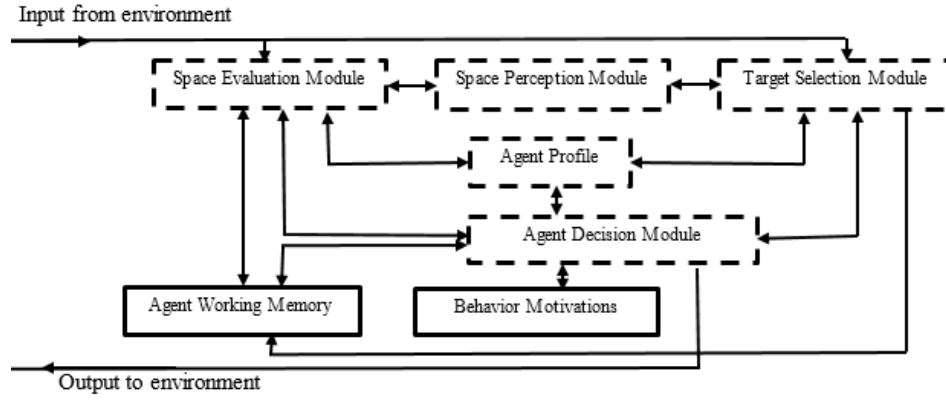


Figure 2 Burglar Agent Architecture

- **he Agent Decision Module** is core central module which model intelligent decisions made by the burglar agent. Most of these decisions are made based on Agent profile and agent's behavior motivations. This module is using naïve artificial intelligent (AI) methods to perform burglar's agent decisions.
- **The Agent Working Memory** is a collection of facts/beliefs/knowledge which change dynamically over time. In our agent architecture, a central memory module used to store temporal information needed for routes and target selection.
- **The Space Perception Module** is special kind of memory stores environment view such as the places that the agent has been often visited. The Space Evaluation Module and The Target Selection Module are continuously update its records.
- **The Agent Profile component** represents the burglar's agent attributes and stores information such as preferences, skills, home geocoded and demography factors. The profile information is taken into account for determining the burglar agent behavior when move in space and time.

Table 3 behavior state variables and their intensity functions

Behavior State Variable	intensity functions	Variables
Sleep	$m_{sleep} = \left(\frac{1}{s_{sleep}}\right) * f(t, offset_{sleep}) + \left(\frac{1}{s_{sleep}}\right) * b_{sleep}$	$f(t, offset_{sleep})$: sleeping time of day function. t : day time $offset_{sleep}$: sleeping time offset constant
Socialize	$m_{social} = \left(\frac{1}{s_{social}}\right) * g(t, offset_{social}) + \left(\frac{1}{s_{social}}\right) * b_{social}$	$g(t, offset_{social})$: socialize time of day function. t : day time $offset_{social}$: social time offset constant
Drug addiction	$m_{drug} = \left(\frac{b_{drug}}{s_{drug}}\right)$	
Wealth	$m_{wealth} = \left(\frac{1}{s_{wealth}}\right) * w(t, offset_{wealth}) + \left(\frac{1}{s_{wealth}}\right) * b_{wealth}$	$w(t, offset_{wealth})$: wealth time of day function. t : day time $offset_{wealth}$: wealth time offset constant

Behavior Motivations component is used to model the intensity of motives using intensity function based on behavior state variables which are represented by most common needs for burglary. Therefore, the state variables which provide sufficient variety to create realistic daily behaviors are wealth, socialize, sleep, and drugs addiction. These state variables represent an improvement over existing agent based models of crime in terms of encapsulation agent behavior. The burglar agent with low drugs, social, or sleep level needs to behave in such a way as to increase the value of the state variable by taking drugs, socializing, or sleeping.

In general the intensity of motive, m is inversely proportional to the size of its state variable, s . Since the population of burglar agents does not need to be homogeneous, which is one of benefits of agent based modeling, different agents can be affected by state variables and motives differently. This feature is incorporated into the model by including a burglar parameter, b , that affects motive intensities such that: $m \propto b \frac{1}{s}$. Therefore, the burglar agent model with a large value of b is more affected by a specific motive than a burglar agent model with a low value of b even if both agents have the same state variable level. The rate that state variable decrease can be used to configure the amount of time that burglar agent should spend to satisfy the motive. Table 3, shows behavior state variables and their intensity functions.

Figure 3, shows how the time of day effects the overall intensity of the wealth, socialize, and sleep motives. The motive to generate wealth is strong during the day time, while the needs for sleep is the strongest during the night time. Also, the burglar agent is more motivated to socialize during the evening than the day time.

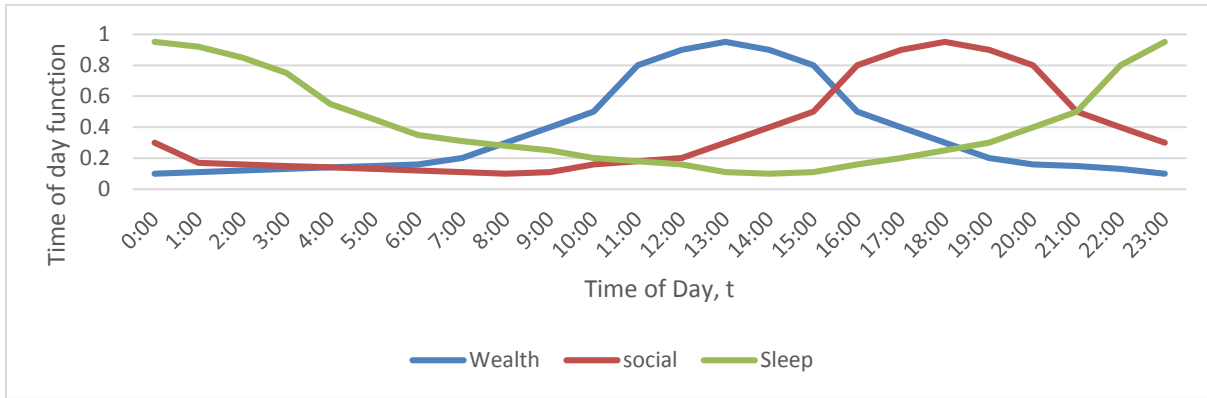


Figure 3 how the behavior intensity motive varies with time of day, assuming a constant value for the state variable.

Figure 4, illustrate how the behavior motive of burglar agent drives their behavior actions. The intensity function which utilize the burglar agent parameter, time of day function, and behavior state variable to calculate which behavior motive has the greatest chance at each time. For each agent, the intensity function take into account the current levels of wealth, sleep, socialize, and drug and the current time of the day. The social interactions among the burglar agents are deemed too complex to be of use in the burglar agents' models. Each motive has an association goal which the agent can accomplish to increase the value of the appropriate state variable and subsequently lower the strength of the motive. For example, accomplish the goal of drug addiction will increase the drug state variable, reduce the size of the motive and cause another motive to start controlling the burglar agent's behavior. To satisfy agent's behavior goal, there are often sub-goals that must be accomplish first based on burglar agent's preferences. For example, wealth is required for socialize, and drug buy behaviors.

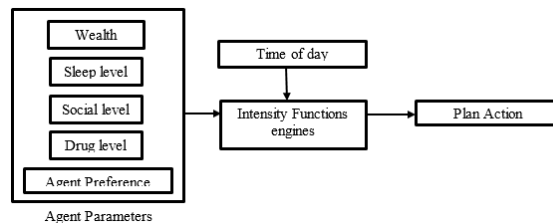


Figure 4 Burglar agent's behavior module.

Burglar agent does not has global knowledge of their environment, for that it build up its own awareness of the environment based on daily routine activity. This introduces the importance of agent characteristic which is the agent ability to build its own cognitive map. This characteristic brings the burglar agent much more with the crime theory and means that the urban form of an area will have an influence on burglary crime patterns. For example, the process of going to buy drugs with large number of drug dealers in specific

geographical area has important influences on burglar agent cognitive map by making that geographical area more vulnerable to burglary crime. Agent working memory module is used to store burglar agent cognitive map which is built by itself during daily routines activity.

2.4 The Realistic Urban Environment

The realistic urban environment represents the space that burglar agents inhabit. In ABM the environment is spatial and has two major characteristic: it must allow the burglar agents to travel from one place to another using available transportation networks, and it must include the important factors that form the environment such as neighborhood data. The environment model has been designed to be realistic as it necessary for a burglary simulation. For that, there are three separate layers make up the environment:

- **The building physical layer:** contains the physical form of buildings design and nature features (such as the geographic coordinates of the buildings) which are important for burglar agents to attempt the crime for specific ROI. This layer was generated from OSM geographic and neighborhood demography census data, where each building is a unique object with different physical attributes that are important for understanding the burglary crime patterns during the simulation of the model. The building physical layer defined as an undirected graph as following: Let $G_B = (V_B, E_B)$ be undirected graph represent the physical buildings objects of ROI, where $V_B = \{1, \dots, n\}$ is set of vertices representing the buildings on the ROI, and $E_B \subseteq V_B \times V_B$ is set of undirected edges representing the nearest neighbors network. The vertices specifying the physical attributes of individual buildings that might increase or decrease their burglary risk, while edges specify the distance between the buildings in the ROI. Table 4, lists the building attributes that have been chosen to express the building burglary risk.

Table 4 Building attributes and their description

Building attributes	Description
Accessibility	Measure the difficulty of the entering the building. In the model highly accessible buildings have a higher burglary risk.
Visibility	Measure of how visible the building is to neighbors and guardians. While the distance between the buildings play important rule in visibility variable, in this model the higher visibility reduces the burglary risk.
Traffic volume	Measure of the traffic volume outside the building. In the model higher traffic levels make it difficult to gain access to the building and reduce the burglary risk.
Security	Measure of the physical security of the building. In the model higher security reduces the burglary risk.

- **The transportation physical layer:** a physical layer that represents the transportation network and its physical features (i.e. highway, railway, cardinal directions) for the simulation ROI. This layer shapes the burglar agents awareness of the environment. Therefore, the virtual environment has a realistic transportation network consisting of the roads that can be used by the burglar agents to be driven or walk. The realistic routing behavior is obtained by the speed that burglar agents used, for example agents with cars are driving on the major roads rather than using minor roads. Each road is unique and its physical features used to define the traffic volume for the building physical layer. This layer has generated from OSM geographic and neighborhood demography census data. The transportation physical layer defined as a directed graph as following: Let $G_t = (V_t, E_t)$ be directed graph represent the road network of ROI, where $V_t = \{1, \dots, n\}$ is set of vertices representing the road intersections, and $E_t \subseteq V_t \times V_t$ is set of directed edges representing the connectivity in the road network. The vertices specify the physical attributes of individual road that might be used by burglar agents to travel during the daily routine activities. While edges specify the distance between the road network intersections and direction of traffics in the ROI.

- **The neighborhood layer:** a social demography based layer used to model individual behavior of social community. For example, high levels of neighborhood cohesion have been linked to level of burglary risk because local people are more likely to prevent a crime occurrence. The importance of neighborhood layer is to determine where burglar agent starts their search, and helps to capture the social patterns that surround a burglary crime. This layer was generated from OSM geographic data, historical crime data, and neighborhood demography census data. Table 5, lists neighborhood attributes that have been chosen to define the neighborhood patterns of burglary.

Table 5 Neighborhood attributes and their descriptions

Neighborhood attributes	Description
Consistent	Measure of cohesive in neighborhood. In this model the high neighborhood cohesion reduces the burglary risk.
Attractiveness	Measure the potential rewards within the building. In this model the higher attractiveness increases the burglary risk.
Occupancy	Probability of building being occupied at given time. In this model higher occupancy decreases the burglary risk.
Familiarity	Measure of how similar the neighborhood to the burglar living area. In this model, the higher similarity increases the burglary risk.

The neighborhood layer also defined as an undirected graph as following: Let $G_N = (V_N, E_N)$ be undirected graph represent the neighborhood social demography of ROI, where $V_N = \{1, \dots, n\}$ is set of vertices representing the building household, and $E_N \subseteq V_N \times V_N$ is set of undirected edges representing the households connectivity. The vertices specify the individual household that might be used by burglar agents to commit burgle, while edges specify the likelihood of neighborhood in the ROI.

2.5 Burglary Process and Decision Scenarios

Social systems are incredibly complex due to the large number of interacting elements and many underlying processes that are simply not understood. Moreover, these processes are generally non-linear such that small changes in system parameters can have large effects on the outcomes of the system as a whole. We designed our burglar Agent to be independent, and has a capability of interacting with any surrounded environment. The burglar agent makes assessment of its situation overtime and then make decisions.

Burglar agent start the simulation at home with low motivation levels to satisfy its needs which mean it has no motives to perform an action such as sleep, take drugs or socialize . Overtime, its motivation levels increase and becomes motivated to perform an action. The simulation is configured so that the agent during the day must sleep for H_{sleep} hours, socialize H_{social} hours, and buy drugs N_{drug} times. To purchase drugs or to socialize requires generate wealth which must be sought through burglary. At present, the wealth gained from a single burglary is sufficient to allow the agent to buy drugs and socialize. Therefore, an agent will burgle once per day. However, the probability of committing a burglary depends on the suitability of the houses that the agent passes and how highly motivated the agent is at the time, so agents will have days when no burglary takes place. Again overtime the agent become more desperate as motives increase, hence there will be other days when multiple burglaries take place. In this sense the agents are autonomous and perform different activities depends on its own behavior without any kind of restrictions.

Finding a place to burgle is the most complex of burglar agent's actions, and most important generate burglary patterns. Using modular approach enables different types of burglar to simulate simply by replacing any action with another. Decision scenarios are designed based on criminology theorems which prove that burglars are unlikely travel far, and not usually burgle too close to home location for fear of being recognized. The burglary process can be break down into three decision scenarios as following:

- **Decide where to start looking for victims:** The decision of where to start the search performed by the agent incorporates Geometric Theory of Crime, and Rational Choice Perspective. Where any burglar used to analyze the potential rewards against the risk during the process of choose target areas. The burglar agent decide where to start the search based on the awareness of the area and assign a region likelihood weight (RLW) to every area in its own cognitive map relative to its home location, and its current location. The burglar agent incorporate the Neighborhood layer attributes to make the surrounded cues more close to the realistic burglary environment. Also, the burglar agent calculates the RLW for every area in its cognitive map and then choose the area with higher RLW probability.

$$RLW = \left(\frac{1}{dist} \right) + \text{Attractiveness} + \text{Familiarity} ,$$

Where *dist* represents the distance of target from the burglar's agent current position, Attractiveness, and Familiarity are Neighborhood layer attributes.

- **Search for victims:** based on Geometric Theory of Crime, the burglar does not search randomly for burglary target, but follow known search patterns. Once the burglar agent reaches a suitable area, as defined in last decision scenario, and begins a search the amount of time spent on the search depends on the profile of the burglar. The search time variable T_{search} can be varied in the burglar agent model, and used to control the burglary process. Once T_{search} is reached the burglar agent choose new start location, then travel to there and search for victims. For example, Burglar agent with small search time spends short time in search process and often commit the burglary in the first empty building.
- **Choose a suitable victim:** burglar choose individual victims based on their individual characteristics, so it cannot normally be assumed that a neighborhood is homogeneous with respect to burglary risk. Once the burglar agent decides to commit burglary, it starts to examine all buildings to determine ability for burglary. This process take place during the travel time to the target place as well as during the search time. The burglar agent decide where to commit the burglary based on the buildings physical layer attributes to define the burgle ability weight BAW to commit the burglary.

$$BAW = \text{Accessibility} + \text{Visibility} + \text{Traffic volume} + \text{Security} ,$$

Therefore, If BAW of the building is higher than the intensity of the motive, then burglar agent will not attempt the burglary, and start looking for another target.

3 EXPERIMENTS AND RESULTS

The proposed model was programed and executed using Repast py, a free open source toolkit widely used for ABM crime (<http://repast.sourceforge.net/>). The following section presents analysis and results of model case study.

3.1 Case study: region 23508 analysis and justification

As the proposed framework is generic and can be applied to any ROI for which appropriate data is available. In this model we preformed data analysis to a region of approximately 2.0 miles square and populated by 23,044 individuals in the city of Norfolk, VA, USA its boundary defined by area zip code 23508 as shown in figure 5. The region 23508 contains some of the most ranging from urban core to urban pioneer neighborhoods in city of Norfolk. We decompose this ROI into five district based the demography neighborhood census data, as shown in figure 5.



Figure 5 Left side image show ROI 23508 boundaries, while right side image shows districts' boundaries of the ROI.

Table 6 ROI 23508 statistics per district

District	Total Population	Area (Square/Mile)	density / square mile
Lamberts Point	3,847	0.449	8,173
Larchmont-Edgewater	5,962	1.464	3,463
Highland Park	2,163	0.24	10,075
Colonial Place/Riverview	4,185	0.584	6,681
Park Place	6,887	0.518	9,610

In general, by analyzing the neighborhood demography census data for ROI 23508, we found that approximately 28.11% of residents are under the age of 20. College-aged residents are the second largest age group in with approximately 27.66% of residents between 20 and 24 years old. Approximately 57.18% of people are Caucasian, 33.70% of 23508's residents are African American, making it the second most common ethnicity, and a smaller percentage of Hispanic residents. Men greatly out-populate women in 23508 by 12.0%. Most households in 23508 primarily speak English 89.5% of residents, Spanish is the second most commonly spoken language with 3.5% of households using it as the primary language. Most families (69.2%) are classified as small, meaning they have two or three members. Most residents over the age of 15 have never married, while The second most common marital status in 23508's is married, approximately 26.8% of the population is currently in a marriage. The percentage of housing units rented is higher than the percent owned by 3%. 40.1% of rental units cost between \$500 and \$1,000 every month in 23508. The next most common monthly rent range is between \$1,000 and \$1,500. Approximately 33.2% of 23508's renters pay between \$1,000 and \$1,500.

To make our model more realistic the analysis of the neighborhood census data is applied also to all districts in the ROI as following to provide a deeper understanding of the environment and find more effective burglary crimes patterns that contribute to the burglar agent knowledge, also to find out the high impact environment attributes that shape the burglary crime environment. In general, Lamberts Point's population is characterized as low income, primarily composed of the children & teenagers and college-aged age groups, and being less educated. Larchmont-Edgewater is upper middle class, well educated, and primarily composed of college-aged, which is the most prevalent age group. Highland Park is low income, well educated, and primarily composed of college-aged, which is the most prevalent age group. Colonial Place-Riverview is upper middle class, well educated, and primarily composed of young professionals, which is the most prevalent age group. Park Place is a low income, less educated neighborhood.

Also, we analysis the real historical crime data for ROI 23508 to capture the hypothesis of burglary crime. The data obtained from Norfolk Police Department (NPD) for the period of Jan 1, 2009 to

Dec 31, 2016 were obtained from (<https://www.neighborhoodscout.com/>, and <http://www.city-data.com/>) consist of approximately 2000 burglary crime records.

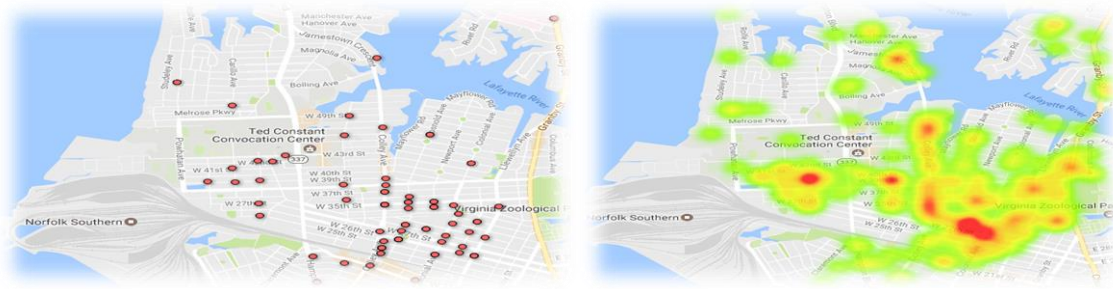


Figure 6 Left image show places where drug dealers are arrested, right side image show burglary crime hotspots.

Figure 6 shows places where drug dealers are arrested, and the burglary crimes in ROI represented by a heat map where the red areas has higher burglary rate, while green areas represent lower rate of crime. The visualization of crime data helped us to rank 23508's districts based on crime as the higher crime rates are found in Park Place district, Lamberts Point district, Highland Park, Colonial Place/Riverview, and the lower crime rates are found in Larchmont-Edgewater area. Figure 7 shows that weekend has higher crime rate more than normal week days, which reflect the effect of occupancy attributes on the burglary crimes, weekend has higher burglary risk rate.

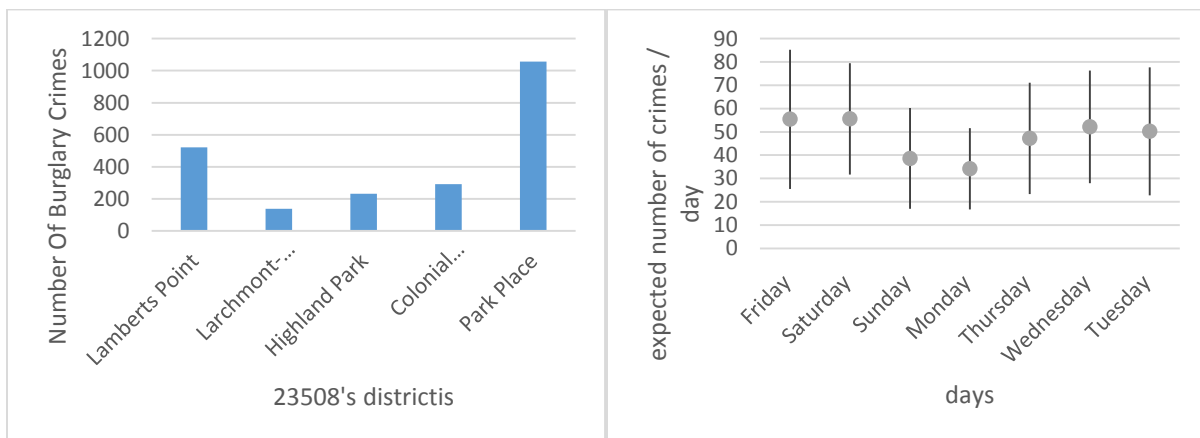


Figure 7 left side chart shows number of crime per district, right side chart shows number of crimes per day for 23508

By mapping the neighborhood demography census data analysis and historical crime data analysis we found that, low income, and less educated neighborhood has higher burglary crime risk, while upper middle class, well-educated neighborhood has lower burglary crime risk. Most of burglaries prefer to live in low income areas. Burglary crimes are linked to drug addiction. Drug dealers' locations always surrounded by high risk burglary crimes, crimes hotspots. Ethnicity has no impact on burglary crimes. Domestic local areas near to highways have higher risk of burglary crimes.

We run our model on ROI 23508 information to regenerate the burglary crime hotspots map, where first four years of historical burglary crime records are used to calibrate our model using 100 iterations, ten burglar agents are conducted, and the simulation run for a year, simulation time. While in this experiment the total burglary crimes are not predicted and there are still differences between historical information and the generated crimes by proposed framework, a match of general pattern can be recognized, figure 8. The match is considered to be sufficient for the propose of the proposed framework. Based on simulation and real burglary records, we found that drug addiction is the key motivation of burglars to engage in burglary. Regions with high income and cohesive community chosen are not usually considered by burglars during

target selection. Regions with high income mean high security which prevent burglars from approaching them.

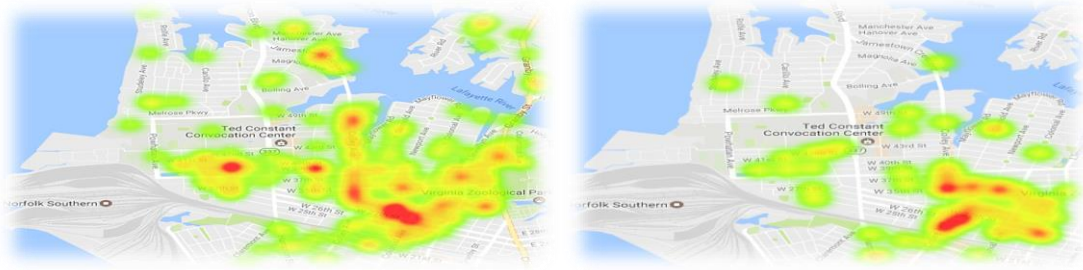


Figure 8 left side image show the real historical burglary crime hotspots map, left side image show the regenerated results using proposed framework.

4 CONCLSION

The proposed framework has shown enormous potential for simulating burglary criminal activity. This model enable burglary crime to be studied from an individual instead of aggregate perspective. This paper proposed an agent-based model and simulator whose main contribution is the ability to formulate burglary criminal behavior in the context of realistic environment strategies.

The proposed framework can be used to test the effectiveness of crime reduction initiatives. This is an important feature of the model since it would enable law enforcement authorities to test strategies before applying them in real situations and observe the effects of crime prevention measures such as displacement of hotspots. Finally, this study should open the door to the simulation of other types of crimes.

REFERENCES

- Cohen, L.E. and Felson, M., 1979. Social change and crime rate trends: A routine activity approach. *American sociological review*, pp.588-608
- Bottoms, A.E. and Wiles, P., 1997. Environmental criminology (pp. 620-656). *The Oxford handbook of criminology*.
- Clarke, R.V. and Cornish, D.B., 1985. Modeling offenders' decisions: A framework for research and policy. *Crime and justice*, 6, pp.147-185
- Birks, D.J., Donkin, S. and Wellsmith, M., 2008. Synthesis over analysis: Towards an ontology for volume crime simulation. In *Artificial crime analysis systems: Using computer simulations and geographic information systems* (pp. 160-192). IGI Global
- Akers, R.L., 2011. Social learning and social structure: A general theory of crime and deviance. Transaction Publishers.
- Melo, A., Belchior, M. and Furtado, V., 2005, July. Analyzing police patrol routes by simulating the physical reorganization of agents. In *International Workshop on Multi-Agent Systems and Agent-Based Simulation* (pp. 99-114). Springer Berlin Heidelberg
- Hayslett-McCall, K., Qui, F., Curtin, K.M., Chastain, B., Schubert, J. and Carver, V., 2008. The simulation of the journey to residential burglary. *Artificial crime analysis systems*, pp.281-300
- Winoto, P., 2002, July. A simulation of the market for offenses in multiagent systems: is zero crime rates attainable?. In *International Workshop on Multi-Agent Systems and Agent-Based Simulation* (pp. 181-193). Springer Berlin Heidelberg
- Brantingham, P., Glässer, U., Jackson, P. and Vajihollahi, M., 2009. Modeling criminal activity in urban landscapes. In *Mathematical methods in counterterrorism* (pp. 9-31). Springer Vienna.
- Malleson, N., Heppenstall, A. and See, L., 2010. Crime reduction through simulation: An agent-based model of burglary. *Computers, environment and urban systems*, 34(3), pp.236-250.

AN AGENT-BASED SIMULATION OF THE IMPACT OF YOGIC BREATHING ADOPTION IN HAMPTON ROADS

Saturnina Nisperos, Sonali Kakde, Frederic McKenzie
Department of Modeling, Simulation & Visualization Engineering
Old Dominion University
Norfolk, VA
snisp001@odu.edu, skakd001@odu.edu, rdmckenz@odu.edu

EXTENDED ABSTRACT

Rising stress levels is a significant problem in the US. People who do not have tools to manage their stress are more likely to acquire cardiovascular disease (CVD) – the leading cause of death in the US. Greater Hampton Roads data indicates that an average of 23% of Medicare beneficiaries in the region were treated for heart disease. For affected individuals, it not only implies health risks but also increased medical expenditures which can lead to financial issues – the top stressor in the US. Research results reveal that yoga reduces the stress perceived and modulates stress response systems. An agent-based model was developed in this study to simulate the micro and macro-level impact of adopting a yogic breathing technique in Hampton Roads, particularly on managing stress, preventing CVD and reducing the associated medical expenses.

Keywords: agent-based model, cardiovascular disease, yogic breathing technique

1 INTRODUCTION

Rising stress levels is a significant problem in the US. When stressed, some people resort to harmful practices like smoking, drugs, alcohol abuse or overeating. Over a period of time, these factors can lead to high risk of cardiovascular disease (CVD) – the leading cause of death in the US. For affected individuals, it not only implies health risks but also increased medical expenditures which can lead to financial issues – the top stressor in the US.

Research studies recommend breathing techniques as a way to deal with stress. Harvard Health Publications (2009) indicates that yoga reduces the stress perceived and modulates stress response systems. Sudarshan Kriya Yoga (SKY), a form of yogic breathing, is unique method for balancing the autonomic nervous system and influencing psychologic and stress related disorders (Brown & Gerbarg, 2005). This study aimed to develop an agent-based model to simulate the adoption of SKY in the Hampton Roads area. Its effect on the stress level, body mass index (BMI), systolic blood pressure (SBP) and medical expenses was examined.

2 METHODS

The model was built using Netlogo, adapting the Spread of Disease model by (Rand & Wilensky, 2008). An agent represents a person and the properties assigned represent the associated CVD risk factors. Stress is modeled as the disease which affects the CVD risk of the agents and SKY is introduced as intervention which stimulates reduction of stress level. The overall flow of the model is divided in two major parts – the initialization and the simulation. The initialization sets the CVD risk matrix, agents' properties, and the interconnection among agents. To determine the agents' corresponding CVD risk index, their properties (age, sex, blood pressure, smoking status, and manifestation of diabetes mellitus) are mapped to the CVD risk matrix which was adopted from the WHO/ISH risk prediction chart (World Health Organization, 2007). Moreover, the agents are either classified as practicing SKY or not, a CVD or non-CVD patient and the

network denotes the relationship that a person may have with others (e.g. family, work or friend). The model does not account for birth and death of agents and the concept of hereditary CVD. Table 1 summarizes the associated activity and rules that occur in every simulation iteration.

Table 1. Simulation rules

Activity	Rules
Interact with other agents	An agent may hassle or uplift linked agents thereby increasing or decreasing their stress level. An agent with higher stress level has greater chance of causing hassle than uplifting others. (DeLongis & Folkman, 1988)
Perform SKY	SKY decreases BMI, SBP and stress level (Narnolia, et al., 2014)
Take hypertension medication	Agents with hypertension take medication. Hypertension medication decreases BMI, SBP (American Heart Association, 2016; Narnolia, et al., 2014)
Update agent properties	The BMI and SBP of an agent are increased when under stress (Lucini, Fede, Parati, & Pagani, 2005; Harding, et al., 2014). BMI identifies the diabetes risk of an agent (National Institute of Diabetes and Digestive and Kidney Diseases, n.d.). The higher the CVD risk, the greater the chances that an agent may incur CVD (World Health Organization, 2007).
Spend on CVD medical expenses	An agent who incur CVD spends amount on medical expenses (Moore, Levit, & Elixhauser, 2014). Increased medical expenditures can lead to financial issues which may affect stress level (American Psychological Association, 2016).

To verify and validate the model, its result, specifically the BMI and SBP parameters were compared to the research results of Narnolia, et al. (2014), Agte, Jahagirdar, & Tarwadi (2011) and Wolff, Sundquist, Lönn, & Midlöv (2013). The unpaired t-test results indicate that the difference between the research studies and simulation results are not statistically significant at 0.05 level.

3 RESULTS

The population profile of Hampton Roads was used as input parameter value of the model to simulate the micro and macro-level impact of SKY in addressing CVD. Data were gathered from government agency websites that provide statistics about Hampton Roads. The sensitivity of the associated CVD output parameters (high risk population, patient population, expenses, BMI, SBP and stress level) were examined by varying the percent population of those who practice SKY.

The simulation results show that the mean risk factors value (SBP, BMI and stress level) of the SKY population is lower in general compared to the non-SKY population. The SBP of the SKY group are mostly normal while the non-SKY group are under hypertension category. Furthermore, the BMI of those who practice SKY generally falls under normal weight while the non-SKY falls under overweight category. As the percentage population of those who practice SKY increases, the percentage of CVD high risk and patient population and the total CVD medical expenses in Hampton Roads decreases. This signifies that the risk of incurring a CVD by the people who practice SKY is reduced, thereby decreasing the overall population of CVD high risk, patients and total medical expenses in the region. Moreover, the analysis of variance of the different percentages of SKY population indicates that the increase in SKY population in Hampton Roads provides an extremely significant effect ($p\text{-value} < 0.001$) in decreasing the percentage of CVD high risk population, associated expenses and the overall stress level in the region.

4 CONCLUSION

The not statistically significant difference between the results of published studies on SKY and the simulation results, supports the validity of the model. Simulation results show that people who practice SKY have generally lower and normal SBP, BMI and stress level while those who don't, have generally higher and falls on the high-risk category. Lastly, the analysis of variance result suggests that SKY can be considered as a feasible tool to manage stress and an intervention to the increasing and alarming rate of CVD.

REFERENCES

- Agte, V., Jahagirdar, M., & Tarwadi, K. (2011). The effects of Sudarshan kriya Yoga on some physiological and biochemical parameters in mild hypertensive patients. *Indian Journal Physiol Pharmacol*, 183-187.
- American Heart Association. (2016, October). *Understanding blood pressure readings*. Retrieved from http://www.heart.org/HEARTORG/Conditions/HighBloodPressure/AboutHighBloodPressure/Understanding-Blood-Pressure-Readings_UCM_301764_Article.jsp#.WJYklvkrLIV
- American Psychological Association. (2016). *Stress in America: The impact of discrimination*. Stress in America™ Survey.
- Brown, R., & Gerbarg, P. (2005). Sudarshan Kriya Yogic breathing in the treatment of stress, anxiety, and depression: part II - clinical Applications and guidelines. *The Journal of Alternative and Complementary Medicine*, 711-717.
- DeLongis, A., & Folkman, S. (1988). The impact of daily stress on health and mood: psychological and social resources as mediators. *Journal of Personality and Social Psychology*, 486-495.
- Harding, J. L., Backholer, K., Williams, E. D., Peeters, A., Cameron, A. J., & Hare, M. J. (2014). Psychosocial stress Is positively associated with body mass index gain over 5 years: evidence from the longitudinal AusDiab study. *Obesity*, 277-286.
- Harvard Health Publications. (2009, April). *Yoga for anxiety and depression*. Retrieved from <http://www.health.harvard.edu/mind-and-mood/yoga-for-anxiety-and-depression>
- Lucini, D., Fede, G. D., Parati, G., & Pagani, M. (2005). Impact of chronic psychosocial stress on autonomic cardiovascular regulation in otherwise healthy subjects. *Hypertension*, 1201-1206.
- Moore, B., Levit, K., & Elixhauser, A. (2014, October). *Costs for hospital stays in the United States*. Retrieved from Agency for Healthcare Research and Quality: <http://www.hcup-us.ahrq.gov/reports/statbriefs/sb181-Hospital-Costs-United-States-2012.pdf>
- Narnolia, P., Binawara, B., Kapoor, A., Mehra, M., Gupta, M., Tilwani, K., & Maharia, S. (2014). Effect of Sudarshan Kriya Yoga on cardiovascular parameters and comorbid anxiety in patients of hypertension. *Scholars Journal of Applied Medical Sciences*, 3307-3314.
- National Institute of Diabetes and Digestive and Kidney Diseases. (n.d.). *Risk factors for Type 2 diabetes*. Retrieved from <https://www.niddk.nih.gov/health-information/health-communication-programs/ndep/am-i-at-risk/diabetes-risk-factors/Pages/diabetesriskfactors.aspx>
- Wolff, M., Sundquist, K., Lönn, S. L., & Midlöv, P. (2013). Impact of yoga on blood pressure and quality of life in patients with hypertension – a controlled trial in primary care, matched for systolic blood pressure. *BMC Cardiovascular Disorders*.
- World Health Organization. (2007). *WHO/ISH Risk prediction charts*.

MEDICAL SIMULATION

- Page 95 Keyi Xue, Debrup Banerjee and Jiang Li
Panther Creek High School and Old Dominion University
Speech Feature Investigation in Transfer Learning for Improved Posttraumatic Stress Disorder Diagnosis
- Page 105 Emily M. Harthley and Matthew C. Hoch
Old Dominion University
Decision Analysis Tree to Determine Injury Prevention Strategy for Ankle Sprain Injury Based on Cost
- Page 111 jing Xu and Andrey N. Chernikov
Old Dominion University
Homeomorphic Tetrahedral Tessellation For Biomedical Images
- Page 120 Md. Shariful Islam and Michel Audette
Old Dominion University
Musculoskeletal Simulation for Geriatric Applications Based on Opensim/simtk
- Page 122 Shrabani Ghosh and Michel Audette
Old Dominion University
Towards Deformable Cranium & Foramen Surface Model
- Page 124 Lucas N. Potter
Old Dominion University
Scoliosis Surgery Planning Through Cadaveric Ligamento-Skeletal Tissue Mapping and Loading Studies, Multi-surface Segmentation, and Finite Element Simulation of the Spine
- Page 126 Austin Tapp and Michel Audette
College of William and Mary and Old Dominion University
Practical Applications of Neurosurgical Ontologies for Various Craniotomic Approaches through Computer Assisted Surgery

SPEECH FEATURE INVESTIGATION IN TRANSFER LEARNING FOR IMPROVED POSTTRAUMATIC STRESS DISORDER DIAGNOSIS

Keyi Xue
Panther Creek High School
6770 McCrimmon Pkwy,
Cary, NC, USA
xkeyi@students.wcpss.net

Debrup Banerjee
Department of Electrical and
Computer Engineering
Old Dominion University
Norfolk, VA, USA
dbane001@odu.edu

Jiang Li, Ph.D.
Department of Electrical and
Computer Engineering
Old Dominion University
Norfolk, VA, USA
JLi@odu.edu

ABSTRACT

In this paper, we investigated speech features for posttraumatic stress disorder (PTSD) diagnosis with deep transfer learning. Three categories of speech features including prosodic features, vocal-tract features, and excitation features were extracted from PTSD patients' speech recordings for diagnosis. In transfer learning, we first trained a deep belief network (DBN) on TIMIT, a large speech recognition data set co-developed by Texas Instrument and Massachusetts Institute of Technology, for phoneme recognition. We then utilized a trained DBN as a feature extractor to transform the PTSD speech features to new representations for improved diagnosis. The goal of this study is to identify which of the speech feature categories are effective in diagnosing PTSD. We conducted hypothesis testing to evaluate if the contribution from a feature category is significant. Our experimental results show that prosodic and vocal-tract features are effective in diagnosing PTSD while contribution from excitation features is negligible.

Keywords: PTSD, speech analysis, pattern recognition, Auto-encoder, Transfer learning

1. INTRODUCTION

Posttraumatic Stress Disorder (PTSD) is a mental disorder that can develop after a person is exposed to a traumatic event, such as warfare, traffic collisions, or other threats on a person's life [1]. Diagnosis of PTSD is mostly based on patient-self reporting during a structured clinical interview. The main challenge in PTSD diagnosis is the disinclination of patients to come to clinics for diagnosis because of the stigma associated with the disease. An alternative solution is to analyze patients' speech patterns to diagnose PTSD.

Our previous study showed that prosodic, vocal-tract and excitation features could be fused for PTSD diagnosis in the context of deep transfer learning [6]. In transfer learning, we transfer knowledge from speech recognition domain to solve problems in PTSD detection application. Transfer learning was usually utilized if both the source domain and target domain are similar but target domain is lack of data for complex modeling. PTSD data is difficult to collect but speech recognition data sets are widely available. We showed that the transfer learning can boost the performance of PTSD diagnosis about 13%

using the fusion of the three feature categories [6]. In this paper, we investigated which of the three categories of features or which feature combination is the most effective feature set for detecting PTSD.

Prosodic features include power, loudness, and intonation etc. of speech signal. Vocal-tract features consist of the MFCCs (human ears' simulation) and Teager energy operator (air flow change detector). Excitation features are composed of Shimmer and Jitter in speech segments. We conducted the feature selection task in the framework of deep transfer learning to identify which of the three feature categories is the most effective feature set in PTSD diagnosis. In deep learning, we used all the three categories of features as inputs to a deep learning model that fused the features for PTSD diagnosis. To test if excluding one of the three categories of features will improve the diagnosis performance, we kept the same deep learning structure but zeroed all features belong to the category as inputs. The deep structure then transferred the modified inputs for PTSD diagnosis. We investigated all possible combinations of the three feature categories aiming to identify the best feature combination for PTSD diagnosis.

The remaining of the paper is organized as follows. Several Technical terms are defined in Section 2 to help readers grasp the concepts. In Section 3, we describe each of the components in our diagnosis system including deep learning model, transfer learning, and feature selection. Hypothesis testing including Paired t-test, McNemar's test, and Wilcoxon Signed-rank test are also introduced. Section 4 presents the feature investigation and statistical testing results, and Section 5 concludes the paper.

2. TERMINOLOGY

Neural network: a widely-used model that simulates the human brain activity for function approximation or pattern classification.

Neurons: a computing unit in neural network that processes information received and transmits it to other units.

Layer: a set of neurons at the same depth in a neural network.

Hidden units: any units located between input and output units in a neural network.

Restricted Boltzmann machine: a specialized neural network structure and training mechanism.

Weights: connections among neurons or units in a neural network that weight incoming information.

Backpropagation: a training algorithm that adjusts the network weights to maximize its performance rate.

Support Vector Machine: a specialized neural network that can be used for function approximation and classification.

Overtraining: an issue that refers to a model will be over-fitted to training data if it contains too few examples so that the trained model cannot be generalized to unseen data.

Transfer Learning: a learning strategy that transfer knowledge from source domain to solve problems in target domain where training examples are limited in source domain but abundant in target domain, as long as the two domains are similar.

SoftMax Layer: usually refers to the final layer in a classification neural network model that output posterior probabilities for the classification purpose.

3. METHOD

3.1 Data Preparation Feature Extraction

Two groups of speech data were collected for this study: 26 recorded speech segments extracted from YouTube, and 26 others collected at a hospital in Ohio, for a total of 52 individuals. Among the individuals, half of them from both sources were diagnosed as having PTSD. During the interview, speech signal was recorded for each subject and the duration of the recordings from the subjects varies approximately between 51 seconds and 480 seconds. Each recording was sampled at 44.1 kHz.

We defined a 3-second speech segment as a frame, and shift 1 second forward to create another frame. For each speech frame, we further divided it into a set of 25ms segments with 10ms overlapping between two consecutive segments. We extracted the three categories of features as listed in Table 1 for each of the 25ms segments and computed averages for the features. In total, we obtained 162 averaged features

for each frame. Finally, we concatenated features from 15 consecutive speech frames as a feature vector for PTSD diagnosis. The feature vector has a size of 2430 (162x15). To create next feature vector, we shift to the next frame such that two consecutive feature vectors have 14 frames overlapping. For a non-PTSD patient, the class ID is 0 while for a PTSD patient is 1. Each subject has a matrix of features with each row represents a feature vector from 15 speech frames. The feature matrix will be used for PTSD diagnosis in the later sections.

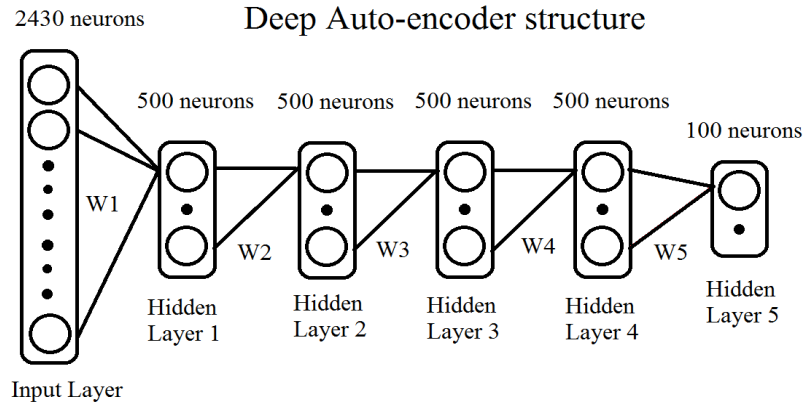
3.2 Deep Learning Model

The deep belief model we developed for PTSD diagnosis consists of several fully connected hidden layers with weights in between neurons on each layer. The model structure is shown in Figure 2. We utilized the restricted Boltzmann machine (RBM) mechanism [6] to pre-train the deep model and ran backpropagation to fine-tune the weights [6]. As a result, a model with the structure of 2430-500-500-500-100 model is obtained for PTSD diagnosis. The deep belief model shown in Figure 2 is a classifier itself but often it is widely used for feature learning, where we extract outputs from these hidden layers as new representations for the original inputs and train a different classifier such as the support vector machine (SVM) [6] for PTSD diagnosis.

Table 1: Features computed from speech signals.

Feature Position in Feature Vector		
Feature	Position	# of features
Prosodic Features		
Short-Time Energy	1	1
Average Power	2	1
Average Magnitude	3	1
# of Zero Crossings	4	1
Mean	5	1
Median	6	1
Standard Deviation	7	1
Minimum	8	1
Maximum	9	1
Range	10	1
Dynamic Range	11	1
Interquartile Range	12	1
Vocal-tract Features		
MFCC	13-51	39
Teager Energy Operator	52	1
Excitation Features		
Jitter	53	1
Shimmer	54	1
Total # of Original Features	54	
1 st Order Time Derivative	55-108	54
2 nd Order Time Derivative	109-162	54
# of Total Feature Per Frame	162	

Figure 1: The structure of Auto-encoder we used in this study. W1, W2, etc. are the weights connecting different hidden layers



3.3 Transfer Learning

The model shown in Figure 2 has millions of parameters and it is very likely to over-fit the small dataset we collected. We utilized transfer learning to resolve this issue. In transfer learning, we used the TIMIT speech database [5] to train the deep belief model as shown in Figure 3 for phoneme recognition. TIMIT consists of speeches from 630 individuals with 8 dialects in American English. To train the deep model for speech recognition, we attached SoftMax layer on top of the model to recognize the 39 phoneme classes. We used the same set of features as for PTSD diagnosis on the TIMIT dataset. Note that the features we used for speech recognition may not be a good choice but the purpose here is to borrow knowledge from TIMIT to diagnose PTSD. Once the deep model was trained on the TIMIT dataset for speech recognition, we used the deep model as a feature extractor on PTSD dataset to extract new representations from the hidden layers, and trained a SVM classifier on the representations for PTSD diagnosis.

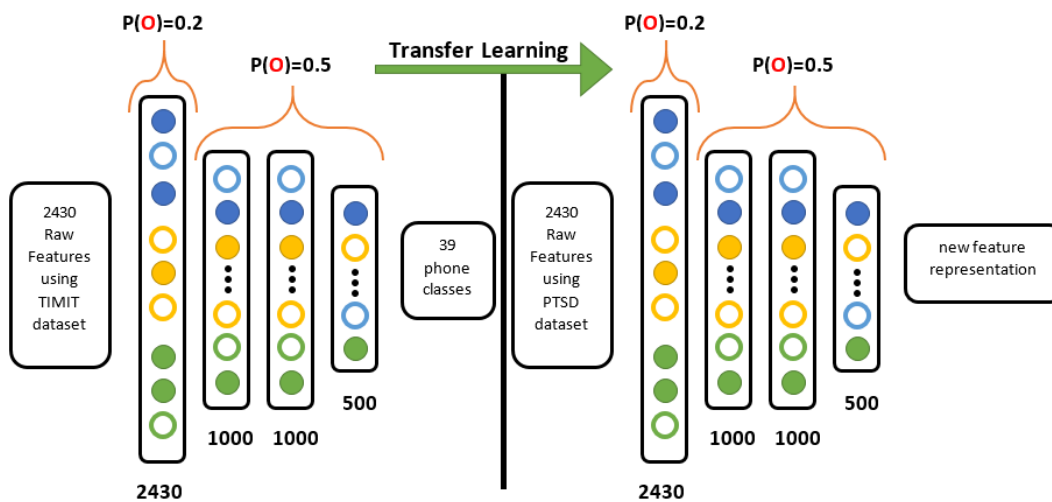


Figure 3: The visual representation of the Transfer learning method. After using the TIMIT dataset to optimize the network architecture, we substituted the PTSD dataset on this pre-trained network to obtain new representations for PTSD diagnosis

3.4 Feature Selection in Transfer Learning

In Section 2.1, we present three categories of features as shown in Figure 1. The goal of feature investigation and selection is to assess and identify which feature category would optimize the classification performance. First, we trained the 5-layer deep model on TIMIT and then we used the trained model as a feature extractor to obtain new representations for PTSD features in each of the 5 hidden layers. We then evaluated the performance of each layer representation on both YouTube and Ohio datasets. Secondly, we excluded one feature category at a time by zeroing the features at the input layer in the deep model and repeated the previous steps, and likewise excluded two categories at a time. Since the first hidden layer showed the best overall performance in classification, all data presented in Section 4 are obtained from first hidden layer 1 only.

3.5 Performance Evaluation

We utilized the Leave One Subject Out—Cross Validation (LOSO-CV) method to evaluate a model for PTSD diagnosis. In LOSO-CV, we left one subject for testing and remaining data were used to train a SVM classifier. This process will be repeated until each subject was tested once, and then training and testing accuracies were computed. We have used two performance metrics in this paper: segment-wise accuracy and subject-wise accuracy. If a subject’s segment-wise accuracy surpassed 50%, the subject is considered as correctly classified.

3.6 Hypothesis Testing

Different feature categories have different performances. We applied the following hypothesis tests to identify if the difference is significant.

3.6.1 Paired t-test

Paired t-test [2] compares two quantitative population means where observations from one population are paired with the observations from the other population. The assumption is that the observations follow normal distribution. In this case, we paired up the initial test accuracies (the control group) and accuracies from each individual attempts of removal (the paired groups), assuming that the data are normally distributed. Paired t -test is used to assess if two sets of segment-wise accuracies from the control group and the paired group are significantly different. The t statistic is determined based on the mean and standard deviation of the group difference, and the degree of freedom. Based on the t statistic, we calculate the two-tailed p -value using the t -distribution with degree of freedom. If p -value is smaller than 0.05, we would reject the null hypothesis that the means of two distributions are equal.

3.6.2 Wilcoxon Signed-Rank Test

Since the normal distribution assumption of Paired t-test may not be true, we introduce a non-parametric statistical hypothesis test—Wilcoxon signed-rank test [3], which does not require any distribution assumptions. For two paired samples, it is to assess whether the population mean ranks differ. After calculating the test statistic W , also called as the sum of the signed ranks, we compare it to the critical value from the reference table based on the degree of freedom. The null hypothesis is that there is no significant difference of the mean ranks in the two paired samples, which denotes that W statistic is smaller than the critical value at certain degree of freedom. However, if the test statistics W is greater than or equal to the critical value, we would reject the null hypothesis, and it is considered there is a significant difference between the means of the two paired groups.

3.6.3 McNemar’s Test

McNemar’s test is a statistical test used on paired nominal data. It is applied to a 2x2 contingency table with two dichotomous traits, with matched pairs of subjects, to assess the significance of the difference between two correlated proportions [4]. We utilized McNemar’s test to assess if two sets of subject-wise accuracies are significantly different. In our experiment, we define the classification results obtained from

the initial transfer learning, i.e. with all three feature categories (Prosodic, Vocal-Tract, and Excitation features) present (as the control group), and the classification results obtained from each combination of transfer learnings with only one category or two categories present as paired groups; hence, we compare each paired group with the control group in the 2x2 contingency table to find out whether the marginal distribution differs significantly. The null hypothesis is that the accuracy of the control group is equal to that of each paired group, which means the percentage of the subjects with True classification in the control group and that in each paired group are essentially the same. If the absolute percentage difference exceeds an extreme value at the significance level (0.05), we would reject the null hypothesis and declare the removal of the features is significantly impacting the accuracy of the classification.

4. RESULT AND DISCUSSION

In the tables provided below, we abbreviated Prosodic features as “P”, Vocal-tract features as “V”, and Excitation features as “E”. For the six experimental groups, we defined “Out” as the feature category being excluded, and “Only” as the feature category remains while the others being excluded. To provide a more comprehensible view, we bolded all data that show statistical significance, or that fit the criterion of rejecting the null hypothesis, for every table in this section. Table 2 summarizes segment-wise accuracies using different feature combinations computed from transfer learning for PTSD diagnosis. While other average test accuracies are approximately between 60% and 80%, the results obtained by using excitation features only are potential outliers in this data set, with test accuracies less than 50%. Table 3 summarizes subject-wise accuracies using different feature combinations computed from transfer learning for PTSD diagnosis. It can be easily observed that the fractions of subjects correctly classified using excitation features only are only half that of using other feature categories in the same dataset. When converted to percentages, the results for Excitation features only are less than 50%.

Table 2: Summaries of segment-wise accuracies of PTSD diagnosis.

Segment-wise Accuracies		Original	P Out	V Out	E Out	P Only	V Only	E Only
YouTube	Average(Test)	78.82	73.96	80.05	80.87	78.49	75.36	49.56
	Std(Test)	36.380	35.995	33.646	33.622	36.478	34.477	37.562
Ohio	Average(Test)	61.46	63.71	61.80	63.49	67.69	65.26	35.69
	Std(Test)	34.446	34.900	37.930	33.512	36.104	33.560	32.660

Table 3: Summaries of subject-wise accuracies of PTSD diagnosis.

Subject-wise Accuracies		Original	P Out	V Out	E Out	P Only	V Only	E Only
YouTube	In Fractions	21/26	20/26	22/26	22/26	21/26	20/26	12/26
	In %	80.77	76.92	84.62	84.62	80.77	76.92	46.15
Ohio	In Fractions	16/26	16/26	15/26	19/26	18/26	16/26	6/26
	In %	61.54	61.54	57.69	73.08	69.23	61.54	23.08

In summary, excitation features are least effective in PTSD diagnosis in majority of the experiments, and prosodic features seem to be the most effective feature category. To find out whether any of the results are statistically significant, we carried out Paired T-test and McNemar’s test for the individual results using layer 1 for classification.

4.1 Paired t-test

For paired t-test, since the subjects remain the same but the features change, we consider the two sets of data “repeated measure”, that is, are dependent upon each other. The advantage of using this test is that we can calculate the difference between the mean of two normal distributions, and find out whether their difference is large enough by comparing to the confidence level. It is straightforward to carry out because we assumed both sample data are normally distributed, and hence we can use this equation shown below:

$$t = \frac{\bar{x} - \mu_0}{\frac{s}{\sqrt{n}}}$$

Where \bar{x} is the mean of the differences between two sets of data, s is the standard deviation of the differences, and n is the number of pairs, in this case, the number of subjects. μ_0 is zero since we just want to know whether there is a significant difference. The degree of freedom is $n-1$. Then we calculated the p -value through t distribution. Since the common confidence level is 0.05, we expected p -value to be less than 0.05, if the difference is significant.

We used a matched-pair design to find the t -value, corresponding the classification performance before and after the removals of certain feature categories. All results are compared to namely the control group. We can see from Table 4 that the absolute value of t -values for E only results are much larger than the rest of the t -values, and the p -values are within the range that enables the rejection of the null hypothesis. Using the confidence level of 0.05, we observed that the p -values for excitation features are clearly less than 0.05. Thus, we declared this result as significant.

Table 4. Paired t -test results on segment-wise accuracies.

Segment-wise Accuracies		P Out	V Out	E Out	P Only	V Only	E Only
YouTube	t-value	-1.53	0.14	0.87	-0.04	-1.00	-2.98
	p-value	0.14	0.89	0.39	0.97	0.33	0.01
Ohio	t-value	0.80	0.05	1.57	0.82	1.55	-2.75
	p-value	0.43	0.96	0.13	0.42	0.13	0.01

4.2 Wilcoxon Signed-Rank Test

Still, only two pairs of data show significance. We then deduced that the original data is not normally distributed as we assumed, so we conducted Wilcoxon's signed-rank test as an alternative for paired t -test since it does not require the sample data to be normally distributed. The method is calculating the difference between two sets of data, as we did in paired t -test. Then we ranked them by the absolute value of the differences. After that, we multiplied the ranks by the signs of the corresponding original differences. For instance, a subject's segment-wise accuracy increased after excluding some features, and then the sign of difference must be positive. Similarly, if subject's segment-wise accuracy decreased, and then the sign of difference must be negative. Finally, we sum up all the signed rank together to get the w test statistic and compare it to the critical value in Table 6 based on the degree of freedom, which is the number of subjects that have a sign in their difference, that is, comparing to the subjects that have a change of 0 in testing accuracy. Surprisingly, not only E only tests for both YouTube and Ohio but also V only test for Ohio dataset, show significance as shown below in the table, for the individual w -value surpasses the corresponding critical value for the w statistic. Since w -value for V only test for Ohio dataset is slightly above its critical value, and that same test does not show significance for YouTube dataset, we can safely assume that the significance is due to sampling variability. By comparing the experimental w -values with the critical values for the w -statistic, we can see that only when Excitation features remains, the critical value falls far below the w -value, which suggests significance of this result.

Table 5: Wilcoxon test results on segment-wise accuracies.

Segment-wise Accuracies		P Out	V Out	E Out	P Only	V Only	E Only
YouTube	w-value	55	9	8	2	22	190
	deg. of freedom	14	21	12	21	15	24
	critical value for the w statistic	21	58	13	58	25	81
Ohio	w-value	26	12	47	58	90	179
	deg. of freedom	23	23	22	23	23	25
	critical value for the w statistic	73	73	65	73	73	89

Table 6: The critical value table based on degree of freedom and significance level.

deg. of freedom	Two Tailed significance levels		
	0.05	0.02	0.01
12	14	10	7
13	17	13	10
14	21	16	13
15	25	20	16
16	30	24	20
17	35	28	23
18	40	33	28
19	46	38	32
20	52	43	38
21	59	49	43
22	66	56	49
23	73	62	55
24	81	69	61
25	89	77	68

4.3 McNemar's Test

Since binomial outcomes are obtained from the testing accuracy: correct or incorrect, we can use the McNemar's test to see whether there is difference existed in feature changes. Some subjects showed improvement in accuracy and turned from incorrect to correct, some showed diminishment and went the other way, and some showed the same results after both tests. But overall improvement is what we want to see. Setting the threshold of McNemar's Test based on 50% segment-wise accuracy, we define any subject which has test accuracy above 50% as "True", while below 50% as "False". The number in the corresponding row is the number of subjects that fit in the condition. We can clearly see that more subjects turned from true to false, and that change showed significance in the statistical mean.

In order to calculate the p-value with a small data size, i.e. with the total number of discordant subjects being less than 25, we used a chi-square distribution with an X^2 statistics that is corrected for continuity as shown below [7].

$$X^2 = \frac{(|b - c| - 1)^2}{b + c}$$

With a degree of freedom of 1, we found the p-value by using the X^2 statistics calculated above. The results are shown below in Table 7 and 8, for YouTube and Ohio dataset respectively.

We have defined the confidence level as 0.05, which indicates that if the p -value is less than 0.05, the corresponding results suggest significance. From the table, we found the test accuracies for Excitation features only have the p -value clearly below 0.05. Hence, we reckon the removal of prosodic features and vocal-tract features significantly changed the classification performance.

Table 7: McNemar's test for results on subject-wise accuracies of YouTube subjects.

YouTube		Control		p-value
		True	False	
P Out	True	19	1	1.000
	False	2	4	
V Out	True	19	3	1.000
	False	2	2	
E Out	True	21	1	1.000
	False	0	4	
P Only	True	18	3	0.683
	False	3	2	
V Only	True	19	1	1.000
	False	2	4	
E Only	True	10	2	0.027
	False	11	3	

Table 8: McNemar's test for results on subject-wise accuracies of Ohio subjects.

Ohio		Control		p-value
		True	False	
P Out	True	15	1	0.480
	False	1	9	
V Out	True	11	4	1.000
	False	5	6	
E Out	True	16	3	0.248
	False	0	7	
P Only	True	12	6	0.752
	False	4	4	
V Only	True	15	1	0.480
	False	1	9	
E Only	True	4	2	0.016
	False	12	8	

4.4 Discussions

Since all of the statistical tests show that the reduction of both prosodic features and vocal-tract features has a considerable effect on the accuracy of classification, we concluded that excitation features are the least significant among all three categories. However, considering there are only 2 features in this category, while 12 in prosodic features and 40 in vocal-tract features, there may be biases confounded in the results. Also, we noticed that after excluding the vocal-tract Features, the performance rate increased by a little for YouTube data. It may be by chance, but it may also imply that there are other lurking variables hidden in the scenes.

All three tests suggested that excitation features are the least important among the three categories. We rejected all three null hypotheses when comparing the control group with all features present and the experimental group with only excitation features present. It is because the data obtained from this experimental group showed significant diminishments in the classification performance that could not have happened by chance. When we only kept prosodic features or vocal-tract features, the classification performances did not show significant changes, which means the changes observed can be explained by chances. Moreover, when excitation features were excluded, the performance improved slightly as shown in Tables 1 & 2. Though the changes are not significant, we can still deduce from the results that this feature category might be a disturbing factor to the overall classification performance. Tables 1 & 2 included both segment-wise accuracy and subject-wise accuracies during testing, providing a more perceptible evidence to this implication.

5. CONCLUSION

In this experiment, we conducted transfer learning through DBN to identify whether a patient has PTSD based on speeches, and obtained up to 80% accuracy for the PTSD data set. We also investigated the importance of the features used in order to eliminate the disturbing factors, by removing specific features and comparing the results. Finally, we concluded that excitation feature category is the least significant one as indicated by multiple statistical tests. Our future work includes setting individual feature to zero at a time, instead of setting the whole category zero, since the number of features in each category varies. By comparing the accuracy obtained from that, we might get a higher accuracy.

REFERENCE

- [1] American Psychiatric Association (2013). Diagnostic and Statistical Manual of Mental Disorders (5th ed.). Arlington, VA: American Psychiatric Publishing. pp. 271–280. ISBN 978-0-89042-555-8.
- [2] David, H. A.; Gunnink, Jason L. (1997). "The Paired t Test Under Artificial Pairing". *The American Statistician*. 51 (1): 9–12. doi:10.2307/2684684. JSTOR 2684684.
- [3] Wilcoxon, Frank (Dec 1945). "Individual comparisons by ranking methods". *Biometrics Bulletin*. 1 (6): 80–83.
- [4] McNemar, Quinn (June 18, 1947). "Note on the sampling error of the difference between correlated proportions or percentages". *Psychometrika*. 12 (2): 153–157. doi:10.1007/BF02295996. PMID 20254758.
- [5] Garofolo, John, et al. TIMIT Acoustic-Phonetic Continuous Speech Corpus LDC93S1. Web Download. Philadelphia: Linguistic Data Consortium, 1993.
- [6] Banerjee, D. (2017). *Emotional State Recognition and PTSD Detection Based on Speech Signals*.
- [7] Edwards, A (1948). "Note on the "correction for continuity" in testing the significance of the difference between correlated proportions". *Psychometrika*. 13: 185–187. doi:10.1007/bf02289261.

DECISION ANALYSIS TREE TO DETERMINE INJURY PREVENTION STRATEGY FOR ANKLE SPRAIN INJURY BASED ON COST

Emily M. Hartley
College of Health Sciences
Old Dominion University
1121 Health Sciences Building
Norfolk, VA, USA
ehart001@odu.edu

Matthew C. Hoch
School of Physical Therapy and Athletic Training
Old Dominion University
3120 Health Sciences Building
Norfolk, VA, USA
mhoch@odu.edu

ABSTRACT

Ankle sprains are one of the most common injuries amongst the physically active population. Injury prevention strategies such as injury prevention programs and injury risk screening assessments have been developed to aid in injury prevention. However, few studies have investigated the most cost-effective strategy to prevent ankle sprain injuries. A decision tree was created to assist in determining which strategy to choose based on cost. The use of an injury prevention program without the use of injury risk screening was the best decision based on cost in the collegiate athletic setting.

Keywords: Cost-Effectiveness, Injury Prevention, Ankle Sprain, Collegiate Athletes

1 INTRODUCTION

Ankle sprains are one of the most common injuries among physically active individuals, accounting for 14% of collegiate athletic injuries (Hootman, Dick, and Agel 2007). Individuals who sustain an ankle sprain are likely to develop a condition known as chronic ankle instability (CAI)(Gerber et al. 1998). CAI is denoted by repetitive ankle sprains, sensations of “giving way”, and residual symptoms (Hubbard et al. 2007). These individuals are likely to suffer from long term consequences such as an early development of osteoarthritis and a decreased health related quality of life (Houston, Van Lunen, and Hoch 2014, Arnold, Wright, and Ross 2011). In addition to the short and long term consequences, ankle sprains have a substantial economic impact on the patient and healthcare system (Knowles et al. 2007).Based on these long term consequences, it is important to prevent ankle sprains to avoid the additional negative sequelae associated with the injury.

Injury prevention programs and injury risk screening assessments have been developed to prevent ankle sprain injuries. Neuromuscular injury prevention programs have been effective at reducing the risk of ankle sprain injury (Hübscher et al. 2010). However one of the major limitations of this method of injury prevention is compliance of the participants to complete the prescribed exercises (Sugimoto et al. 2012). Some screening assessments, such as the start excursion balance test (SEBT), have been used to first stratify individuals into a high or low risk group (Gribble et al. 2016). The high-risk group would then participate in an injury prevention program to reduce the risk of injury. The low risk group would proceed with no treatment.

Very few studies have investigated the cost-effectiveness of ankle sprain prevention strategies (Hupperets et al. 2010, Verhagen et al. 2005) One study found that a home balance training program used in those with a recent ankle sprain was more cost-effective at preventing further ankle sprain injuries than usual care alone (Hupperets et al. 2010) Oppositely, an additional study found that no treatment was more cost-effective in the prevention of ankle sprains when compared to an injury prevention program within volleyball athletes (Verhagen et al. 2005) There is a lack of research suggesting which injury prevention strategy is the most cost-effective at preventing ankle sprain injuries within collegiate athletes. Therefore, the purpose of this

study was to use a decision analysis tree to determine which strategy should be chosen to prevent ankle sprain injuries within a collegiate athletic setting based on cost.

2 METHODS

2.1 Study Design

An economic evaluation of ankle sprain prevention strategies was created using decision analysis software (TreeAge Pro; TreeAge Software, Williamstown, Massachusetts). The hypothetical cohort of this model was male and female collegiate athletes. The model compared three different strategies for the prevention of ankle sprain injuries by predicting the outcomes (ankle sprain/no ankle sprain) and the associated costs of each.

2.2 Model Design

A decision tree was created (Figure 1) to compare three strategies based on cost for the prevention of ankle sprain injuries: 1) injury prevention program 2) injury screening session 3) no prevention. The injury prevention program consisted of neuromuscular training programs completed on a regular basis. The injury screening session was based on the utilization of the SEBT to stratify individuals into a high risk and low risk group. The high-risk group would then participate in the injury prevention program while the low risk group did not. The no prevention group did not participate in the injury prevention program or injury risk assessment.

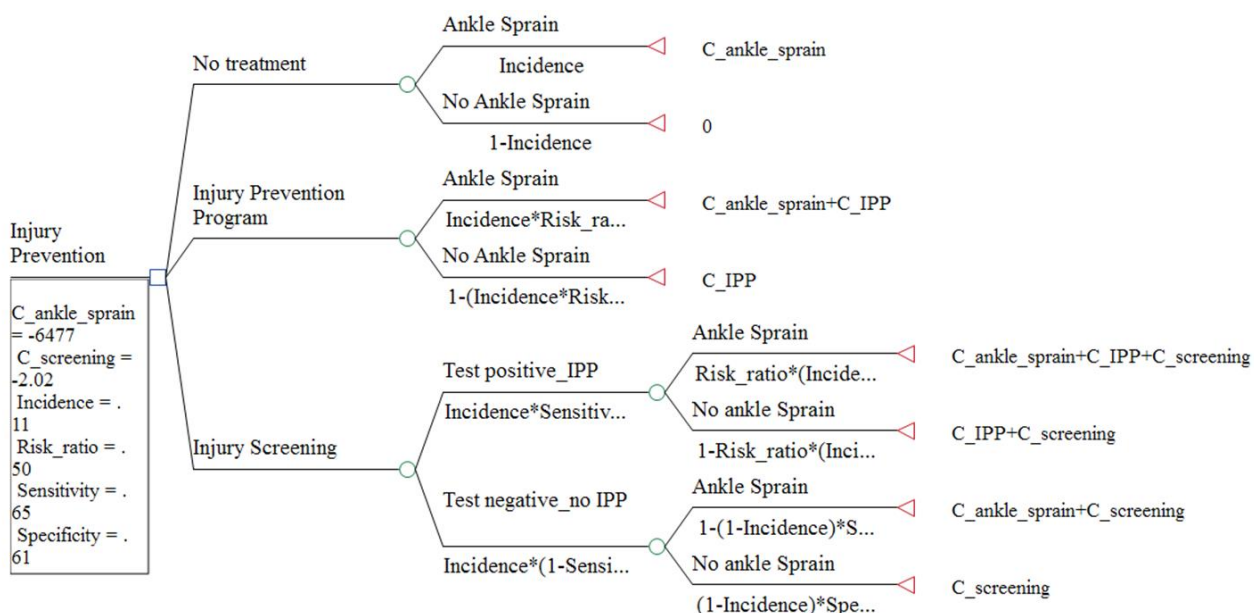


Figure 1. Decision Analysis Tree to Determine Optimal Choice of Ankle Sprain Prevention Based on Cost

2.3 Model Inputs

2.3.1 Costs

The main cost inputs within this model were the costs associated with ankle sprain treatment, implementation of the injury prevention program, and administration of the lower extremity injury screening. The costs of ankle sprain injury treatment varied from \$723.00 for an emergency room visit to \$11925.00 for overall costs associated with the healthcare visits, treatment, loss of future earnings, and costs related to reduction in quality of life [11, 12].

The costs associated with the implementation of the injury prevention program were based on previous prospective studies that implemented neuromuscular training programs. The approximate costs of implementing an injury prevention program were similar within all three studies ranging from \$28.48-\$29.23 per athlete [5, 9]. These approximate costs only included the equipment necessary to complete the exercises and instructional materials. The use of a facilitator to implement and instruct the injury prevention program was not considered.

The costs associated with the administration of the lower extremity injury screening session were based on the sole use of the SEBT. It was assumed that an athletic trainer would be administering the functional test and that the athletic trainer would be employed at the collegiate level. The mean salary of these individuals was \$42,000 (\$36,000-\$49,000) per year. The pay per minute was then calculated and used to multiple that amount by the time necessary to administer the test. It takes approximately 6 minutes per athlete to complete the SEBT. The mean cost of administering the SEBT as an injury risk screening tool was \$2.02 per athlete.

2.3.2 Clinical Outcome Probabilities

The clinical outcome probabilities utilized within the decision tree consisted of the incidence of ankle sprain injury and the potential risk reduction gained by participating in an injury prevention program. Several studies found an incidence of ankle sprain injury within collegiate athletes of 0.11. (Gribble et al. 2016) The average risk ratio reduction found within the literature was 0.5 (0.31-0.79). (Hübscher et al. 2010) The incidence was used to determine the probability that an individual would sustain an ankle sprain had they not participated in an injury prevention program. Additionally, the incidence and risk ratio reduction were utilized together to determine the probability of sustaining an ankle sprain injury given the individual participated in an injury prevention program.

2.3.3 Sensitivity and Specificity of Screening

The sensitivity and specificity of the SEBT was utilized to determine if individuals would correctly be stratified into the high risk or low risk group. The sensitivity of the SEBT was 0.65 and specificity was 0.61. (Gribble et al. 2016) The sensitivity and specificity were combined with the risk ratio reduction and incidence of injury to determine the probability of sustaining an ankle sprain given the individual was stratified into the high-risk group. The sensitivity and specificity were also paired with incidence of injury to determine the probability of sustaining an ankle sprain given the individual was stratified to the low risk group.

2.4 Data Analysis

The decision analysis tree was run to determine the optimal choice between injury prevention program participation, injury screening session, and no treatment as a prevention strategy for ankle sprain injuries within collegiate athletes based on cost. A one-way sensitivity analysis was used to determine if the cost of the ankle sprain injury treatment would influence the optimal choice. An additional one-way sensitivity analysis was used to determine if the risk reduction ratio influenced the optimal choice based on cost.

3 RESULTS

The results of the decision tree indicated the optimal choice for ankle sprain prevention was participation in the injury prevention program seen in Figure 2. The estimated cost per athlete of those who participated in an injury prevention program was \$385.46 while the cost of injury screening was \$495.17. The no treatment choice was the least optimal costing \$712.47 per athlete. The sensitivity analysis (Figure 3.) revealed when the cost of ankle sprain treatment drops below \$777.87, injury screening becomes the most optimal method. However, when the cost of ankle sprain treatment is above \$777.87, injury prevention program participation is the optimal choice. Additionally, the difference between costs continues to increase as the cost of ankle sprain treatment increases further indicating injury prevention program participation is

the optimal choice. The risk reduction ratio did not play a role in the optimal choice of ankle sprain prevention (Figure 4).

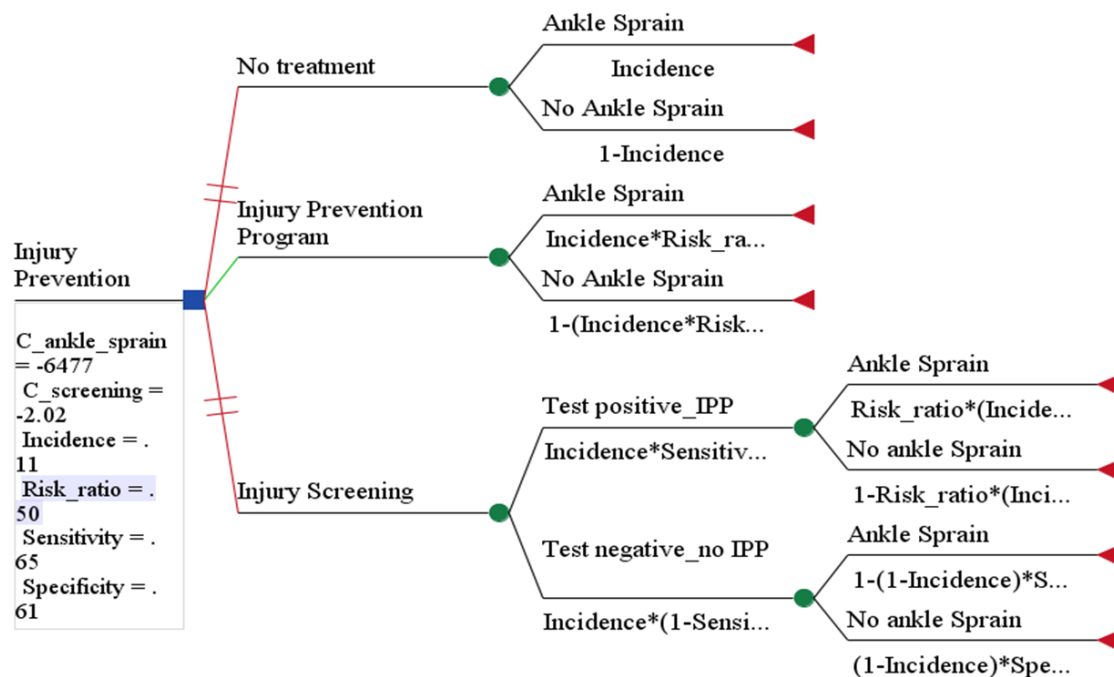


Figure 2. Results of Decision Analysis Tree

4 DISCUSSION

There is a lack of studies investigating the most cost-effective method of preventing ankle sprain injuries. The results of our decision analysis tree indicated the best method for preventing ankle sprain injury within the collegiate athletic setting based on cost is implementing an injury prevention program with all athletes. However, the sensitivity analysis revealed when the potential cost of ankle sprain treatment drops below \$777.87, injury screening session with a neuromuscular training for high risk individuals becomes the best choice.

There are several limitations within this study. The data utilized to create the probabilities and costs are gained from previous collected data within the literature. Additionally, there were limited resources related to the estimated cost of ankle sprain treatment. The costs within the literature varied widely. Future research should further investigate the true cost of ankle sprain injury treatment. Currently there is a gap in the literature regarding the health related quality of life in individuals with a history of ankle sprain injury. Future research should investigate the quality of life in individuals who sustain an ankle sprain injury.

The results of this study indicate the dominant strategy for ankle sprain prevention within collegiate athletes is participation in an injury prevention program. However, if the cost of ankle sprain treatment is lower than \$777, injury screening may be a more optimal choice based on cost. Clinicians should also consider the impact on quality of life that lower extremity injuries introduce when deciding which prevention technique to choose for their patients. Additionally, additional lower extremity injury screening tools with better sensitivity and specificity may need to be considered in the future.

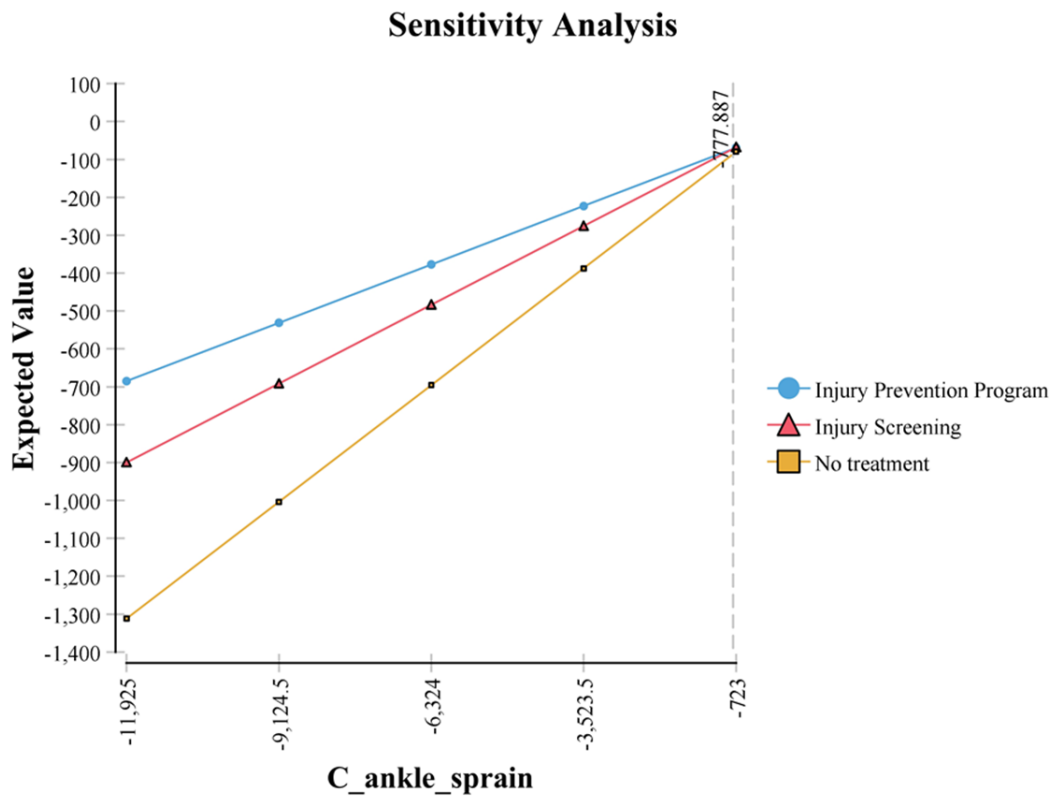


Figure 3. Sensitivity Analysis of Effect of Ankle Sprain Treatment Cost

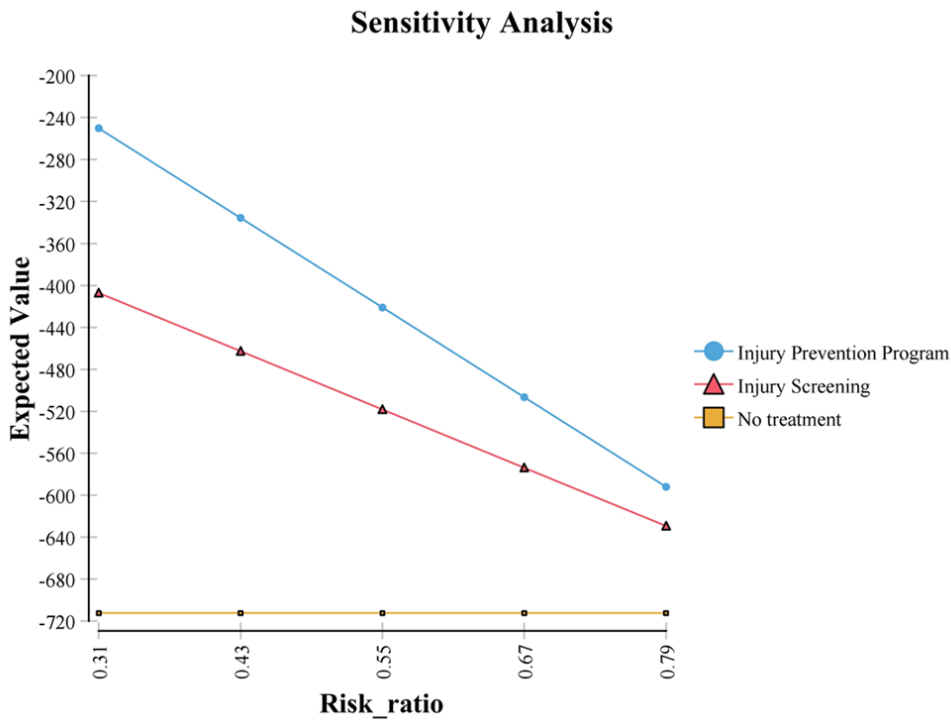


Figure 4. Sensitivity Analysis of Effect of Risk Reduction Ratio

Acknowledgements: This work was supported by the modeling and simulation fellowship at Old Dominion University

REFERENCES

- Arnold, Brent L., Cynthia J. Wright, and Scott E. Ross. 2011. "Functional Ankle Instability and Health-Related Quality of Life." *J Athl Train* 46 (6):634-641.
- Gerber, J. P., G. N. Williams, C. R. Scoville, R. A. Arciero, and D. C. Taylor. 1998. "Persistent disability associated with ankle sprains: a prospective examination of an athletic population." *Foot Ankle Int* 19 (10):653-660.
- Gribble, Phillip A., Masafumi Terada, Megan Q. Beard, Kyle B. Kosik, Adam S. Lepley, Ryan S. McCann, Brian G. Pietrosimone, and Abbey C. Thomas. 2016. "Prediction of Lateral Ankle Sprains in Football Players Based on Clinical Tests and Body Mass Index." *Am J Sports Med* 44 (2):460-467 8p. doi: 10.1177/0363546515614585.
- Hootman, J. M., R. Dick, and J. Agel. 2007. "Epidemiology of collegiate injuries for 15 sports: summary and recommendations for injury prevention initiatives." *J Athl Train* 42 (2):311-319.
- Houston, Megan N., Bonnie L. Van Lunen, and Matthew C. Hoch. 2014. "Health-Related Quality of Life in Individuals With Chronic Ankle Instability." *J Athl Train* 49 (6):758-763. doi: 10.4085/1062-6050-49.3.54.
- Hubbard, T. J., L. C. Kramer, C. R. Denegar, and J. Hertel. 2007. "Contributing factors to chronic ankle instability." *Foot Ankle Int* 28 (3):343-354.
- Hübscher, M., A. Zech, K. Pfeifer, F. Hänsel, L. Vogt, and W. Banzer. 2010. "Neuromuscular training for sports injury prevention: a systematic review." *Med Sci Sports Exerc* 42 (3):413-421. doi: 10.1249/MSS.0b013e3181b88d37.
- Hupperets, M. D. W., Ealm Verhagen, M. W. Heymans, J. E. Bosmans, M. W. van Tulder, and W. van Mechelen. 2010. "Potential Savings of a Program to Prevent Ankle Sprain Recurrence: Economic Evaluation of a Randomized Controlled Trial Maarten." *American Journal of Sports Medicine* 38 (11):2194-2200 7p. doi: 10.1177/0363546510373470.
- Knowles, S. B., S. W. Marshall, T. Miller, R. Spicer, J. M. Bowling, D. Loomis, R. W. Millikan, J. Yang, and F. O. Mueller. 2007. "Cost of injuries from a prospective cohort study of North Carolina high school athletes." *Inj Prev* 13 (6):416-421.
- Sugimoto, Dai, Gregory D. Myer, Heather M. Bush, Maddie F. Klugman, Jennifer M. Medina McKeon, and Timothy E. Hewett. 2012. "Compliance With Neuromuscular Training and Anterior Cruciate Ligament Injury Risk Reduction in Female Athletes: A Meta-Analysis." *J Athl Train* 47 (6):714-723 10p. doi: 10.4085/1062-6050-47.6.10.
- Verhagen, E. A. L., M. van Tulder, A. J. van der Beek, L. M. Bouter, and W. van Mechelen. 2005. "An economic evaluation of a proprioceptive balance board training programme for the prevention of ankle sprains in volleyball." *Br J Sports Med* 39 (2):111-115 5p.

HOMEOMORPHIC TETRAHEDRAL TESSELLATION FOR BIOMEDICAL IMAGES

Jing Xu

Department of Computer Science
Old Dominion University
4700 Elkhorn Ave
Norfolk, VA, USA
jxu@cs.odu.edu

Andrey N. Chernikov

Department of Computer Science
Old Dominion University
4700 Elkhorn Ave
Norfolk, VA, USA
achernik@cs.odu.edu

ABSTRACT

We present a novel algorithm for generating three-dimensional unstructured tetrahedral meshes for biomedical images. The method uses an octree as the background grid from which to build the final graded conforming meshes. The algorithm is fast and robust. It produces meshes with high quality since it provides dihedral angle lower bound for the output tetrahedra. Moreover, the mesh boundary is a geometrically and topologically accurate approximation of the object surface in the sense that it allows for guaranteed bounds on the two-sided Hausdorff distance and the homeomorphism between the boundaries of the mesh and the boundaries of the materials. The theory and effectiveness of our method are illustrated with the experimental evaluation on synthetic and real medical data.

Keywords: tetrahedralization, quality, fidelity, homeomorphism

1 INTRODUCTION

With medical data sets, one can generate conforming quality meshes of the spatially realistic domains that help producing computer aided visualization, manipulation, and quantitative analysis of the multi-dimensional image data. The domain of focus often possesses heterogeneous materials that specify functionally different characteristic properties, so it is usually segmented into multiple regions of interest (materials). In finite element analysis, each material of interest is assigned an individual attribute or a material coefficient. Thus, meshes with conforming boundaries describing each of the partitioned material regions are generated for these purposes.

The generation of geometric discretizations from segmented multi-material images presents many challenges. In particular, a mesh should meet constraints on both the shape and the size of its elements, and must conform at the material boundaries. In addition, the algorithm must handle arbitrary topology. In this paper we present an algorithm for constructing tetrahedral volume meshes that are suitable for real-time finite element analysis, i.e., they satisfy the following criteria:

1. Elements with arbitrarily small angles which cause the stiffness matrix in FE analysis to be ill-conditioned do not appear in the mesh. Specifically, we guarantee that all dihedral angles are above a user-specified lower bound which can be set to any value up to 19.47° .

2. The mesh offers a reasonably close representation (fidelity) of the underlying materials. In particular, the two-sided Hausdorff distance between the boundaries of the mesh and the boundaries of the materials respects the user specified fidelity bounds.
3. The mesh is able to offer a faithful topology of the materials. In other words, mesh boundary is homeomorphic to the object surface. We discuss the concept of homeomorphism, and give a sufficient condition for the approximation to offer homeomorphism.
4. The mesh contains a small number of elements that comply with the three guarantees above.

There has been a significant amount of work on guaranteed quality mesh construction. One approach for generating tetrahedral meshes of smooth and piecewise smooth surfaces is that the input is assumed to be an implicit function $f : R^3 \rightarrow Z$ such that points in different regions of interest evaluate f differently. One guaranteed-quality technique is based on the Delaunay refinement (Boissonnat and Oudot 2005, Cheng et al. 2010, Foteinos et al. 2014). In Foteinos et al. (2014)’s work, the authors present a Delaunay meshing algorithm with several mathematical guarantees. They proved that the tetrahedra in the output mesh have the radius-edge ratio less than 1.93. The two-sided Hausdorff distance between the object surface and mesh boundary is bounded by a user-specified parameter. Using a strategy called ε -Sample (Amenta et al. 2000, Amenta et al. 2003), the mesh boundary is proved to be ambient isotopic to the object surface. However, the method only supplies the circumradius-to-shortest-edge ratio bound. Even if this ratio is very small, it can not avoid the almost flat tetrahedra. Another guaranteed-quality technique employs a space-tiling background grid to guide the creation of a mesh (Labelle and Shewchuk 2007, Bronson et al. 2014, Chernikov and Chrisochoides 2011), the focus of this paper. Isosurface Stuffing (Labelle and Shewchuk 2007) is a guaranteed-quality tetrahedral meshing algorithm for general surfaces under the assumption that the surface is a smooth 2-manifold. It offers the one-sided Hausdorff distance guarantee from the mesh to the model. If the surface is a smooth manifold with bounded curvature, it also provides the one-sided Hausdorff distance guarantee from the model to the mesh. Using interval arithmetic, the dihedral angles for the mesh with uniform sized boundary are proved to be bounded between 10.7° and 164.8° . However, our method depends on a standard octree and decimation, and is able to achieve the minimum dihedral angle bound 19.47° .

This work builds upon the Lattice Decimation (LD) method (Chernikov and Chrisochoides 2011). LD is a tetrahedral image-to-mesh conversion algorithm that allows for guaranteed bounds on the smallest dihedral angle (up to 35.26°) and on the Hausdorff distance between the boundaries of the mesh and the boundaries of the materials. The algorithm constructs an octree and splits the octree until no leaf contains voxels from multiple materials. Then it fills the octree leaves with high-quality elements. This initial mesh has a large number of elements because it satisfies the highest dihedral angle and fidelity bounds. The authors designed a post-processing decimation step using vertex removal operation. It coarsens the mesh to a much lower number of elements while at all times maintaining the required fidelity and quality bounds. However, the decimation step is a greedy algorithm which was not designed for smooth transition in element size. In fact, it can produce clusters of small elements surrounded by much larger elements. In this work, we refine the octree to a lower level and the element size does not have so big difference, therefore, the issue can be mitigated. However, the relatively coarse mesh brings a new issue: the topological faithfulness may not be guaranteed. We designed a new technique called single manifold condition that solves this problem.

In this paper, we approximate the boundary of the materials with a set of triangular patches in octree leaves using a pre-defined surface look-up table. The triangular patches all together form a waterproof surface mesh which is homeomorphic to the boundary of the materials by proof. We achieve the two-sided Hausdorff distance bound by constructing a sequence of such surface meshes until the fidelity condition in each leaf is satisfied. The octree leaves are filled with high-quality elements from a pre-defined volume look-up table.

During decimation, the initial mesh is coarsened to a much lower number of elements while at all times the quality, fidelity and topological requirements are maintained.

2 METHODOLOGY

The proposed algorithm is described as follows: the algorithm takes a two- or a three-dimensional multi-material image as its input. The user also specifies as input the target fidelity bounds (two-sided Hausdorff distance) and the desired angle lower bound. The algorithm outputs a quality mesh which is a good geometric and topological approximation of the underlying object. We first define some concepts relative to the topological faithfulness.

Definition 2.1. Homeomorphism A topological notion of equivalence. A mapping $\mu : X \rightarrow Y$ defines a homeomorphism between two compact Euclidean subspace X and Y if μ is continuous, one-to-one and onto. The inverse function μ^{-1} is also continuous.

Definition 2.2. Single manifold condition The intersection of the image boundary with each of the boundary leaves is a $(n - 1)$ -manifold with boundary, n being the dimension.

Specificly, in two dimensional case, the image boundary edges form a 1-manifold with boundary, a chain composed by image boundary edges. On this chain, every image boundary vertex has two neighbors, except the first and last ones, which have only one neighbor. In three dimensional case, the degree of the image boundary edge is defined by the number of boundary faces that incident upon the edge. The image boundary faces form a 2-manifold with boundary, a sheet on which the degree one image boundary edges form a cycle and all the other image boundary edges are degree two edges.

2.1 Initial mesh construction

The algorithm first constructs an octree that completely encloses all the materials (except for the background voxels) from the bitmap. The boundaries between the octree leaves correspond exactly to the boundaries between the voxels. Besides that, an extra space between the materials and the exterior boundaries of the octree should be equal to or larger than the user specified fidelity bounds. Initially, the algorithm splits the octree such that the leaves satisfy the single manifold condition, and the sizes of the leaves respect the 2-to-1 ratio. Then it generates a waterproof surface mesh which is homeomorphic to the boundaries of the materials from a pre-defined surface look-up table. It computes the two-sided Hausdorff distance between the boundaries of the mesh and the boundaries of the materials in each leaf. It splits the leaf into 8 children if the Hausdorff distance of either side is larger than the user specified fidelity bounds. If at least one of the leaves was split because of the conflicting of the fidelity condition, the surface mesh is discarded. The algorithm iteratively checks the single manifold condition and the 2-to-1 ratio, and generates waterproof surface meshes until the two-sided Hausdorff distance respects the user specified fidelity bounds.

2.1.1 Waterproof surface mesh

When the octree leaves respect the single manifold condition and the 2-to-1 ratio, the algorithm generates triangular patches from each octree leaf such that all the patches form a waterproof surface mesh. To generate the waterproof surface mesh, we use an approach reminiscent of the Marching Cubes algorithm (Lorensen and Cline 1987). In contrast to the Marching Cubes algorithm, our algorithm evaluates the vertices to three values: positive, negative and zero. A vertex of a leaf is evaluated to be positive if the vertex is located inside a material, to be negative if the vertex is located outside the material, and to be zero if the vertex is located exactly on the boundary of the materials. The templates on cubes would be cumbersome for our algorithm, because there would be 3^8 cases need to be analyzed. Instead, we designed a surface look-up

table on squares for each cube face, so there are totally only 3^4 entries in the table. The surface look-up table generates intersection edges on the cube faces, and the triangular patches are generated by connecting those intersection edges with the center of the cube (see Figure 3a). We list all possible stencils in Figure 1 by grouping cases with the opposite relations to the vertex value in all corners into one case and also grouping rotationally symmetric cases.

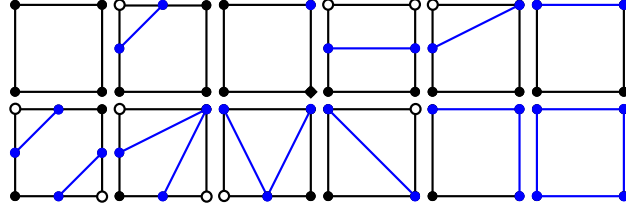


Figure 1: All possible stencils for creating intersection vertices and intersection edges by grouping cases with the opposite relations to the vertex value in all corners into one case and also grouping rotationally symmetric cases. Black circles show the positive vertices, white circles show the negative vertices, and blue circles show the zero vertices. Blue segments show the intersection edges created by the algorithm.

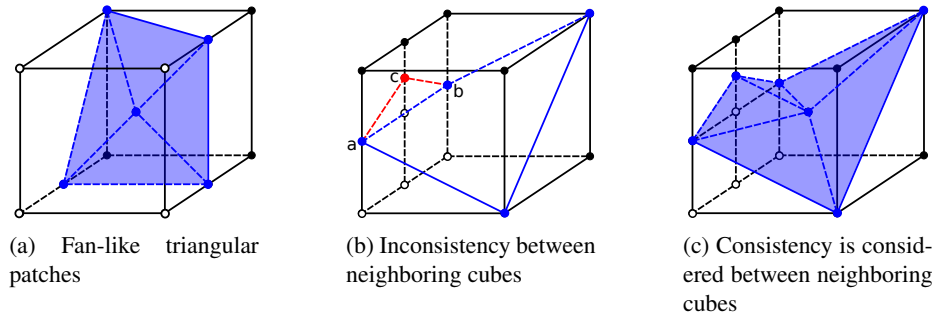


Figure 2: An illustration of generating triangular patches to approximate the image boundary in a leaf. Blue circles show the zero vertices, blue edges on cube faces show the intersection edges and blue triangles show the triangular patches.

One important issue that needs extra concern is the consistency between neighboring leaves. In the case that a face of the octree leaf shared by four smaller neighboring cubes, the consistency is not always guaranteed. For example, in Figure 3b we show a leaf whose left face is shared by four smaller neighboring cubes. From the leaf side of the view, the intersection edge generated for the left face of the leaf should be edge \overline{ab} (shown in blue dashed segment). However, from the four small neighbors side of the view, the intersection edges generated for the same face should be edge \overline{ac} and edge \overline{cb} (shown in red dashed segments). To solve this problem, we process the octree leaves in the order of their size, starting with the smallest. We duplicate the intersection edges from the processed neighbors (including the neighbors of same size and of smaller size), and generate new intersection edges for the faces that shared by the neighboring cubes which have not been processed (see Figure 2c). Then we check if the intersection edges from the six cube faces form a cycle. If all the intersection edges do not form a cycle, the leaf needs to be split. By this way, the surface mesh of our algorithm is guaranteed to be water-tight.

2.1.2 Two-sided Hausdorff distance

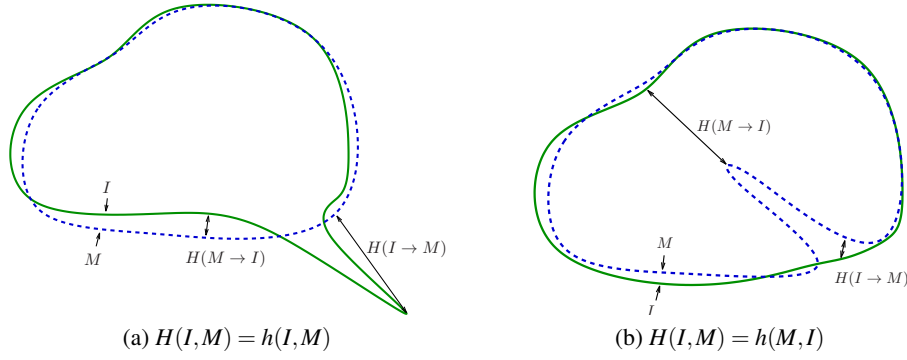


Figure 3: An illustration of the Hausdorff distance.

The mesh has to provide a close approximation of the object shape. We measure the closeness by the fidelity tolerance, quantified by the two-sided Hausdorff distance from the mesh to the image and the image to the mesh. In mathematics, the Hausdorff distance measures how far two subsets of a metric space are from each other. For image boundary I and mesh boundary M , the one-sided distance from I to M is given by

$$h(I, M) = \max_{i \in I} \min_{m \in M} d(i, m),$$

where $d(\cdot, \cdot)$ is the regular Euclidean distance. The one-sided distance from M to I is given similarly by

$$h(M, I) = \max_{m \in M} \min_{i \in I} d(m, i).$$

The two-sided distance is:

$$H(I, M) = \max\{h(I, M), h(M, I)\}.$$

To measure the Hausdorff distance from the surface mesh to the boundaries of the image, we compute the Euclidean distance transform (Maurer, Qi, and Raghavan 2003) (EDT) of the extended image (same size as the octree), and split the octree until no leaf has the distance of EDT both larger and within the input fidelity bound. We mark the leaves that are within the input fidelity tolerance. If one of the triangular facets of the surface mesh in the leaf intersects at least one of the leaves that marked as outside the fidelity tolerance, the fidelity condition is violated and the leaf is split. To measure the Hausdorff distance from the boundaries of the image to the corresponding surface mesh, we compute the shortest distance from each image boundary vertex in the leaf to the triangular patches of the surface mesh in the leaf. If one of the image boundary vertices has a shortest distance larger than the fidelity tolerance, the fidelity condition is violated and the leaf is split.

2.1.3 Filling in the octree

When the waterproof surface mesh respects the user specified fidelity bounds, the octree is constructed. We fill the octree leaves with high quality tetrahedra using the pre-defined volume look-up table. The volume look-up table is designed based on the same idea as the surface look-up table, however, instead of being used to generate edges on cube faces, it is used to generate triangles on cube faces. The template triangles of the volume look-up table should respect the intersection edges of the surface look-up table, which means that the intersection edges should be the edges of the template triangles. By analyzing all possible shapes of the initial tetrahedra filling a cubic leaf of the octree, the minimum dihedral angle bound is 19.47° .

2.2 Mesh decimation

Similar to the LD method (Chernikov and Chrisochoides 2011), the vertex removal operation is used to coarsen the mesh. We maintain a queue of mesh vertices that are candidates for merging. A vertex is merged to a destination vertex if the vertex and the edge between the vertex and its destination are removed from the mesh. As a consequence, all tetrahedra (triangles) incident upon the removed edge are also removed from the mesh. The detailed operation can be consulted in the paper (Chernikov and Chrisochoides 2011), here we only discuss the merging conditions. A vertex can not be merged along an edge to another vertex if it violates the following requirements:

1. The quality requirement, i.e., if at least one of the newly created elements, as a result of a sequence of merges, is inverted or its dihedral angle is smaller than the input quality angle bound, the merge is discarded.
2. The fidelity requirement, i.e., if at least one of the newly created mesh boundary facets has at least one-sided Hausdorff distance larger than the input fidelity bound, the merge is discarded. Same as the fidelity check of the octree construction, this check also consists of two parts, for each of the one-sided Hausdorff distances.
3. The topological equivalence requirement, i.e., if the homeomorphism is not maintained during an operation of merge, the merge is discarded. We apply the following rules to maintain the original structure of the inter-material boundaries: (1) boundary vertices only merge to boundary vertices, (2) a vertex cannot merge to a non-boundary vertex of a different material, (3) for each boundary vertex we also maintain a cumulative list of the octree leaves it belongs to. The merge happens if the union of the set of octree leaves of the merged vertex and the set of the octree leaves of the destination respect the single manifold condition, i.e., the image boundary edges (or faces) of the union form only one $(n - 1)$ -manifold with boundary, and (4) each tetrahedron keeps the original color even after it changes shape due to vertex merge.

3 EXPERIMENTAL RESULTS

We apply the proposed octree refinement and decimation algorithm (ORD) to both synthetic and real medical data in the following sections. All the experiments were conducted on a 64 bit machine equipped with two 3.06 GHz 6-Core Intel Xeon CPU and 64 GB main memory. The algorithm was implemented in C++, in both two and three dimensions.

For the 3D visualization of the final meshes, we used ParaView (Ahrens et al. 2005), an open source visualization software. In Figure 4, we compare the final mesh generated by a Delaunay open source mesh generator Computational Geometry Algorithms Library (CGAL) and our ORD method on topology on *head neck* image. The size of the *head neck* is $255 \times 255 \times 229$ voxels. Each voxel has side lengths of 0.97, 0.97, and 1.4 units in x, y, and z directions. It is clear that there are two manifolds in the original image (see Figure 4a), however, CGAL only generated one manifold (see Figure 4b). We highlight the missing manifold of CGAL mesh in green in Figure 4a. On the contrary, our method generated two manifolds (see Figure 4c).

Figure 5 shows the final mesh produced on *abdomen* atlas and its corresponding cut view. The size of the *abdomen* atlas is $512 \times 512 \times 219$ voxels, each voxel has side lengths of 0.96, 0.96, and 2.4 units in x, y, and z directions. In Table 1 we show the comparison of the output mesh size of the ORD and LD for the *abdomen* atlas. We fixed the two-sided Hausdorff distance bound parameters, and vary the dihedral angle bound to the interested value 5° , 10° , 15° , and 19.47° spread through the range of its feasible values (0° to 19.47°). As we can see from Table 1, the output mesh size is low when $H(I, M)$ is high for all configurations. Also, the final number of tetrahedra decreases as the dihedral angle bound decreases. The tetrahedra generated by ORD before decimation is much fewer than the tetrahedra generated by LD before decimation. When

$H(I, M) = 1$, the sizes of ORD meshes are slightly smaller than the sizes of LD meshes. However, when $H(I, M) = 2$, $H(I, M) = 3$ and $H(I, M) = 4$ the ORD algorithm generated significant smaller number of tetrahedra compared to the number of tetrahedra generated by the LD algorithm. We conclude that the ORD algorithm performs better in terms of the number of tetrahedra than LD algorithm even though it maintains the topology.

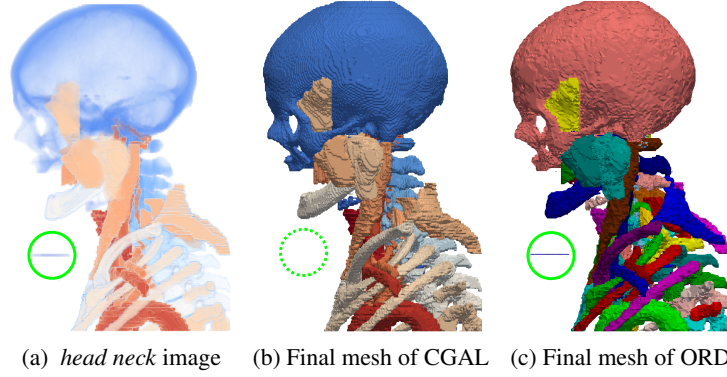


Figure 4: The topology comparison of the ORD mesh and the CGAL mesh on *head neck*.

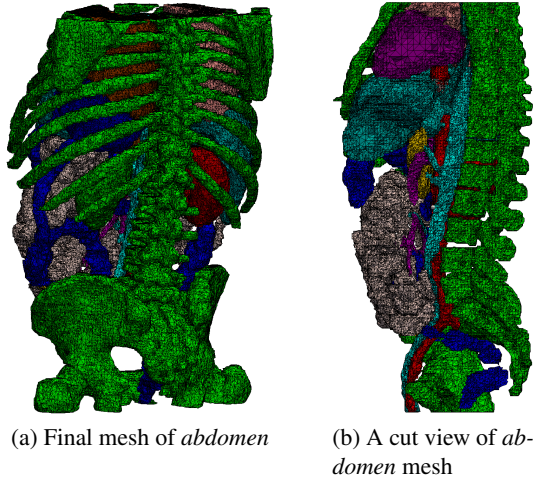


Figure 5: The ORD mesh and a cut view of the *abdomen* for $h(I, M) = 2$, $h(M, I) = 2$ and $\theta = 19.47^\circ$.

Table 1: The comparison of final number of tetrahedra for ORD and LD

Hausdorff distance	$h(I, M) = 1, h(M, I) = 1$		$h(I, M) = 2, h(M, I) = 2$		$h(I, M) = 3, h(M, I) = 3$		$h(I, M) = 4, h(M, I) = 4$	
Algorithm	ORD	LD	ORD	LD	ORD	LD	ORD	LD
Before decimation	30,361,224	42,344,304	18,492,773	51,284,042	19,637,745	57,723,164	19,968,161	61,672,648
After decimation with $\theta = 19.47^\circ$	13,027,911	17,608,762	4,062,233	12,434,028	4,164,029	13,672,489	4,243,052	13,238,118
After decimation with $\theta = 15.00^\circ$	10,275,628	11,994,099	1,593,965	3,426,170	1,560,986	3,035,179	1,628,551	2,940,484
After decimation with $\theta = 10.00^\circ$	9,319,252	10,646,685	856,455	2,081,046	825,285	1,592,756	885,394	1,383,768
After decimation with $\theta = 5.00^\circ$	8,510,971	9,488,867	580,192	1,644,226	559,969	1,188,149	604,498	1,079,160

We also conducted an experiment using CGAL and compare the performance with the performance of our ORD algorithm. Table 2 presents the experimental evaluation of the I2M conversion functionality

make_mesh_3 offered by CGAL and by ORD algorithm. We vary the value of parameter *facet_distance*, and show $h(M, I)$ and $h(I, M)$. They are measured by resampled voxel unit. We also list the final number of tetrahedra produced by CGAL, the minimum dihedral angle (measured by degree) and the total running time (measured by seconds). In some cases, when the value of parameter *facet_distance* was increased, the value of $h(M, I)$ also was increased, however, there is no obvious relationship between *facet_distance* and $h(M, I)$. Further more, the values of $h(I, M)$ in all the cases are unreasonably large. We conclude that CGAL can only approximate one-sided Hausdorff distance, while the proposed algorithm can always guarantee that the Hausdorff distance bound is two-sided. CGAL improves mesh properties such as the dihedral angles using various combinations of optimization algorithms, however, we could only obtain the best minimum dihedral angles 5° . On the contrary, our best minimum dihedral angle bound is 19.47° .

Table 2: The comparison of final number of tetrahedra and run time for ORD and CGAL on *abdomen*

<i>abdomen</i>											
Opt.		<i>no_lloyd()</i> , <i>no_odt()</i> , <i>perturb()</i> , <i>exude()</i>					<i>lloyd()</i> , <i>no_odt()</i> , <i>perturb()</i> , <i>exude()</i>				
Algorithm	<i>facet_dist.</i>	$h(M, I)$	$h(I, M)$	dih. angle	# of tets	Total time	$h(M, I)$	$h(I, M)$	dih. angle	# of tets	Total time
CGAL	0.2	3	23	1.49	10,349,057	532.06	6	18	1.61	9,979,982	4558.32
ORD	N/A	3	23	1.49	1,165,177	1024.49	6	18	1.61	1,160,524	1233.29
CGAL	0.4	3	23	2.26	2,042,684	95.80	4	23	1.74	2,000,988	762.00
ORD	N/A	3	23	2.26	1,180,896	1024.59	4	23	1.74	1,163,455	1022.03
CGAL	0.6	4	23	2.35	862,668	39.23	4	24	2.66	849,445	299.09
ORD	N/A	4	23	2.35	1,176,734	1021.54	4	24	2.66	1,193,001	1185.63
CGAL	0.8	4	23	3.05	487,277	27.99	4	24	3.20	481,107	163.09
ORD	N/A	4	23	3.05	1,201,106	1022.16	4	24	3.20	1,207,731	1184.07
Opt.		<i>no_lloyd()</i> , <i>odt()</i> , <i>perturb()</i> , <i>exude()</i>					<i>lloyd()</i> , <i>odt()</i> , <i>perturb()</i> , <i>exude()</i>				
Algorithm	<i>facet_dist.</i>	$h(M, I)$	$h(I, M)$	dih. angle	# of tets	Total time	$h(M, I)$	$h(I, M)$	dih. angle	# of tets	Total time
CGAL	0.2	5	24	1.01	10,178,967	2132.82	5	18	2.01	9,985,320	5079.16
ORD	N/A	5	24	1.01	1,149,725	1179.20	5	18	2.01	1,173,067	1233.10
CGAL	0.4	7	18	0.86	2,031,462	379.13	5	18	2.17	2,002,825	800.18
ORD	N/A	7	18	0.86	1,144,908	1231.44	5	18	2.17	1,172,660	1229.75
CGAL	0.6	6	19	2.04	862,156	152.53	6	19	2.07	851,196	354.28
ORD	N/A	6	19	2.04	1,164,331	1025.44	6	19	2.07	1,170,672	1026.11
CGAL	0.8	8	19	3.05	487,879	85.16	8	23	5.01	481,937	178.30
ORD	N/A	8	19	3.05	1,195,433	1033.67	8	23	5.01	1,257,769	1022.35

4 CONCLUSION

We presented a novel approach for automatic construction of two- and three-dimensional unstructured meshes of multi-material images characterized by (i) guaranteed dihedral angle bound for the output tetrahedra, (ii) guaranteed bounds on two-sided Hausdorff distance between the boundaries of the mesh and the boundaries of the materials, (iii) the mesh boundary is proved to be homeomorphic to the object surface, and (iv) a small number of mesh elements. The applications of the algorithm include mechanical modeling for finite element simulation, computational medicine and computational biology such as medical imaging, image registration, surgical simulation, and image-guided intervention.

5 ACKNOWLEDGMENTS

This work was supported (in part) by the Modeling and Simulation Graduate Research Fellowship Program at the Old Dominion University and NSF grant CCF-1439079. We thank the anonymous reviewers for helpful comments.

REFERENCES

- Ahrens, J., B. Geveci, and C. Law. 2005, 01. “ParaView: An End-User Tool for Large Data Visualization.”
- Amenta, N., S. Choi, T. K. Dey, and N. Leekha. 2000. “A Simple Algorithm for Homeomorphic Surface Reconstruction”. In *Proceedings of the Sixteenth Annual Symposium on Computational Geometry*, SCG '00, pp. 213–222.
- Amenta, N., T. J. Peters, and A. C. Russell. 2003. “Computational topology: ambient isotopic approximation of 2-manifolds”. *Theoretical Computer Science* vol. 305 (1), pp. 3–15.
- Boissonnat, J.-D., and S. Y. Oudot. 2005. “Provably good sampling and meshing of surfaces”. *Graphical Models* vol. 67, pp. 405–451.
- Bronson, J., J. A. Levine, and R. Whitaker. 2014, Feb. “Lattice Cleaving: A Multimaterial Tetrahedral Meshing Algorithm with Guarantees”. *IEEE Transactions on Visualization and Computer Graphics* vol. 20 (2), pp. 223–237.
- CGAL. “CGAL, Computational Geometry Algorithms Library”. <http://www.cgal.org>.
- Cheng, S.-W., T. K. Dey, and E. A. Ramos. 2010. “Delaunay Refinement for Piecewise Smooth Complexes”. *Discrete & Computational Geometry* vol. 43 (1), pp. 121–166.
- Chernikov, A., and N. Chrisochoides. 2011. “Multitissue tetrahedral image-to-mesh conversion with guaranteed quality and fidelity”. *SIAM Journal on Scientific Computing* vol. 33, pp. 3491–3508.
- Foteinos, P., A. Chernikov, and N. Chrisochoides. 2014. “Guaranteed quality tetrahedral Delaunay meshing for medical images”. *Computational Geometry Theory and Applications* vol. 47, pp. 539–562.
- Labelle, F., and J. R. Shewchuk. 2007. “Isosurface Stuffing: Fast Tetrahedral Meshes with Good Dihedral Angles”. *ACM Transactions on Graphics* vol. 26 (3), pp. 57.1–57.10. Special issue on Proceedings of SIGGRAPH 2007.
- Lorensen, W. E., and H. E. Cline. 1987. “Marching cubes: A high resolution 3D surface construction algorithm”. *COMPUTER GRAPHICS* vol. 21 (4), pp. 163–169.
- Maurer, Jr., C. R., R. Qi, and V. Raghavan. 2003, February. “A Linear Time Algorithm for Computing Exact Euclidean Distance Transforms of Binary Images in Arbitrary Dimensions”. *IEEE Trans. Pattern Anal. Mach. Intell.* vol. 25 (2), pp. 265–270.

MUSCULOSKELETAL SIMULATION FOR GERIATRIC APPLICATIONS BASED ON OPENSIM/SIMTK

Md Shariful Islam, Dr. Michel Audette
Modeling, Simulation & Visualization Engineering
1300 Engineering & Computational Sciences Building

E-mail address
misla003@odu.edu, maudette@odu.edu

ABSTRACT

The development of Healthcare technology has provided us the opportunity to minimize the age old problem of musculoskeletal disorder i.e. falling of elderly persons while walking. A musculoskeletal system can be simulated by the registration of full body MRI images. 3D models of a human body, which can be simulated by structure sensor, can be fitted into the above musculoskeletal model so that we can potentially predict the musculoskeletal impairments in elderly persons to plan the effective surgical and rehabilitation treatments.

Keywords: Musculoskeletal impairments, rehabilitation, registration, structure sensor.

1 INTRODUCTION

There is a growing interest in the diagnosis, the surgical and rehabilitation treatments. These disorders are very common and costly, especially for elderly persons. These unnecessary economic and human costs of musculoskeletal disorders can be minimized by using musculoskeletal modeling and simulation. The patient specific modeling and simulation of musculoskeletal system can elucidate the cause and effect relationships in elderly persons with musculoskeletal impairments.

2 METHODS

We have a plan to use MRI images from BodyParts3D which segments of a 3D whole-body model for an healthy adult human male to simulate and visualize the musculoskeletal that replicate human musculoskeletal system. An interface will be built that encapsulates the main parts of a normal human body by registering data from full body MRI images. At first, the skeletal system will be built by attaching and organizing appropriate bones on different planes that separate the skeletal system in six regions i.e. Head, left hand, right hand, left leg, right leg and torso. Then the muscles will be attached to the bones according to their relative positions matching with human musculoskeletal system. Then OpenSIM/SIMTK will be used to simulate musculoskeletal system that can walk or run like a human being. Then the Structure Sensor will be used to create high fidelity 3D models of senior subjects with high resolution texture. Then, this 3D range surface will be fitted with the aforementioned musculoskeletal model to predict the musculoskeletal disorder in a senior subject. From this two models, body mass and motion of the senior subjects will be compared and their imminent fall events can be predicted.

3 RESULTS

We expect to encapsulate the main parts of the musculoskeletal system with both bones and muscles attached to it. We expect to get outputs like the following sample figures:

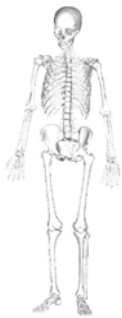


Figure 1: Simulated musculoskeletal (bones only) system from MRI image.



Figure 2: Musculoskeletal system with muscles.

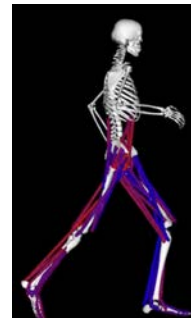


Figure 3: Musculoskeletal system that can move

The following figures (4 and 5) shows the original image and created 3D range surface model using structure sensor.



Figure 4: Original image



Figure 5: Created 3D range surface

4 CONCLUSION AND DISCUSSION

Biological complexity is also a major challenge to simulate the musculoskeletal model. For simplicity, we are only considering the main muscles of a human body. Despite all these, we hope that musculoskeletal simulations will maximize treatment efficacy, limit undesired consequences and reduce cost. Because, prediction of fall events of the senior subjects will minimize or avoid injuries which may lead to financial and human cost.

TOWARDS DEFORMABLE CRANIUM & FORAMEN SURFACE MODEL

Shrabani Ghosh, Dr. Michel Audette
Modeling, Simulation & Visualization Engineering
Old Dominion University
1300 Engineering & Computational Sciences Building
Norfolk, VA 23529
Email: {sghos003, maudette }@odu.edu

ABSTRACT

Recent development in segmentation on medical imaging has advanced medical diagnosis in a broader way. In this paper, we are representing cranial nerve segmentation technique that is combination of surface contouring and deformable surface model. At first, we have already created a tissue blob by stacking 2D images which will be used for contouring using marching cubes in the next step and at last we will end up with deformable surface model using 2 simplex.

Keywords: Cranial nerve, Contour model, Marching cubes, Simplex Mesh.

1 INTRODUCTION

Human skull has twelve pairs of cranial nerve which controls our sensory system. There are numerous holes in which many nerves exit the skull. It's hard to differentiate the holes and nerves. It is important to identify cranium and foramina in CT for neurosurgical planning and operations.

2 METHODS

Our whole procedure is divided into three parts. The first part is to create a volumetric tissue blob unifying the CT image volume of axials, sagittal and coronal. We have generated 3D image volume of axials, sagittal and coronal by stacking 2D images. Now, we need to unify these three-image volumes. Our second step will be the contouring of the 3D blob using marching cubes. In the last step, the contour will be incorporated into 2-simplex which will generate the deformable surface model.

3 RESULTS

So far, we have generated 3D stacked images of axials, sagittal and coronal from skull based CT 2D images. Next I will unify these three results all together. That unified isotropic image will be used for contouring using marching cubes and the contour model will be incorporated into 2 simplex.

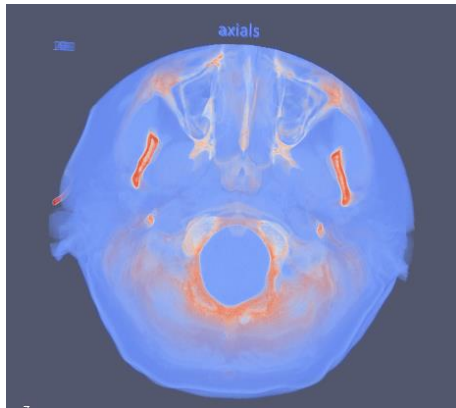


Figure 1 :3D axials image volume

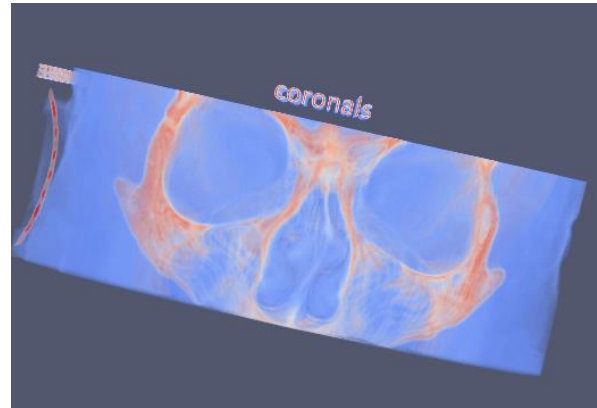


Figure 2 :3D coronals image volume

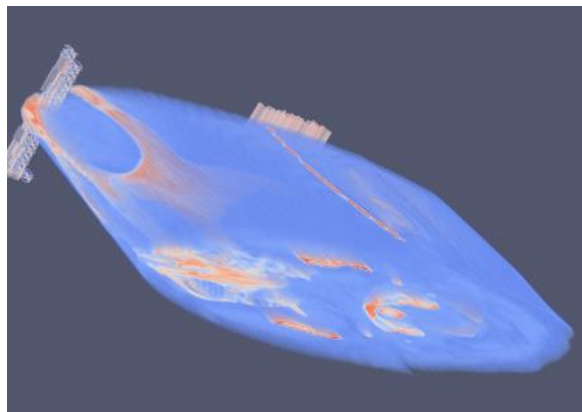


Figure 3 :3D sagittal image volume

4 CONCLUSION AND DISCUSSION

Our technique will help us to differentiate the cranial nerves and holes in the skull which will be helpful for surgical operations. Our next plan is to improve the performance using super resolution to identify the cranial and foramens.

SCOLIOSIS SURGERY PLANNING THROUGH CADAVERIC LIGAMENTO-SKELETAL TISSUE MAPPING AND LOADING STUDIES, MULTI-SURFACE SEGMENTATION, AND FINITE ELEMENT SIMULATION OF THE SPINE

Lucas N. Potter

Frank L. Batten College of Engineering and Technology

Biomedical Engineering

5115 Hampton Blvd, Norfolk, VA, USA

Lpott005@odu.edu

ABSTRACT

Currently scoliosis surgeons do not utilize vital information about the corrective forces applied to the spine in the course of surgery. These forces are reliant on patient-specific anatomy, and therefore require specific segmentations of that patient. These segmentations may require multiple modalities and robust segmentation methods. By applying shape-statistic models to create patient-specific segmentations, the forces involved in the corrective surgery could be predicted and thus workflow can be increased.

Keywords: Scoliosis, medical simulation, patient-specific, patient outcomes.

1 INTRODUCTION

Scoliosis effects 6-9 million people in the United States¹. Yet, the exact amount of force applied during this corrective procedure is, as of yet, unknown. This is true for multiple procedures, including the positioning and the removal of vertebral processes and faces, and the fixation of the vertebral bodies. However, with patient-specific segmentations, supported by cadaveric validation studies, this could change- especially by applying mechanical Finite Element (FE) methods. These methods applied to musculoskeletal loading (in the field of spinal medicine) have a strong emphasis on the mechanics behind pedicle screws^{2,3}. Therefore, patient specific anatomy is not usually part of these studies- as it would add an additional level of complexity. Furthermore, these studies do not directly apply to patient outcomes, but instead inform the manufactures of devices and surgical operators.

2 PROJECT AIMS

The end goal of this project is to create a foundational scoliosis surgery planning and simulation system. With such a relatively large end goal, the emphases will be on three primary objectives. 1) A minimally supervised segmentation of vertebrae, intervertebral discs, and ligaments, 2) A estimated of the force applied during corrective surgery through Finite Element (FE) methods, and 3) Cadaveric studies (as a means of validation).

The difficulties associated with the segmentation of spinal structures are two-fold. Firstly, the shape statistics model required for adequately segmenting the vertebrae and intervertebral discs are only available for the lumbar region⁴. Secondly, the ligaments, which contribute to the mechanical properties of the entire vertebrae, do not show up well either with CT (computed tomography) scanning, since X-rays do not attenuate well with soft tissues such as ligaments, or with MRI (Magnetic Resonance Imaging) because ligaments have little contrast with respect to nearby tissues. Additionally, some spinal structures (such as

hemi-vertebrae) are so rare that they may not have enough examples of their structure to form an adequate Shape Statistic Model (SSM). All of these factors will be taken into account while designing the segmentation.

The forces applied during these procedures will be acquired by instrumentation of surgical instruments, and the forces from before and after can be estimated by imaging of the patient and comparing the maximum deflection (also known as a bend test) to the computer mode of the spine and the associated ligaments.

To get around this problem, cadaveric studies will be done in order to generate shape-statistic models that will hopefully make the segmentations easier. This will use cadavers with radio-detectable thread within the ligament to more easily identify it, and to create a generalizable model for shape-statistics. This will allow for a kind of ground truth to be applied to the generated segmentations.

Once these cadaveric studies are done, we can proceed with FE studies to better analyze the forces present both on the scoliotic spine tissues normally and the forces involved during surgery. These FE studies will in turn be validated by cadaveric studies on the spine on a testing rig capable of applying forces on the spine both laterally and transversely.

3 CONCLUSION

The aims of this project are to generate a semi-automatic segmentation for the human spine to be used in scoliosis surgeries, and to quantify and predict the forces present on different parts of the skeletal anatomy before, during, and after corrective surgery. It is also a springboard towards further predictive simulations and can be used in the future as a start of haptics-driven interactive simulation. The main goal of the author's is to create a segmentation using shape statistics for the structures of the spine, including vertebrae, intervertebral discs, and spinal ligaments. Further goals, such as the FE testing of the scoliotic spinal column and the cadaveric validation of that testing will be done as part of a larger research team.

REFERENCES

1. *Scoliosis*, May 2016, American Association of Neurological Surgeons. Retrieved from: <http://www.aans.org/Patient%20Information/Conditions%20and%20Treatments/Scoliosis.aspx>
2. *The biomechanics of pedicle screw-based instrumentation*. Cho W, Cho SK, Wu C. s.l. : J Bone Joint Surg Br., 2010 Aug; Vols. 92(8):1061-5. doi: 10.1302/0301-620X.92B8.24237.
3. *Pedicle Screw Fixation Under Non-Axial Loads: a Cadaveric Study*. Bianco RJ, Aubin CE, Mac-Thiong JM, Wagnac E, Eng P, Arnoux PJ. s.l. : Spine (Phila Pa 1976), 2015 Oct 15., Vol. [Epub ahead of print].
4. Haq, R., (2015). *MULTI-SURFACE SIMPLEX SPINE SEGMENTATION*

PRACTICAL APPLICATIONS OF NEUROSURGICAL ONTOLOGIES FOR VARIOUS CRANIOTOMIC APPROACHES THROUGH COMPUTER ASSISTED SURGERY

Austin Tapp, The College of William and Mary

Michel Audette, Ph.D, Old Dominion University

Introduction

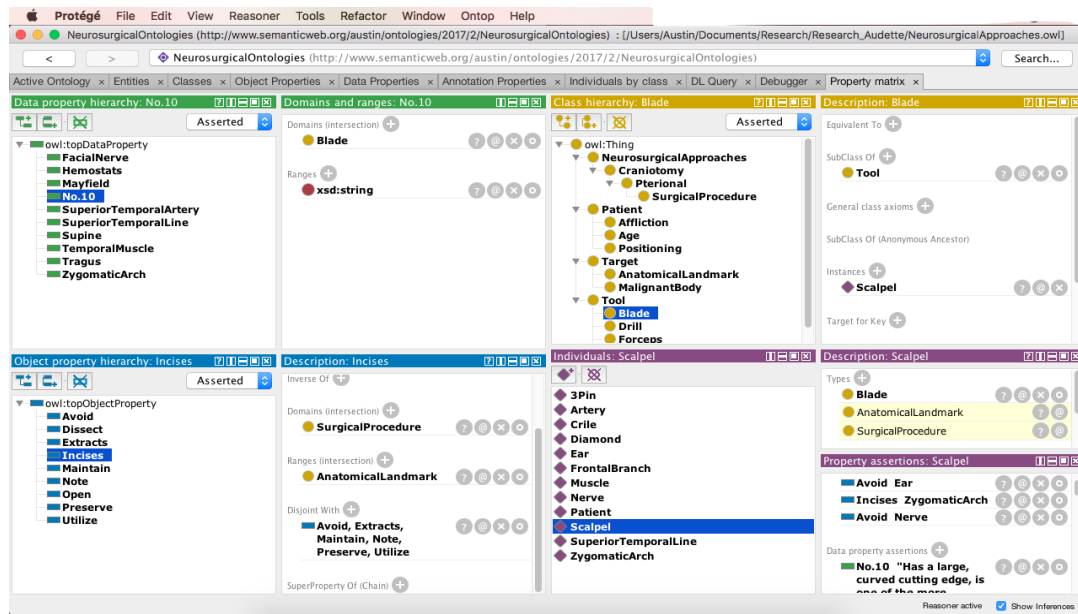
The development of multimodal imaging devices, patient specific data sets and well-defined surgical procedures now provide intense influxes of information that may be processed and utilized in Computer Assisted Surgery (CAS). As CAS becomes increasingly common for a variety of surgical fields, there is a need for the representation, storage, and processing of surgical knowledge in a structured, integrated and standardized format, which can easily adapt to match case specific needs. This format, represented as an ontological model, is characterized by explicit descriptions of concepts in a domain of discourse, with the properties of each concept described by various features and attributes of the ontology's classes, which are organized as a hierarchy. Ontologies have real time operating room applications that display discrete and well-defined steps to guide surgical actions based on patient information, imaging or otherwise, obtained during the surgical procedure. In an operating room, surgical ontologies and models are obtained by witnessing a set of interventions, noting each action by computer, and fusing the descriptions together to provide the ideal surgical approach for a specific patient. The aim of this study is to build a variety of coarsely defined craniotomic neurosurgical approach ontologies that will be used to assist surgeons in the operating room while training, utilizing integrative imaging and/or considering best practices. The scope of the ontologies is limited to craniotomic approaches for the removal of malignant bodies but may extend far beyond single linear surgical approaches. These ontologies will strive to integrate other current ontologies, including those for maxillofacial and orthopedic surgery, pre/post-op patient care, and clinical examination and diagnosis. An essential complement to these ontologies will be their compatibility with modern imaging analysis and medical modeling software that provides live feedback and utilizes algorithms to suggest surgical approaches with the highest rates of success based on the information obtainable and analyzed in real time.

Methods

It is critical to understand the way in which ontologies interact with other computer software and imaging hardware. Particularly, learning about the language of open web ontologies (OWL) files, which are the most common type of ontology file created, is extremely necessary. These ontologies were developed based on the neurosurgical approach data obtained from 3D – Neuroanatomy. The data mined here was then cross-referenced and confirmed for terminological and anatomical accuracy on sites such as NeuroLex, SNOMED, OpenGALEN, and DBpedia. The ontologies themselves are created through a program known as Protégé, developed by Stanford. Within Protégé, there are a variety of ontological “reasoners” that further assist in the confirmation of accuracy of the ontologies as well as reveal their inconsistencies. Once the ontologies are completed, they will be cultured and ratified by neurosurgeons and other medical practitioners to ensure the ontology depicts the highest current standards of medical care and surgical practices.

Results and Conclusion

Currently, the ontologies are still in development. The figure below represents an example of current ontological production. The approaches themselves are condensed into a single ontological map, which will allow a communication of sorts between each of the approaches, particularly in regard to anatomical landmarks and appropriate interventions.



Discussion

In the context of artificial intelligence and CAS, ontologies may provide a specific form of conceptualization that represents formal descriptions of concepts and the relationship between particular domains of interest. Semantically mapping neurosurgical approaches that can be interpreted by both humans and machines opens a vast avenue of future real time data analysis, integration, and processing of intra-operative medical images, which would greatly expand the ability for surgeons to provide ideal care for any patient. Every step of any surgical procedure could be explained with the annotated medical images of each patient, and provide the surgeon with additional information, allowing them to elect for the best approach.

References

- Mudunuri, R., Burgert, O., & Neumuth, T. (2009). Ontological Modelling of Surgical Knowledge. In *GI Jahrestagung*. Retrieved from <http://subs.emis.de/LNI/Proceedings/Proceedings154/gi-proc-154-61.pdf>
- Noy, N. F., & McGuinness, D. L. (2001). *Ontology Development 101: A Guide to Creating Your First Ontology*. Stanford Knowledge Systems Laboratory, 25. <https://doi.org/10.1016/j.artmed.2004.01.014>
- Perrone, R. (2012). *Ontological modeling for Neurosurgery : application to automatic classification of temporal and extratemporal lobe epilepsies*.
- Stanford Center for Biomedical Informatics Research. *Protégé*. Stanford University, 2003. 24 Mar. 2017. <http://protege.stanford.edu/>
- 3D-Neuroanatomy. Alicante, Spain. 2016. 24 Mar. 2017. <http://3dneuroanatomy.com/>

STUDENT CAPSTONE CONFERENCE

2017

GENERAL SCIENCES & ENGINEERING

- Page 129 Nathan Y. Li, Debrup Banerjee and Jiang Li
Princess Anne High School and Old Dominion University
A Comparative Study of Classification Schemes in Transfer Learning for PTSD Diagnosis
- Page 141 Christos Tsolakis, Andrey N. Chernikov and Nikos P. Chrisochoides
Old Dominion University
Parallel Constrained Delaunay Meshing Algorithm in Three Dimensions
- Page 143 Gary Lawson and Robert Baurle
Old Dominion University and NASA Langley Research Center
Enhancing Application Performance Using Mini-apps: Comparison Of Hybrid Parallel Programming Paradigms
- Page 150 Gaya Gnanalingam, Mark J Butler and Holly Gaff
Old Dominion University
Protecting the Big Ones: Harvest Slot Limits and Marine Protected Areas for the Management of Caribbean Spiny Lobster
- Page 152 Beau H. Branch, Samantha C. Collins, Lee C. Dumaliang, Nathan D. Gonda, Timothy P. Lane, Kari Miles, Melissa Periman and Dominic A. Scerbo
Old Dominion University
USV Simulation in the Rapid USV Model Prototyping System
- Page 163 Anthony Williams and Ruhai Zhou
Old Dominion University
Instability and Patterns of Active Suspensions of Liquid Crystals
- Page 172 Sai Dangan and Yuzhong Shen
Old Dominion University
A Simulation and Visualization Approach of Air Pollution in China
- Page 174 Evan Coleman and Masha Sosonkina
Naval Surface Warfare Center and Old Dominion University
Fault Tolerance for Fine-Grained Parallel Iterative Methods
- Page 186 Aaron Walden
Old Dominion University
A Vector Intrinsic Point-implicit Linear Solver for Unstructured Grid Applications on Intel Xeon Phi Knights Landing
- Page 188 Erik Jensen and Masha Sosonkina
Old Dominion University
Modeling Task Migration for Fault Tolerance in Matrix-Matrix Multiplication
- Page 190 Tunazzina Islam
Old Dominion University
Protein Secondary Structure Detection Using Pattern Recognition and Modeling
- Page 193 Michelle E. Pizzo and Fang Q. Hu
Old Dominion University
Simulation of Sound Absorption by Scattering Bodies Treated with Acoustic Liners and the Assessment of its High-Performance Parallel Computing Capabilities

A COMPARATIVE STUDY OF CLASSIFICATION SCHEMES IN TRANSFER LEARNING FOR PTSD DIAGNOSIS

Nathan Y. Li
The IB Program
Princess Anne High School
4400 Virginia Beach Boulevard
Virginia Beach, VA
nathanli2014nl@gmail.com

Debrup Banerjee
Department of Electrical And
Computer Engineering
Old Dominion University
Norfolk, VA
dbane001@odu.edu

Jiang Li
Department of Electrical And Computer Engineering
Old Dominion University
Norfolk, VA
JLi@odu.edu

ABSTRACT

In this paper, we investigated two classification models, support vector machines (SVM) and random forests (RF) for PTSD diagnosis through transfer learning. We extracted three categories of features including prosodic, vocal-tract and excitation from speech signals, and utilized a deep belief network (DBN) to diagnose PTSD. The DBN model has a number of parameters and were initialized through transfer learning to mitigate the over-fitting problem. In transfer learning, we first trained the DBN on a large data set, TIMIT, for speech recognition. We then used the DBN as a feature extractor to transfer the three raw speech features to new representations, and finally, we trained SVM and RF classifiers for PTSD diagnosis. We tested the two classifiers on 52 subjects and experimental results show that SVM classifier is more effective than the RF classifier, achieving a best accuracy of 76.92% for PTSD diagnosis.

Keywords: Transfer learning, Random Forests, speech, diagnosis, PTSD

1 INTRODUCTION

Post-traumatic stress disorder (PTSD) is a mental health problem that some people develop after experiencing or witnessing a life-threatening event, like combat, a natural disaster, a car accident, or sexual assault (DVA, 2017). Several factors can increase the chance that someone will have PTSD, of which many are not under that person's control. PTSD can happen to anyone and is a serious problem for the military. Currently at clinics, PTSD diagnosis are normally conducted through structured interviews. However, the success is limited because of embarrassment and social norms as well as memory and self-perception associated with the disease.

Analyzing human speech is an alternative method for the diagnosis of PTSD as it influences people's voices. Moreover, speech is non-invasive and can be relatively easier to obtain through recordings. The effectiveness of treatment can also be monitored through speech. Much effort has been done to study speech features and classification schemes for PTSD diagnosis. Vergri and others (2015) explored three feature categories: lexical, spectral and longer-range prosodic features based on the PTSD recordings from Clinician-Administered PTSD Scale (CAPS) interview, a gold standard in PTSD diagnosis, and an overall accuracy of 77% was achieved. Zhuang et al (2014) conducted research on multi-view learning, combining other signals like electroencephalograms (EEG) with speech and found a net relative increase between 20% and 37% in speech-based PTSD detection. A number of other scientists utilized different data modalities and features for PTSD diagnosis such as functional magnetic resonance imaging (fMRI) by Liu et al. (2014), magnetic resonance imaging (MRI) by Zhang et al. (2016), video and physiology by Brown et al.(2015), event characteristics, emergency department observations and early symptoms by Levy et al. (2014).

Utilizing other modalities may increase the accuracy of PTSD diagnosis but the acquisition of these modalities other than speech is so difficult/expensive that it prevents them from being used for screening in a large population. Speech based diagnosis has many advantages but the shortage of PTSD speech data is a serious challenge for the building of diagnostic models (Banerjee, 2017). In this paper, we presented a deep transfer learning strategy to handle the training data shortage issue with the goal of studying the performances of RF and SVM in the context of transfer learning on the diagnosis of PTSD.

SVM is a discriminative classifier, based on decision planes (Cortes and Vapnik, 1995). A decision plane separates sets of objects that have different classes of data. SVM undertakes classification tasks by finding an optimal hyperplane that gives the largest minimum distance to the training examples or maximizes the margin of the data. It supports both regression and classification tasks and can handle multiple continuous and categorical variables.

RF is a tree-based model (Breiman, 2001). The RF model is a combination of tree predictors where each tree depends on the values of a random vector sampled independently with the same distribution for all trees in the forest. The basic principle is that a group of "weak learners" can come together to form a "strong learner". RF model is an effective tool for making predictions because it usually does not overfit data. The purpose of this paper is to compare the effectiveness of SVM and RF on PTSD diagnosis in the context of deep transfer learning.

DBN is probabilistic generative graphic models that are composed of multiple layers of latent variables or hidden units (Hinton, 2009). The connections between the layers are provided, but no connections exist between units within each layer. For a given set of data, a DBN can learn to probabilistically reconstruct its inputs. After being trained, each layer is regarded as a feature detector for inputs. Through this top-down layer by layer procedure, a DBN can be further trained in a supervised way to perform classification.

Transfer learning was proposed as a learning system with the ability to apply knowledge and skills learned in previous tasks to novel tasks. Transfer learning was designed to address and mitigate the small data problem. A large amount of training data is normally needed by deep learning networks, which can automatically learn features from raw data without prior knowledge and have been shown to perform well (Srivastava et al., 2014; Deng et al., 2013; Dieleman and Schrauwen, 2014; Bengio, 2007). But this requires a large amount of data that is not always available. Obtaining the new training data and rebuilding the models may be prohibitively expensive (Hinton and

Salakhutdinoy, 2006). For this reason, transfer learning becomes a very useful tool to improve learning between task domains. In transfer learning for PTSD diagnosis, we trained a deep learning structure on a large speech reorganization dataset, TIMIT, and transferred the knowledge learned from the data set to diagnose PTSD.

2 METHODOLOGY

2.1 Speech Data Preparation and Feature Extraction

2.1.1 Data Source

PTSD Data Set: The data source for this study is from Youtube and an Ohio hospital. From each of the two sources, we collected 13 normal and 13 PTSD subjects' audio data files, respectively. All these recordings were sampled at the same frequency of 44.1 kHz, and each recording was from only one subject. A majority of these recordings were between 120 and 140 seconds, though the duration could range from 51 seconds to 480 seconds (Banerjee, 2017).

TIMIT Data Set: An acoustic-phonetic speech corpus, TIMIT, was employed for transfer learning. This corpus was developed by Texas Instruments (TI), SRI International (SRI) and Massachusetts Institute of Technology (MIT) to provide speech data for the acquisition of acoustic-phonetic knowledge, the development of mathematic models and verification of automatic speech recognition systems. There are more 900 speakers with 9 different dialects of English speaking styles recorded in the data set.

2.1.2 Raw Speech Features

Three categories of raw features were computed as shown in Table 1. The features were calculated based on a 25ms segments with 10ms overlapping between two consecutive segments. The mathematic models for these features were explained below.

Table 1: Description of speech frame features extracted across all the different speech corpora.

Feature Type	Number Of Features
Prosodic features	
Short-term energy, Average power, Average magnitude, Zero crossings, Mean, Standard deviation, Median, Max, Min, Range, Dynamic range, Interquartile range	12
Vocal-tract features	
MFCC (Mel Frequency Cepstrum Coefficients)	39
Teager Energy Operator	1
Excitation features	
Jitter	1
Shimmer	1
Total number of raw features per frame	54
First order time derivative of raw features	54
Second order time derivative of raw features	54
Total number of features per frame	162

Prosodic features are defined as following:

- Short-term energy is the energy associated with a short-speech segment. It can be used to classify voiced, unvoiced and silence speech.
- Average power is the short-term energy divided by the number of speech samples for a short speech segment.
- Average magnitude function is used for energy contour calculation. Instead of the squares of individual values, their absolute values are summed over shifting short-time window (Rabiner and Schafer, 1978).
- Zero-Crossing is speech signals crossing the zero axes during each frame. Mathematically, zero-crossing occurs if successive samples have different algebraic signs. The rate of zero crossings is a measure of the frequency content of a signal.
- Dynamic range is the difference in the base-10 logarithm of the maximum and the minimum amplitudes of the speech signal frame.
- The interquartile range is the difference between the 75th and 25th percentile.

Vocal-tract features consist of Mel Frequency Cepstrum Coefficients (MFCC) and Teager Energy Operator.

- MFCC is a representation of the short-term power spectrum of a sound, based on a linear cosine transform of a log power spectrum on a nonlinear mel scale of frequency. It is calculated based on frequency bins on the Mel-scale (Patel and Prasad, 2013), which is a frequency binning method based on the human ear's frequency resolution. Based on the frequency that is received and processed, the mel-scale can be utilized to model the functions of human ear.
- Teager's energy operators are defined in both the continuous and discrete domains and are very useful feature for analyzing single component signals from an energy point-of-view. It is a non-linear energy tracking operator. It has shown that airflow separates in the vocal-tract when it propagates, instead of just flowing as a plane wave. Under stress, the vocal system physiology behaves differently, which affects the vortex-flow interactions in the vocal tract (Hansen et al., 2011). This phenomenon has been proved to occur in speech production. This feature can be detected to be responsive to speech under stress using audio from the speech under the simulated and actual stress database (SUSAS) corpus (Hansen et al., 2000), which is a comprehensive speech database to be recorded under stressful conditions and is the database of choice for stressed emotion recognition.

Excitation features include jitter, shimmer, and the time-derivative of the feature vector.

- Jitter is a measurement of vocal stability or frequency perturbation. It is defined as average absolute difference of consecutive pitch periods.
- Shimmer is the same as frequency perturbation, but analogous to amplitude. It is defined as the average absolute difference between amplitudes of consecutive periods.

The time derivative of a feature vector is an important signal processing characteristic used when studying robust features. A set of *delta* functions computes the first-order time derivative of an input feature vector sequence (Center for Spoken Language Understanding, 2017). The second order time derivative is computed based on the first order difference approximation by changing the order parameter of the *delta* function to 2. The combined feature vector is augmented with its first and second order time derivatives resulting in two different feature dimensions.

2.1.3 Speech Features for PTSD Diagnosis

For this study, we chose a long time duration for PTSD diagnosis. We utilized frame sizes of 1, 2 and 3 seconds and a frame shift of 1 second to extract frame-level features for PTSD diagnosis. For each speech frame, we computed the 162 raw features for each of the 25ms segments (there are 10ms overlapping between two consecutive ones) and used the averaged 162 raw features across all the segments as the features for the frame. Then, we concatenated the 162 averaged features from 15 consecutive frames for PTSD diagnosis, resulting in a feature vector with a size of 2430 (162x15). Finally, we moved to the next frame to obtain the next feature vector for the subject.

2.2 Speech Based PTSD Diagnostic Models

The proposed model is displayed in Figure 1. The speech signal was first processed to extract features. The transfer learning was then applied to obtain new feature representations. After this step, the system performed PTSD diagnosis either through RF or SVM.

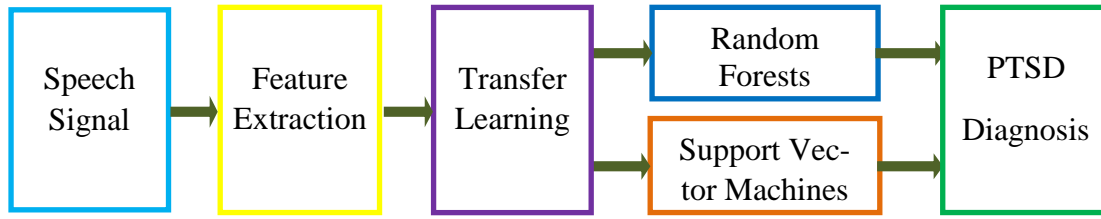


Figure 1. The proposed model for PTSD diagnosis

Transfer learning was designed to focus on transferring knowledge gained from source domain to solve problems in a different but related target domain. It is especially useful when the target task is in shortage of high-quality training data. In the paper, we utilized DBN to diagnose PTSD. DBN are probabilistic generative models that usually have a large number of layers of hidden variables. In order to obtain competitive performances from DBN, it usually needs a large amount of data for training that is not available for our case. We developed a transfer learning strategy to resolve this challenge.

In transfer learning, we first trained the DBN on TIMIT, a large data set for phoneme recognition (there are 39 phonemes in the data set). And then, the Deep belief networks were used as a feature extractor to transfer the speech features to new representations. The new representations were then used to train a SVM classifier or a RF classifier for PTSD diagnosis. The concept of transfer learning is depicted in Figure 2.

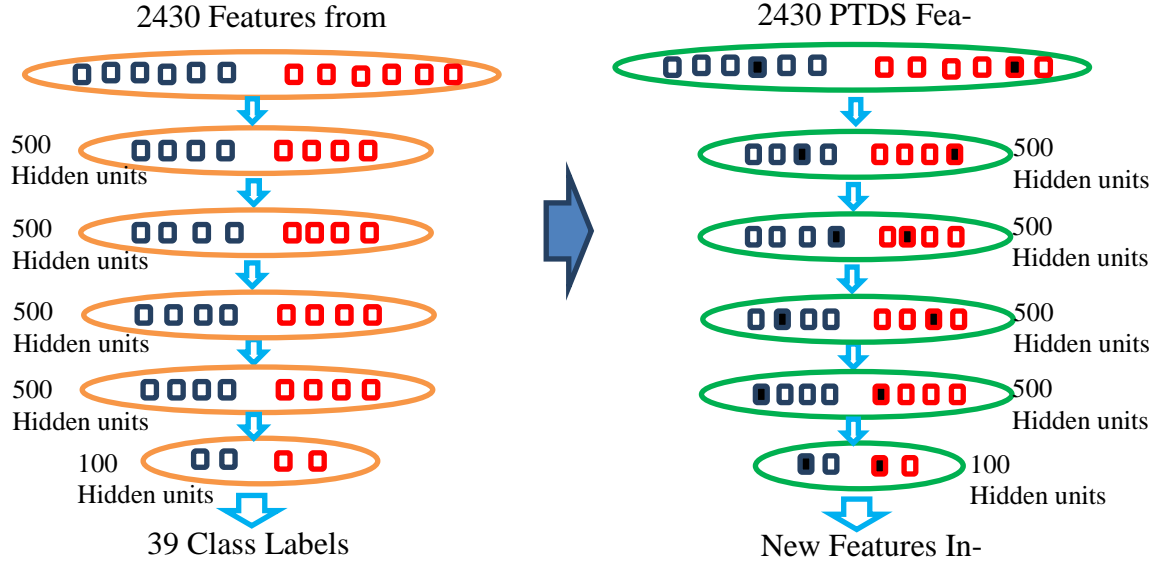


Figure 2: Concept of transfer learning.

3 EXPERIMENTAL SETUP

Experiments were carried out to evaluate the performance of RF and SVM in transfer learning for PTSD diagnosis. Initially, a deep belief network with an architecture of 2430-500-500-500-500-100 was trained on TIMIT data, and then the pre-trained model was transferred to new representations for RF and SVM to conduct PTSD diagnosis.

After the feature transfer learning, RF were trialed with different numbers of decision trees and feature predictors. Under the first hidden layer of 500, 5 groups of decision trees of 50, 100, 200, 500, and 1000 were selected. For each group of trees, 6 group of predictors of 5, 10, 20, 50, 100, and 200 were trialed, resulting in 30 (5x6) sets of data for each time frame and shift under one network architecture. For the Youtube data, a frame of 1.0 second and a shift of 1.0 second was tested. For the Ohio data, a frame of 3.0 seconds and a shift of 1.0 second was tested.

Under the identified configuration for RF, we performed PTSD diagnoses using the new representations extracted from different layers in the 2430-500-500-500-500-100 architecture, using leave-one-subject-out cross validation. The results are shown in the next section.

4 RESULTS AND DISCUSSIONS

4.1 Results of SVM without Transfer Learning

Experiments were first undertaken to evaluate the performance of SVM on PTSD diagnosis without transferring data. Based on the 2430 feature vectors, the direct application of SVM on the multiple frame and multi-category raw feature produced the best diagnostic accuracy of 61.53% on the Youtube data, and 53.84% on the Ohio data (Table 2). These data will be used as a reference to evaluate the effectiveness of the transfer learning strategy.

Table 2. Classification results of SVM directly on 2430 feature vectors by applying leave-one-subject-out cross validation and without applying transfer learning.

Data and Feature			SVM		
Data source	Frame Length Sec	Frame Shift Sec	Subject Wise Accuracy on Youtube, %	Subject Wise Accuracy on Ohio, %	Overall Subject Wise Accuracy %
26 Youtube	3	1	50.00	42.30	46.15
	2	1	61.53	53.84	57.69
26 Ohio	1	1	53.84	53.84	53.84
Average			55.12	49.99	52.56

4.2 Result of RF Configuration Identification

When investigating the optimal trees and predictors as stated in Section 3, the obtained results showed that the configuration of 200 trees and 50 feature predictors performed optimally for both the Youtube and the Ohio data and this pair was then used for further calculation for all data sets.

4.3 Results of Transfer Learning by SVM and RF

In Table 3, when the first hidden layer of 500 features were used for classification, the average subject-wise test accuracy of RF for the Youtube data is 62.82% and for the Ohio data 55.13, and the overall subject-wise accuracy is 58.98%. For SVM, the average subject-wise test accuracy for the Youtube data is 82.04% and for the Ohio data 65.43, and the overall accuracy is 73.71%. It is obvious that SVM performs better than RF here.

Moreover, it seems that in Random Forests, the overall subject-wise accuracy increases with the decreasing frame length from 3 seconds to 1 second, consistent with the results in SVM.

Table 3. RF and SVM classification results using the first hidden layer of 500 features generated after applying transfer learning and applying leave-one-subject-out cross validation.

Data and Feature			RF			SVM		
Data source	Frame Length Sec	Frame Shift Sec	Subject Wise Accuracy on Youtube, %	Subject Wise Accuracy on Ohio, %	Overall Subject Wise Accuracy, %	Subject Wise Accuracy on Youtube, %	Subject Wise Accuracy on Ohio, %	Overall Subject Wise Accuracy, %
26 Youtube	3	1	57.69	53.85	55.77	80.76	61.53	71.14
	2	1	61.54	57.69	59.62	80.76	65.53	73.07
26 Ohio	1	1	69.23	53.85	61.54	84.61	69.23	76.92
Average			62.82	55.13	58.98	82.04	65.43	73.71

Table 4 displays the results when the second hidden layer of 500 features was used. In RF, the average subject-wise test accuracy for the Youtube data is 70.51% and for the Ohio data 57.69%, with the overall subject-wise accuracy is 64.10%. When compared the results in the Table 3, it

can be seen that SVM had a slightly higher accuracy than RF for the Youtube data and on the overall performance, but there was no difference for the Ohio data. Also, the frame length made little difference.

Table 4. RF and SVM classification results using the second hidden layer of 500 features generated after applying transfer learning and applying leave-one-subject-out cross validation.

Data and Feature			RF			SVM		
Data source	Frame Length Sec	Frame Shift Sec	Subject Wise Accuracy on Youtube, %	Subject Wise Accuracy on Ohio, %	Overall Subject Wise Accuracy, %	Subject Wise Accuracy on Youtube, %	Subject Wise Accuracy on Ohio, %	Overall Subject Wise Accuracy, %
26 Youtube 26 Ohio	3	1	69.23	57.69	63.46	73.07	57.69	65.38
	2	1	69.23	57.69	63.46	80.76	57.69	69.22
	1	1	73.08	57.69	65.39	73.07	57.69	65.38
Average			70.51	57.69	64.10	75.63	57.69	66.66

When the third hidden layer of 500 features was used for classification, the results of RF and SVM (Table 5) were similar to those achieved under the second layer. RF had a lower average subject-wise test accuracy for the Youtube data (62.82%) and lower overall accuracy (62.82%) than SVM, but the same accuracy for the Ohio data (62.82%) as SVM. The frame length makes little difference to the accuracy in both RF and SVM.

Table 5. RF and SVM classification results using the third hidden layer of 500 features generated after applying transfer learning and applying leave-one-subject-out cross validation.

Data and Feature			RF			SVM		
Data source	Frame Length Sec	Frame Shift Sec	Subject Wise Accuracy on Youtube, %	Subject Wise Accuracy on Ohio, %	Overall Subject Wise Accuracy, %	Subject Wise Accuracy on Youtube, %	Subject Wise Accuracy on Ohio, %	Overall Subject Wise Accuracy, %
26 Youtube 26 Ohio	3	1	61.54	61.54	61.54	73.07	65.38	69.22
	2	1	61.54	65.38	63.46	76.92	65.38	71.15
	1	1	65.38	61.54	63.46	69.23	57.69	63.46
Average			62.82	62.82	62.82	73.07	62.82	67.94

When the fourth hidden layer of 500 units was used for classification in RF, it is noted from Table 6 that the average subject-wise test accuracy for the Youtube data was identical to that for the Ohio data (67.95%), thus the overall subject accuracy is 67.95%. Interestingly, SVM gave the same diagnostic accuracy to RF for the Youtube data (67.95%), but lower accuracy (61.53%) for the Ohio data. Therefore, The overall performance of SVM is slightly worse than RF under this situation.

Table 6. RF and SVM classification results using the fourth hidden layer of 500 features generated after applying transfer learning and applying leave-one-subject-out cross validation.

Data and Feature			RF			SVM		
Data source	Frame Length Sec	Frame Shift Sec	Subject Wise Accuracy on Youtube, %	Subject Wise Accuracy on Ohio, %	Overall Subject Wise Accuracy, %	Subject Wise Accuracy on Youtube, %	Subject Wise Accuracy on Ohio, %	Overall Subject Wise Accuracy, %
26 Youtube 26 Ohio	3	1	69.23	69.23	69.23	69.23	61.53	65.38
	2	1	65.38	69.23	67.31	69.23	65.38	67.30
	1	1	69.23	65.38	67.31	65.38	57.69	61.53
Average			67.95	67.95	67.95	67.95	61.53	64.74

The results of RF classification using the final hidden layer of 100 are shown in Table 7. The average subject-wise test accuracy for the Youtube data is 66.66% and for the Ohio data is 61.54%, and the overall subject-wise accuracy is 64.10%. Compared with the results of SVM, it is clear that SVM outperformed RF under all scenarios.

Table 7. RF and SVM Classification results using the fifth hidden layer of 100 features generated after applying transfer learning and applying leave-one-subject-out cross validation.

Data and Feature			RF			SVM		
Data source	Frame Length Sec	Frame Shift Sec	Subject Wise Accuracy on Youtube, %	Subject Wise Accuracy on Ohio, %	Overall Subject Wise Accuracy, %	Subject Wise Accuracy on Youtube, %	Subject Wise Accuracy on Ohio, %	Overall Subject Wise Accuracy, %
26 Youtube 26 Ohio	3	1	65.38	57.69	61.54	73.07	65.38	69.22
	2	1	65.38	57.69	61.54	73.07	69.23	71.15
	1	1	69.23	69.23	69.23	73.07	65.38	69.22
Average			66.66	61.54	64.10	73.07	66.66	69.86

Comparing the results in Table 2 with those in Tables 3-7, it can be clearly seen that the diagnostic performance of SVM is much enhanced when the transfer learning model is used. This observation suggests that transfer learning strategy greatly improve the quality of data training.

The accuracy of the Youtube data, the Ohio data and their overall accuracy for both RF and SVM are summarized in Figure 3. From the figure, it can be seen that for RF, the best overall subject-wise accuracy of 70.51% is achieved by the second layer of the network architecture of 2430-500-500-500-100 for the Youtube data, while the best result of 67.95% obtained for the Ohio data is from the fourth layer. For the overall subject-wise accuracy in RF, the fourth layer of 500 gives the best result of 67.95% accuracy. SVM achieved the best accuracy of 82.04% in the first layer for the Youtube data, but for the Ohio data, the best accuracy of 66.66% is obtained in the fifth layer. Comparing the results of the model in RF and SVM, Figure 3 clearly shows that SVM performs better than RF in the overall accuracy, particularly for the first layer. However, when modelled by fourth layer, the RF gives a little better performance.

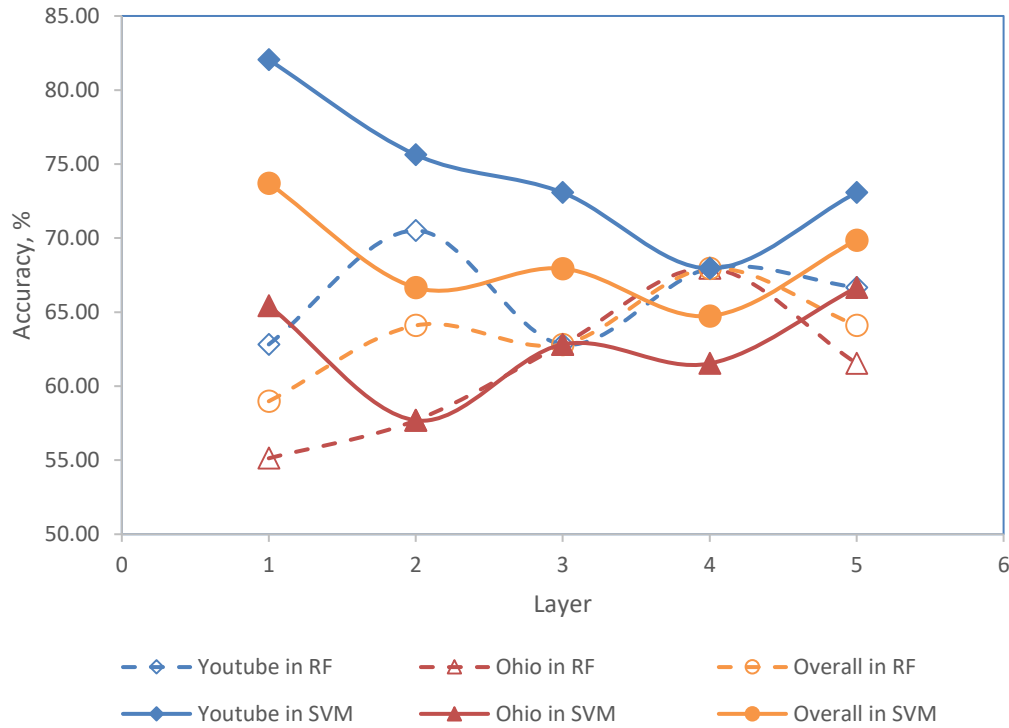


Figure 3. Performance of Random Forests and Support Vector Machines through transfer learning based on the Youtube and the Ohio data sets.

It is interesting to note that SVM performs better than RF on most of the layers, but not on the fourth layer. Practically, both of these classifiers have their advantages and disadvantages, depending on problems, datasets and data distribution. Effort is needed to understand the difference of their performance under this scenario.

The results were obtained by using a configuration of 200 decision trees and 50 feature predictors, under the network architecture of 2430-500-500-500-500-100. Therefore, it can only be said that under this particular scenario, the SVM performs better than RF.

5 CONCLUSIONS

A comparative study was conducted to evaluate the performance of Random Forests and Support Vector Machines classifiers on PTSD diagnosis through transfer learning. TIMIT data was first used to train deep belief networks on raw speech features, and then the pre-trained model was transferred and generated new representative features for RF and SVM to perform PTSD diagnosis. We focused our investigation on multiple frame features and deep networks on 26 subjects each from a Youtube data and an Ohio hospital, and the results showed SVM outperformed RF, as it possessed the highest diagnostic accuracy of 76.92%.

ACKNOWLEDGMENTS

The author would like to thank Dr. Jiang Li for his patience, support and help during the study of this project.

REFERENCES

- Banerjee, D. 2017. "Emotional state recognition and PTSD detection based on speech signals", Thesis, Old Dominion University.
- Bengio, Y. 2007. "Learning deep architectures for AI," *Foundations and trends in Machine Learning*, vol. 2, no 1, pp. 1-127.
- Breiman, L. 2001. "Random Forests", *Machine Learning*, vol. 45, issue 1, pp5-32.
- Brown, S., Webb, A., Mangoubi, R. S., Dy, J. G. 2015. "A Sparse Combined Regression-Classification Formulation for Learning a Physiological Alternative to Clinical Post-Traumatic Stress Disorder Scores", *Proc. of the Twenty-Ninth AAAI Conference on Artificial Intelligence*, pp. 1700-1706.
- Center for Spoken Language Understanding, <http://www.cslu.edu/toolkit/old/old/version2.0a/documentation/csluc/node5.html>. Accessed, Feb. 25, 2017.
- Cortes, C. and Vapnik, V. 1995. "Support Vector Networks," *Machine Learning*, vol. 20, no. 3, pp 273-297.
- Deng, L., Li, J., Huang, J., Yao, K., Yu, D., Seide, F., Seltzer, M. L., Zweig, G., He, X., Williams, J., Gong, Y. and Acero, A. 2013. "Recent Advances in Deep Learning for Speech Research at Microsoft", ICCASP.
- Dieleman, S. and Schrauwen, B. 2014. "End-to-end Learning for Music Audio", ICCASP.
- Hansen, J. H. L., Kim, W., Rahurkar, M., Ruzanski, E. and Myerhoff, J. 2011 "Robust Emotional Stressed Speech Detection Using Weighted Frequency Subbands," Hindawi Publishing Corporation, vol. 2011, no. 10.
- Hansen, J. H. L., Swail, C., South, A., Moore, R. K., Steeneken, H., Cupples, E. J., Anderson, T., Vloeberghs, C. R. A., Trancoso, I., and Verlinde, P. 2000. "The impact of speech under stress on military speech technology," NATO IST/TG-01.
- Hinton, G. (2009). "Deep belief networks". Scholarpedia. 4 (5): 5947. [doi:10.4249/scholarpedia.5947](https://doi.org/10.4249/scholarpedia.5947). Accessed, April 9, 2017.
- Hinton, G. E. and Salakhutdinov, R. R. 2006. "Reducing the Dimensionality of Data with Neural Networks", *Science*, vol. 313, pp. 504-507.
- Jiang, X. 2015 "representation Transfer in Deep Belief Networks", Springer, pp. 338 – 342.
- Levy, I. R. G., Karstoft, K. I., Statkinov, A. and Shalev, A. Y. 2014. "Quantitative Forecasting of PTSD from Early Trauma Responses: A Machine Learning Application", *Journal of Psychiatric Research*, pp. 68-76.
- Liu, F., Xie, B., Wang, Y., Guo, W., Fouche, J. P., Long, Z., Wang, W., Chen, H., Li, M., Duan, X., Zhang, J., Qiu, M., and Chen, H. 2014 "Characterization of Post-traumatic Stress Disorder Using Resting-State fMRI with a Multi-level Parametric Classification Approach", *Brain Topography*, Springer, pp. 221-237.
- Patel, K. and Prasad, R. K. 2013. "Speech recognition and verification using MFCC and VQ," *Int. J. Emerging Science and Engineering*, vol. 1, no. 7, pp. 33-37.

- Rabiner, L. R. and Schafer, R. W. 1978. *Digital Processing of Speech Signals*, Prentice-Hall.
- Srivastava, N., Hinton, G. E., Krizhevsky, A., Sutskever, I., and Salakhutdinov, R. 2014. “Drop-out: A Simple Way to Prevent Neural Networks from Over Fitting”, *Journal of Machine Learning Research*, vol 15.
- US Department of Veterans Affairs, <http://www.ptsd.va.gov/public/PTSD-overview/basics/what-is-ptsd.asp>. Accessed, Feb., 19, 2017.
- Vergyri, D., Knoth, B., Shriberg, E., Mitra, V., McLaren, M., Ferrer, L., Garcia, P., and Marmar, C. 2015. “Speech-based assessment of PTSD in a military population using diverse feature classes”, *Interspeech*, pp. 3729-3733.
- Zhang, Q., Wu, Q., Zu, H., He, L., Huang, H., Zhang J., and Zhang, W. 2016. ” Multimodal MRI-Based Classification of Trauma Survivors with and without Post-Traumatic Stress Disorder“, *Frontiers in Neuroscience*.
- Zhuang, X., Rozgic, V., Crystal, M. and Marx, B. P. 2014. “Improving Speech-Based PTSD Detection via Multi-View Learning”, *IEEE Spoken Language Technology Workshop*, pp. 260-265.

PARALLEL CONSTRAINED DELAUNAY MESHING ALGORITHM IN THREE DIMENSIONS

Christos Tsolakis, Andrey N. Chernikov and Nikos P. Chrisochoides

Center for Real-Time Computing, Department of Computer Science,
Old Dominion University
Norfolk, VA 23529
{ctsolakis,achernik,nikos}@cs.odu.edu

ABSTRACT

Mesh generation is an integral part of the finite element method which in turn is vital for many applications such as bio-medical image-based simulations and computational fluid dynamics. In order to keep the throughput of the simulation pipeline high, mesh generation needs to be scalable and efficient. In this work we present an initial design focused on the correctness of a three dimensional Parallel Constrained Delaunay Meshing Algorithm that uses mesh decomposition techniques to create work units to be executed in parallel and existing sequential meshing software to refine them (ie, maximum possible code re-use). Asynchronous messages are employed to achieve conformity between shared regions of the work units while keeping the overhead low.

Keywords: parallel mesh generation, Delaunay, asynchronous

1 INTRODUCTION

In our previous work (Chernikov and Chrisochoides 2008) we presented an algorithm for parallel constrained Delaunay mesh generation in two dimensions called Parallel Constrained Mesh Generation algorithm (PCDM). PCDM employs two-dimensional domain decomposition techniques to divide the original problem into smaller subdomains. Every subdomain is meshed separately by a different process. Neighboring subdomains exchange asynchronous messages containing information related to the splits of shared segments necessary to achieve conformity on shared interfaces. The asynchronous nature of the messages together with their small size introduce a small overhead to PCDM allowing the method to scale up linearly up to 100 processes.

In this work we present PCDM3D, a design of the three dimensional version of our previous two dimensional algorithm. PCDM3D starts with an existing mesh and refines it in parallel while keeping the triangulation of the shared interfaces the same across all subdomains resulting thus in a mesh which is Constrained Delaunay per subdomain and consistent across the interfaces. Although, PCDM3D shares many ideas with our previous method, it has also to deal with new issues. In particular, (a) the refinement rules for three dimensional geometries are more complex and depend on the features of the input, (b) points can now be inserted both on faces and segments and finally (c) interfaces can be shared by more than two subdomains.

Moving to three dimensions introduces many new challenges. First, domain decomposition of arbitrary three dimensional geometries is not trivial and there are not known methods for the general case. We need thus to provide an alternative input for the mesher. Second, the meshing method that we use (Constrained

Delaunay refinement) is more computationally challenging in three dimensions and can often require more than one point in order to split a single element while respecting the invariant properties of the mesh (Si and Shewchuk 2014). Finally, the concurrent access on the interface triangulation may create non-conforming interfaces since it is known that the order of processing the elements on a triangulation can result in a different triangulation.

2 METHODS

To avoid the complexities introduced by small angles in three dimensional objects we will assume for the rest of the paper that the input satisfies the projection condition as described in (Shewchuk 1997). This condition guarantees that the meshing algorithm will terminate and it will produce a mesh of a certain quality while inserting one point at a time. In practice, to achieve the requirements of the projection condition we will use surface meshes with dihedral angles of 90 degrees obtained from a hexahedral mesh composed entirely by cubes. Having right angles throughout the mesh takes away the need for high quality partitions since any partition will have 90 degrees angles. This enables us to use generic graph partitioning software like METIS (Karypis and Kumar 1998) for decomposing the mesh into subdomains which are then distributed to the processes for sequential refinement using TetGen (Si 2013).

Concurrent access of the shared boundary was not an issue in our previous work because in two dimensions the triangulation of the interface reduces to the discretization of the edges which are split on predefined positions. On the other hand, in three dimensions we need to enforce the same order of processing the elements on the boundary. To achieve this, a token will be assigned to each interface. The token will be unique and only the subdomain that owns it can modify the interface whereas, the rest of the mesh can be refined in parallel with no communication. As soon as a point is inserted on an interface, it is sent together with the token to the next of the neighbors and when the token completes a cycle, the modifications on the interface have been synchronized. The token can also be requested by a subdomain, if it is not present, allowing thus all the subdomains to use it. The uniqueness of the token will enforce sequential modifications on the shared boundaries assuring thus the correctness of the method. While, on the other hand, asynchronous communication will allow to overlap communication with computation.

3 FUTURE WORK

The three dimensional version of PCDM is currently under active development. The loose coupling between the subdomains is expected to allow the algorithm to scale up reasonably well with low communication overhead. Methods that proved useful in the previous implementation such as, message aggregation and over-decomposition will be also evaluated. Moreover, we plan to formalize and prove the correctness of the discussed communication scheme.

4 ACKNOWLEDGMENTS

This work in part is funded by the Modeling and Simulation fellowship, NSF grant no. CCF-1439079, NASA grant no. NNX15AU39A and DoD's PETTT Project PP-CFD-KY07-007.

REFERENCES

- Chernikov, A. N., and N. P. Chrisochoides. 2008, January. "Algorithm 872: Parallel 2D Constrained Delaunay Mesh Generation". *ACM Trans. Math. Softw.* vol. 34 (1), pp. 6:1–6:20.
- Karypis, G., and V. Kumar. 1998. "A Fast and High Quality Multilevel Scheme for Partitioning Irregular Graphs". *SIAM Journal on Scientific Computing* vol. 20 (1), pp. 359–392.
- Shewchuk, J. R. 1997. *Delaunay Refinement Mesh Generation*. Ph. D. thesis. Available as Technical Report CMU-CS-97-137.
- Si, H. 2013. "TetGen, A Quality Tetrahedral Mesh Generator and a 3D Delaunay Triangulator". <http://wias-berlin.de/software/tetgen/>.
- Si, H., and J. R. Shewchuk. 2014. "Incrementally constructing and updating constrained Delaunay tetrahedralizations with finite-precision coordinates". *Engineering with Computers* vol. 30 (2), pp. 253–269.

ENHANCING APPLICATION PERFORMANCE USING MINI-APPS: COMPARISON OF HYBRID PARALLEL PROGRAMMING PARADIGMS

Gary Lawson
Michael Poteat
Masha Sosonkina

Robert Baurle*
Dana Hammond⁺

Department of Modeling,
Simulation and Visualization Engineering
Old Dominion University
Norfolk, VA, USA

Hypersonic Air Breathing Propulsion Branch*
Computational Aero-Sciences Branch⁺
NASA Langley Research Center
Poquoson, VA, USA

ABSTRACT

In many fields, real-world applications for High Performance Computing have already been developed. For these applications to stay up-to-date, new parallel strategies must be explored to yield the best performance; however, restructuring or modifying a real-world application may be daunting depending on the size of the code. In this case, a mini-app may be employed to quickly explore such options without modifying the entire code. In this work, several mini-apps have been created to enhance a real-world application performance, namely the VULCAN code for complex flow analysis developed at the NASA Langley Research Center. These mini-apps explore hybrid parallel programming paradigms with Message Passing Interface (MPI) for distributed memory access and either *Shared MPI* (SMPI) or *OpenMP* for shared memory accesses. Performance testing shows that MPI+SMPI yields the best execution performance, while requiring the largest number of code changes. A maximum speedup of $23\times$ was measured for MPI+SMPI, but only $10\times$ was measured for MPI+OpenMP.

Keywords: Mini-apps, Performance, VULCAN, Shared Memory, MPI, OpenMP

1 INTRODUCTION

In many fields, real-world applications have already been developed. For established applications to stay up-to-date, new parallel strategies must be explored to determine which may yield the best performance, especially with advances in computing hardware. However, restructuring or modifying a real-world application incurs increased cost depending on the size of the code and changes to be made. A mini-app may be created to quickly explore such options without modifying the entire code. Mini-apps reduce the overhead of applying new strategies, thus various strategies may be implemented and compared. This work presents the authors experiences when following this strategy for a real-world application developed by NASA.

VULCAN (Viscous Upwind Algorithm for Complex Flow Analysis) is a turbulent, nonequilibrium, finite-rate chemical kinetics, Navier-Stokes flow solver for structured, cell-centered, multiblock grids that is maintained and distributed by the Hypersonic Air Breathing Propulsion Branch of the NASA Langley Research Center (NASA 2016). The mini-app developed in this work uses the Householder Reflector kernel for solving systems of linear equations. This kernel is used often by different workloads, and is a good candidate to decide what strategy type to apply to VULCAN. VULCAN is built on a single-layer of MPI and

Author emails: { glaws003, mpote001, msosonki } @ odu.edu and { robert.a.baurle, dana.p.hammond } @ nasa.gov

the code has been optimized to obtain perfect vectorization, therefore two-levels of parallelism are currently used. This work investigates two flavors of shared-memory parallelism, OpenMP and Shared MPI, which will provide the third-level of parallelism for the application. A third-level of parallelism increases performance, which decreases the time-to-solution.

MPI has extended the standard to MPI version 3.0, which includes the Shared Memory (SHM) model (Mikhail B. (Intel) 2015, Message Passing Interface Forum 2012), known in this work as Shared MPI (SMPI). This extension allows MPI to create memory windows that are shared between MPI tasks on the same physical node. In this way, MPI tasks are equivalent to threads, except Shared MPI is more difficult for a programmer to implement. OpenMP is the most common shared-memory library used to date because of its ease-of-use (OpenMP 2016). In most cases, only a few OpenMP pragmas are required to parallelize a loop; however, OpenMP is subject to increased overhead, which may decrease performance if not properly tuned.

The major contributions of this work are as follows:

- Created mini-apps to solve $AX = B$ using the Householder reflector kernel from NASA VULCAN real-world code
- Applied MPI+OpenMP scheme to create the OpenMP mini-app
- Applied MPI+SMPI scheme to create the Shared MPI mini-app
- Validated numerical output of each mini-app
- Compared execution performance of all mini-apps

1.1 Related Work

As early as the year 2000, the authors in (Cappello and Etiemble 2000) found that latency sensitive codes seem to benefit from pure MPI implementations whereas bandwidth sensitive codes benefit from hybrid MPI+OpenMP. Also, the authors found that faster processors will benefit hybrid MPI+OpenMP codes if data movement is not an overwhelming bottleneck (Cappello and Etiemble 2000).

Since this time, hybrid MPI+OpenMP implementations have improved, but not without difficulties. In (Drosinos and Koziris 2004, Chorley and Walker 2010), it was found that OpenMP incurs many performance reductions, including: overhead (fork/join, atomics, etc), false sharing, imbalanced message passing, and a sensitivity to processor mapping. However, OpenMP overhead may be hidden when using more threads. In (Rabenseifner, Hager, and Jost 2009), the authors found that simply using OpenMP could incur performance penalties because the compiler avoids optimizing OpenMP loops – verified up to version 10.1. Although compilers have advanced considerably since this time, application users that still compile using older versions may be at risk if using OpenMP. In (Drosinos and Koziris 2004, Chorley and Walker 2010) the authors found that the hybrid MPI+OpenMP approach outperforms the pure MPI approach because the hybrid strategy diversifies the path to parallel execution.

More recently, MPI extended its standard to include the SHM model (Mikhail B. (Intel) 2015). The authors in (Hoeffler, Dinan, Thakur, Barrett, Balaji, Gropp, and Underwood 2015) present MPI RMA theory and examples, which are the basis of the SHM model. In (Gerstenberger, Besta, and Hoeffler 2013), the authors conduct a thorough performance evaluation of MPI RMA, including an investigation of different synchronization techniques for memory windows. In (Hoeffler, Dinan, Buntinas, Balaji, Barrett, Brightwell, Gropp, Kale, and Thakur 2013), the authors investigate the viability of MPI+SMPI execution, as well as compare it to MPI+OpenMP execution. It was found that an underlying limitation of OpenMP is the shared-by-default model for memory, which does not couple well with MPI since the memory model is private-by-default. For this reason, MPI+SMPI codes are expected to perform better, since shared memory is explicit and the memory model for the entire code is private-by-default.

Most recently, a new MPI communication model has been introduced in (Gropp, Olson, and Samfuss 2016), which better captures multinode communication performance, and offers an open-source benchmarking tool to capture the model parameters for a given system. Independent of the shared memory layer, MPI is the *de facto* standard in data movement between nodes and such a model can help any MPI program. The remainder of this paper is organized into the following sections: 2 introduces the Householder mini-apps, 3 presents the performance testing results for the mini-apps considered, and 4 concludes this paper.

2 HOUSEHOLDER MINI-APP

The mini-apps use the householder computation kernel from VULCAN, which is used in solving systems of linear equations. The householder routine is an algorithm that is used to transform a square matrix into triangular form, without increasing the magnitude of each element significantly (Hansen 1992). The Householder routine is numerically stable, in that it does not lose a significant amount of accuracy due to very small or very large intermediate values used in the computation.

The routine works through an iterative process of utilizing Householder transformations to annihilate elements from the column-vectors of the input matrix. The Householder reflector H is applied to a system as:

$$(HA)x = Hb, \quad \text{where} \quad H = (I - 2vv^T)a_i. \quad (1)$$

The Householder operates on a_i , which is a column of A , and v , which is a unit-vector perpendicular to the plane by which the transform is applied. A more detailed discussion of the Householder routine can be found in (Hansen 1992). In this work, the problem to be solved is $AX = B$, where A is a 3-dimensional matrix of size $m \times n \times n$, and X and B are 2-dimensional matrices of size $m \times n$. Each system, represented by m , is independent of all other systems; therefore, this algorithm is embarrassingly parallel.

2.1 Mini-App Design

Mini-apps are designed to perform specific functions. In this work, the important features are as follows:

- Accept generic input,
- Validate the numerical result of the optimized routine,
- Measure performance of the original and optimized routines,
- Tune optimizations.

The generic input is read in from a file, where the file must contain at least one matrix A and resulting vector b . Should only one matrix and vector be supplied, the input will be duplicated for all instances of m . Validation of the optimized routine is performed by taking the difference of the output from the original and optimized routines. The mini-app will first compute the solution of the input using the original routine, and then the optimized routine. This way the output may be compared directly, and relative performance may also be measured using execution time. Should the optimized routine feature one or more parameters that may be varied, they are to be investigated such that the optimization may be tuned to the hardware. In this work, there is always at least one tunable parameter.

One feature that should have been factored into the mini-app design was modularizing the different versions of the Householder routine. In this work, two mini-apps were designed because each implements a different version of the parallel Householder routine; however, it would have been better to design a single mini-app that uses modules to include other versions of the parallel Householder kernel. With this functionality, it would be less cumbersome to work on each version of the kernel.

2.2 Parallel Householder

To parallelize the Householder routine, m is decomposed into separate, but equal chunks that are then solved by each *thread* – shared MPI tasks are equivalent to threads in this work for brevity. However, the original routine varies over m inside the inner-most computational loop (an optimization that benefits vectorization and caching), but the parallel loop must be the outer-most loop for best performance. Therefore, loop blocking has been invoked for the parallel sections of the code. Loop blocking is a technique commonly used to reduce the memory footprint of a computation such that it fits inside the cache for a given hardware. Therefore, the parallel Householder routine has at least one tunable parameter, block size.

In this work, two flavors of the shared memory model are investigated: OpenMP and SMPI. The difference between OpenMP and SMPI lies in how memory is managed. OpenMP uses a public-memory model where all data is available to all threads by default. Public-memory makes it easy to add parallel statements, since the threads will all share this data, but threads are then susceptible to false-sharing, where variables that should otherwise be private are inadvertently shared. Shared MPI uses a private-memory model where data must be explicitly shared between threads, and all data is private by default. Private-memory makes any parallel implementation more complicated, because threads must be instructed to access specific memory for computation. Further, OpenMP creates and destroys threads over the course of execution which is handled internally and is costly to performance. SMPI threads are created upon execution start and persist throughout. This makes managing SMPI threads more difficult, since each parallel phase must be explicitly managed by the programmer. However, the extra work by the programmer may pay off in terms of performance, since less overhead is incurred by SMPI.

3 PERFORMANCE EVALUATION

This section presents the procedure and results of performance testing for the MPI+OpenMP and MPI+Shared MPI Householder Reflector kernel optimizations. For performance testing, it was of interest to vary the number of nodes used for the calculation because many nodes are often used when executing VULCAN with real-world simulations. Up to four nodes have been investigated in this work on a multinode HPC cluster. The number of MPI tasks and OpenMP threads are varied, as well as block size for loop-blocking in the parallel section.

For each mini-app, the optimized version of the Householder routine was validated against the original version by calculating the numerical difference in output. The validation found OpenMP to provide exact numerical solutions (a difference of zero) and SMPI had small numerical discrepancies (10^{-9}).

Computing Platforms The performance evaluation has been conducted on a multinode HPC system *Turing* located at Old Dominion University (HPC Group 2016). Each node on Turing has dual-socket E5-2670 v2 (Ivy-Bridge) CPU's, each socket has 10 cores @ 2.5 GHz and 25 MB cache. A total of 64 GB RAM memory is available on each node. Up to four nodes are used and the network interconnect is Infiniband FDR (Fourteen Data Rate).

Results and Evaluation The performance evaluation varies the size n for the input matrix and the number m of linear systems investigated. Two values of n , 10 and 23, are investigated, which are common sizes based on the VULCAN sample inputs. A third input, nicknamed *Cauchy*, for n is investigated, which is a square Cauchy matrix (Fiedler 2010) of size 10. The number of linear systems m depends on the number of nodes. For the single-node performance tests, m is set to 10k, 100k, 1m, and 5m. For the multinode performance tests, m is set to 100k, 1m, 5m, 10m, 50m, and 100m. The number of threads was varied using powers of two: 1, 2, 4, 8, 16, and 20, because each node on Turing has a total number of 20 cores. The block size is varied using the values 10, 25, 50, 75, 100, 250, 500, 750, and 1000, in order to observe effects on the cache performance and memory latency.

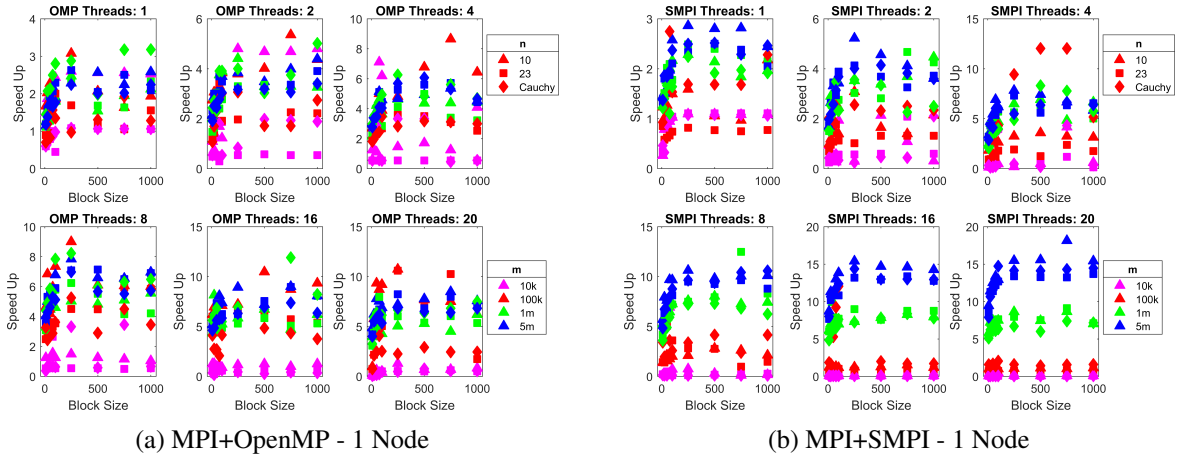


Figure 1: 1 Node Performance Evaluation: Speedup of the optimized mini-apps vs. the original routine with only 1 MPI task.

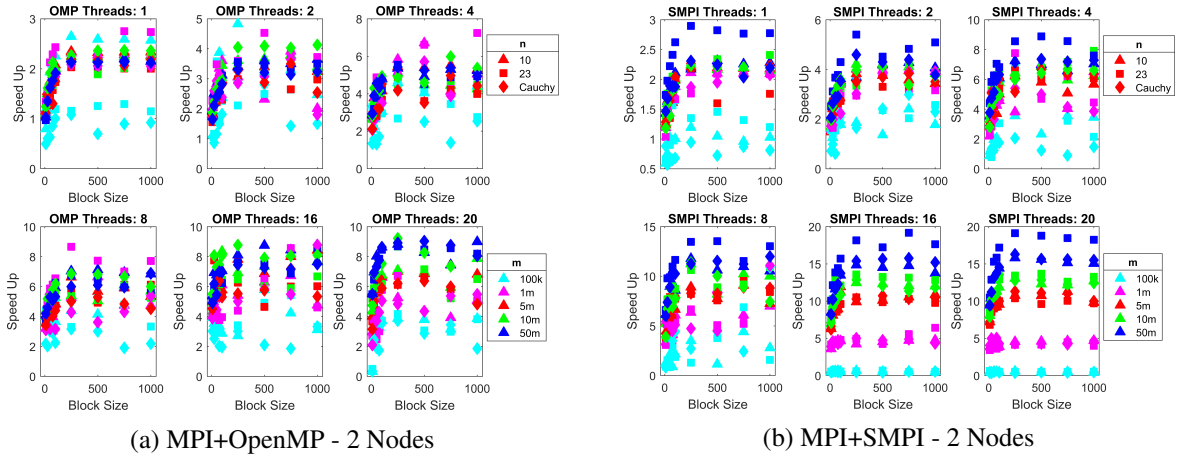


Figure 2: 2 Node Performance Evaluation: Speedup of the optimized mini-apps vs. the original routine with only 1 MPI task per node.

The speedup for the single-node performance tests are shown in Fig. 1. Fig. 1a presents the MPI+OpenMP speedup and Fig. 1b presents the MPI+SMPI speedup where m , n , and Block Size are varied. Speedup is shown on the y-axis, block size on the x-axis, n is represented using a triangle for 10, square for 23, and diamond for Cauchy, and m is represented using colors as shown in the plot legend. Speedup for the multi-node performance tests is shown in: Fig. 2 (2-node) and Fig. 3 (4-node). Notice that the workloads (m) are different for the multi-node tests than for the single-node tests; 100k-50m vs. 10k-5m respectively.

Speedup is a measure of execution performance. A value of one means both versions have equal performance. A value less than one means the optimized version is worse than the original routine, and a value greater than one means the optimized version is better.

The workload, $m=10k$, is only investigated for the single-node case. Notice in Fig. 1 that speedup is consistently one or less. Therefore, the parallel Householder routine must have a sufficient workload to attain any speedup.

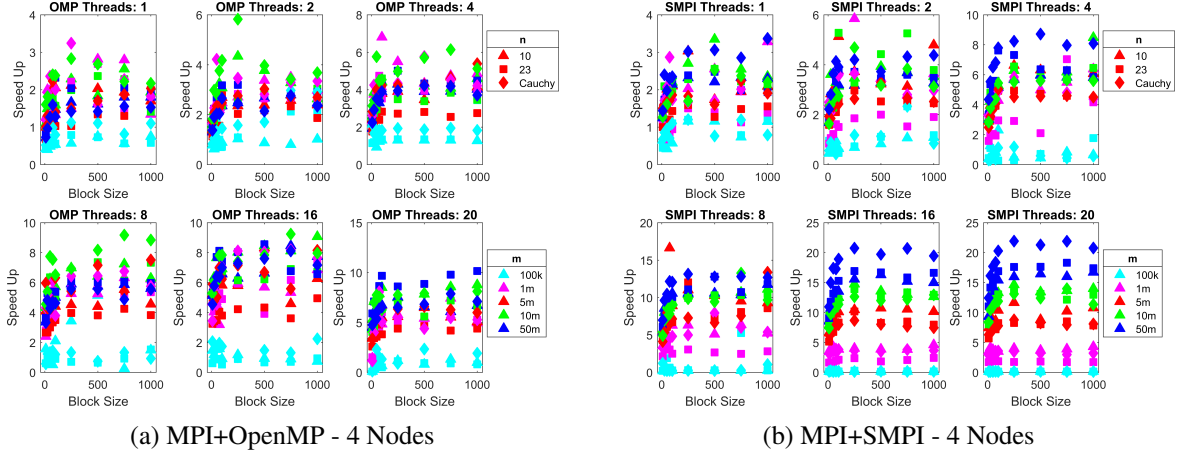


Figure 3: 4 Node Performance Evaluation: Speedup of the optimized mini-apps vs. the original routine with only 1 MPI task per node.

In all performance tests conducted, Figs. 1 to 3, Shared MPI consistently attains the greatest speedup over the original Householder routine. Speedup is normalized against the one thread per node performance for each respective workload (n and m) and block size. The maximum speedup for OpenMP is $12\times$, $9\times$, and $10\times$, and Shared MPI is $18\times$, $19\times$, and $23\times$ for 1, 2, and 4 nodes, respectively. From the performance results, it is apparent that SMPI benefits from less overhead as a result of increased cost to the programmer.

It is interesting to note that the best performing input n varies for SMPI as the number of nodes varies. For one-node and the maximum number of threads, n of 10 has the best speedup. For two nodes, n of 23 has the best speedup, and n of Cauchy has the best speedup for the four-node case. This was an unexpected result, and one that is not obtained when using OpenMP. Further investigation is needed to determine if this is coincidence or a meaningful result.

Blocking performance, without shared memory parallelism, is captured by the one thread tests in Figs. 1 to 3. A max speedup of $3\times$ is consistently measured no matter the mini-app and number of nodes. This finding shows that optimizing the algorithm for cache performance is mildly beneficial and should be considered for performance-bounded computational kernels.

4 CONCLUSION

In this work, mini-apps were developed to optimize the Householder Reflector kernel within NASA real-world application, VULCAN. Two programming paradigms for shared memory parallelism were investigated, OpenMP and Shared MPI, and performance testing was conducted on a multi-node system Turing for up to four nodes. Speedup, the measure of performance, was found to be higher for the Shared MPI version of the Householder mini-app than that for the OpenMP version. Specifically, the speedup for SMPI was up to $1.9\times$ that of OpenMP. With the maximum number of threads, SMPI obtains perfect speedup with sufficiently large workloads ($m=50m$). OpenMP was only able to achieve a speedup of $10\times$, which is half of the expected speedup based on the number of threads used.

ACKNOWLEDGMENT

This effort was supported by NIA subaward activity 2B87 funded through the NASA Langley Computational Digital Transformation initiative and, in part, by the Turing High Performance Computing cluster at Old Dominion University.

REFERENCES

- Cappello, F., and D. Etiemble. 2000, Nov. "MPI versus MPI+OpenMP on the IBM SP for the NAS Benchmarks". In *Supercomputing, ACM/IEEE 2000 Conference*, pp. 12–12.
- Chorley, M. J., and D. W. Walker. 2010. "Performance analysis of a hybrid MPI/OpenMP application on multi-core clusters". *Journal of Computational Science* vol. 1 (3), pp. 168 – 174.
- Drosinos, N., and N. Koziris. 2004, April. "Performance comparison of pure MPI vs hybrid MPI-OpenMP parallelization models on SMP clusters". In *18th International Parallel and Distributed Processing Symposium, 2004. Proceedings.*, pp. 15–.
- Fiedler, M. 2010. "Notes on Hilbert and Cauchy matrices". *Linear Algebra and its Applications* vol. 432 (1), pp. 351 – 356.
- Gerstenberger, R., M. Besta, and T. Hoefer. 2013. "Enabling Highly-scalable Remote Memory Access Programming with MPI-3 One Sided". In *Proceedings of the International Conference on High Performance Computing, Networking, Storage and Analysis, SC '13*, pp. 53:1–53:12. New York, NY, USA, ACM.
- Gropp, W., L. Olson, and P. Samfass. 2016, 9. *Modeling MPI communication performance on SMP nodes: Is it time to retire the ping pong test*, Volume 25-28-September-2016, pp. 41–50. Association for Computing Machinery.
- Hansen, P. B. 1992, June. "Householder Reduction of Linear Equations". *ACM Comput. Surv.* vol. 24 (2), pp. 185–194.
- Hoefer, T., J. Dinan, D. Buntinas, P. Balaji, B. Barrett, R. Brightwell, W. Gropp, V. Kale, and R. Thakur. 2013, December. "MPI + MPI: A New Hybrid Approach to Parallel Programming with MPI Plus Shared Memory". *Computing* vol. 95 (12), pp. 1121–1136.
- Hoefer, T., J. Dinan, R. Thakur, B. Barrett, P. Balaji, W. Gropp, and K. Underwood. 2015, June. "Remote Memory Access Programming in MPI-3". *ACM Trans. Parallel Comput.* vol. 2 (2), pp. 9:1–9:26.
- HPC Group 2016. "Turing Community Cluster General Information". <https://www.odu.edu/facultystaff/research/resources/computing/high-performance-computing>.
- Message Passing Interface Forum 2012. "MPI: A Message-Passing Interface Standard Version 3.0". <http://mpi-forum.org/docs/mpi-3.0/mpi30-report.pdf>.
- Mikhail B. (Intel) 2015. "An Introduction to MPI-3 Shared Memory Programming". <https://software.intel.com/en-us/articles/an-introduction-to-mpi-3-shared-memory-programming>.
- NASA 2016. "VULCAN-CFD". <https://vulcan-cfd.larc.nasa.gov/>.
- OpenMP 2016. "OpenMP: The OpenMP API specification for parallel programming". <http://www.openmp.org/>.
- Rabenseifner, R., G. Hager, and G. Jost. 2009, Feb. "Hybrid MPI/OpenMP Parallel Programming on Clusters of Multi-Core SMP Nodes". In *2009 17th Euromicro International Conference on Parallel, Distributed and Network-based Processing*, pp. 427–436.

PROTECTING THE BIG ONES: HARVEST SLOT LIMITS AND MARINE PROTECTED AREAS FOR THE MANAGEMENT OF CARIBBEAN SPINY LOBSTER

Gaya Gnanalingam, Mark J Butler, Holly Gaff
Department of Biological Sciences
Old Dominion University
ggnan001@odu.edu

ABSTRACT

The Caribbean spiny lobster, *Panulirus argus*, is one of the most iconic species in the Caribbean, supporting some of the region's largest and most economically valuable fisheries. The average size of spiny lobsters, however, has decreased worldwide over the last 30 years with the largest individuals principally targeted by fishers. Given differences in reproductive output relative to size, large lobsters contribute disproportionately to a population's reproductive capacity, thus the loss of these largest individuals is of particular concern to the sustainability of fisheries. Novel management schemes are needed to conserve large breeding lobsters, and a combination of harvest slot limits (maximum and minimum legal sizes) and marine protected areas (areas closed to fishing) are one potential solution. We aim to create a spatially-explicit stage-based model to investigate the potential use of this novel management approach and its effect on lobster abundance, spawning biomass, and harvest.

Keywords: matrix model, spiny lobster, fisheries management, spawning stock

1 INTRODUCTION

The Caribbean spiny lobster, *Panulirus argus*, is one of the Caribbean's most iconic and economically valued species (CRFM, 2013). The species forms the primary fishery for 24 Caribbean nations, employing an estimated 50,000 fishers and an additional 200,000 in fishery related jobs in the region (CRFM, 2011). As a consequence of their high value and market demand however, many populations are currently fully capitalized or overfished (Ehrhardt et al., 2010). Regionally, lobster landings have declined since the early 1990s and the average size of the lobsters caught has also declined (CRFM, 2011). This decrease in size can have a number of significant repercussions particularly in relation to reproduction as there is clear evidence of a positive relationship between body size and reproductive output in both male and female *P. argus* (MacDiarmid & Butler, 1999; Butler et al., 2011; Butler et al., 2015). All things being equal, larger individuals can contribute disproportionately to the reproductive output of the entire population. Sustained and excessive fishing of larger individuals can therefore reduce total stock reproductive success.

Two tools that have the potential to provide increased protection for larger spawning animals are Marine Protected Areas (MPAs; no take fishing zones) and harvest slot limits (combined maximum and minimum size limits that prevent fishers from taking juveniles and large sized adults). The aim of this research is to investigate the potential use of these management tools in populations of *P. argus* connected via larval dispersal and assess their effect on total lobster abundance, spawning biomass and harvest, through the use of a spatially explicit stage-based matrix model.

2 METHODS

Given the difficulties associated with ageing lobsters, the model will be stage-based, with 11 different life history stages represented. Because the focus for this model is the maximization of spawning biomass and the effect of fishing effort on adults of different sizes, the species' many early life history stages will be simplified to just three: one larval stage (L), one juvenile stage (J; 0-55mm carapace length (CL)), and one subadult stage (SA; 55-75mm CL). Eight adult stages (A1-A8; >75 mm CL) defined by 10mm CL increments will be included in the model to capture differences in growth and fecundity relative to body size (Fig. 1). Estimates of mortality (natural and fishing related), and growth will be taken from stock assessments and prior research where they exist, whereas estimates for fecundity will be based on existing data (Bertelsen and Matthews 2001, Butler et al. 2015) and our laboratory experiments conducted in 2015 and 2016.

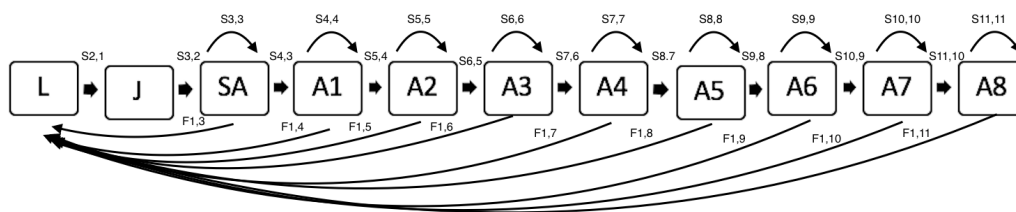


Fig 1. Loop diagram of the matrix model to be implemented for *Panulirus argus*

Initial population sizes (n_0) will be set with reference to lobster landings in the top 10 *P. argus* producing nations in the Caribbean - the assumption being that landings are representative of the relative population size in these areas. Larval dispersal (i.e., movement of the larval stage between neighboring populations) will be based on connectivity matrices from previous work (Kough et al., 2013).

Three management scenarios will be applied: (1) MPAs [0%, 1%, 20% of seascape area], (2) Slot limits (3 different limits applied Caribbean wide), and (3) MPAs and slot limits combined. Each management scenario will also be applied with three different fishing efficiencies (Low, Medium, High). Estimates for fishing intensity will reference existing estimates for *P. argus* in the Caribbean. Ten replicate simulations for each management scenario running for a period of between 20-25 years will be run in R v3.3.1 (R Core Team 2016). The response variables that will be extracted from the simulations for each year include: (a) egg production, (b) spawning biomass, and (c) harvest (fishing mortality). Model sensitivity to variability in parameters will be assessed by way of elasticity analysis.

REFERENCES

- Bertelsen, R.D., and Mathews, T. R. 2001. Fecundity dynamics of female spiny lobster (*Panulirus argus*) in a south Florida fishery and Dry Tortugas National Park lobster sanctuary. *Marine and Freshwater Research*, 52: 1559–1565.
- Butler IV MJ, Heisig-Mitchell J, MacDiarmid AB, Sawnsen RJ. 2011. The effect of male size and spermatophore characteristics on reproduction in the Caribbean spiny lobster, *Panulirus argus*. *New Frontiers in Crustacean Biology* 15: 69-84
- Butler IV MJ, MacDiarmid AB, Gnanalingam G. 2015. The effect of parental size on spermatophore production, egg quality, fertilization success and larval characteristics in the Caribbean spiny lobster, *Panulirus argus*. *ICES Journal of Marine Science* 72 (S1).
- CRFM 2011. Baseline review of the status and management of the Caribbean spiny lobster fisheries in the CARICOM region. CRFM Technical & Advisory Document 2011/5. CRFM Secretariat, Belize pp64
- CRFM 2013. CRFM Annual Report April 1 2012-March 13 2013. CRFM Administrative Report. CRFM Secretariat, Belize pp43
- Ehrhardt NM, Puga R, Butler IV MJ. 2010. Large ecosystem dynamics and fishery management concepts: The Caribbean spiny lobster, *Panulirus argus* fisheries. In: *Towards Marine Ecosystem-Based Management in the Wider Caribbean*. Fanning L, Mahon R, McConney P (eds). Amsterdam
- Kough AS, Paris CB, Butler VI MJ. 2013. Larval connectivity and the international management of fisheries. *PLoS One* <http://dx.doi.org/10.1371/journal.pone.0064970>.
- MacDiarmid, A.B. and M. J. Butler IV. 1999. Sperm economy and limitation in spiny lobsters. *Behavioral Ecology and Sociobiology*. 46: 14-24

USV SIMULATION IN THE RAPID USV MODEL PROTOTYPING SYSTEM

Beau H. Branch
Samantha C. Collins
Lee C. Dumaliang
Nathan D. Gonda
Timothy P. Lane
Kari Miles
Melissa Periman
Dominic A. Scerbo

Department of Modeling, Simulation, and Visualization Engineering
Old Dominion University

5115 Hampton Blvd, Norfolk, VA, USA

{bbran016@odu.edu, scollo03@odu.edu, lduma002@odu.edu, ngondo02@odu.edu,
tlane006@odu.edu, kmile006@odu.edu, mperio03@odu.edu, dscero01@odu.edu}

ABSTRACT

Rapid prototyping can be an extreme benefit to any engineering project, especially in fields where testing the final product is very costly and time-consuming. Most current methods to simulate unmanned surface vehicles (USVs) are high in resolution and computational complexity which may incur too much cost compared to the amount of information it provides about USV operation. Modeling and simulation of simpler, lower fidelity models may help provide the necessary amount of information to solve a certain problem at low cost. The Rapid Unmanned Surface Vehicle (USV) Modeling Prototyping System (RUMPS) utilizes a USV Generation Tool, which aids the user in generating parameters needed in order to construct a model that appropriately represents a prototype of a USV. The Simulator in RUMPS is designed as a framework that connects lower level models within the simulation that may be developed separately or else have different sets of requirements. This architecture allows for components to be changed or replaced to extend functionality or introduce compatibility with existing or future environment or control models. Thus, actual control algorithms may be tested under real life like conditions while still offering the benefits of rapid prototyping.

Keywords: naval design, water vehicles, rapid prototyping

1 INTRODUCTION

Unmanned surface vehicles (USVs), are water vessels capable of operating without an onboard crew. The lack of onboard crews allows USVs to have greater versatility compared to manned counterparts. The Naval Surface Warfare Center, Combatant Craft Division (CCD) currently designs and builds USVs and is interested in expanding the organization's design capabilities, in particular rapid prototyping.

There are several other notable works on testing autonomous control and USV operation using lower fidelity models. A project at the Clark School of Engineering at the University of Maryland uses a machine learning approach to test different autonomous control policies in different environment and situational scenarios (Gupta, 2013). The system takes advantage of lower fidelity motion models in order to generate a large number of possible paths for the USV very quickly. Another rapid prototyping utility is the Marine Systems Simulator (MSS), a MATLAB / Simulink library for simulating ships, underwater vessels, and floating structures. The library consists of a number of Simulink "blocks" that provide the dynamic, guidance, and

control functionality and can be pieced together in different ways to construct the full simulation (Fossen, 2014).

This document focuses mainly on the Simulator and the functionality it provides for the system. A high-level view of the USV system architecture is first presented to provide background to the Simulator architecture. Afterward, a technical description of the Simulator software architecture is discussed in detail. Each major component is listed and described within the context of the USV system. Issues such as time and state management are also highlighted as well as other differences between the high-level and software architectures. Finally, a discussion of the advantages of the architecture are described in relation to rapid prototyping of USVs.

2 SIMULATOR OVERVIEW

2.1 Section Overview

This section presents a high-level view of the RUMPS and Simulator. An overall view of the system is first presented that describes how RUMPS is decomposed. A high-level view of the architecture is then presented where each component is described and interactions are specified. The entire Simulator software architecture is then described in relation to the USV model.

2.2 System Overview

The Rapid USV Model Prototyping System (RUMPS) facilitates the rapid prototyping of USVs through modularity and ease of use. RUMPS consists of three major components: Generation Tool, Simulator, and Visualizer. The Generation Tool (GT) is responsible for allowing the user to quickly generate USV files based on the specified criteria. The GT automatically calculates difficult to determine values, such as center of mass, using information given by the user. The Simulator is responsible for setting up and carrying out the USV simulation. This includes accepting exported data from a USV Generation Tool, setting run simulation parameters, running the simulation for a specified time interval, and communicating simulation data (e.g. state variables) to external components that require it. Finally, the Visualizer presents a visualization and analysis of the data that is produced by the simulation. The Visualizer is customizable and able to perform different analytical operations for the user to compare simulation runs.

2.3 High Level USV Architecture

Figure 1 illustrates a model that represents the authors' understanding of the USV system. These components interact to perform the USV's basic functions for maneuvering in its environment. USV control is comprised of high-level and low-level control. Its purpose is to evaluate sensor data and generate a set of commands that are passed to the physical system to perform some action. The sensor components observe the state of the environment and USV, then produce an interpretation of the data. The environment component represents the external factors of the USV that may impact the motion of the vessel (e.g. wind, waves, current). The Motion Model is the component representing the physical vessel and its motion due to control commands and environmental effects.

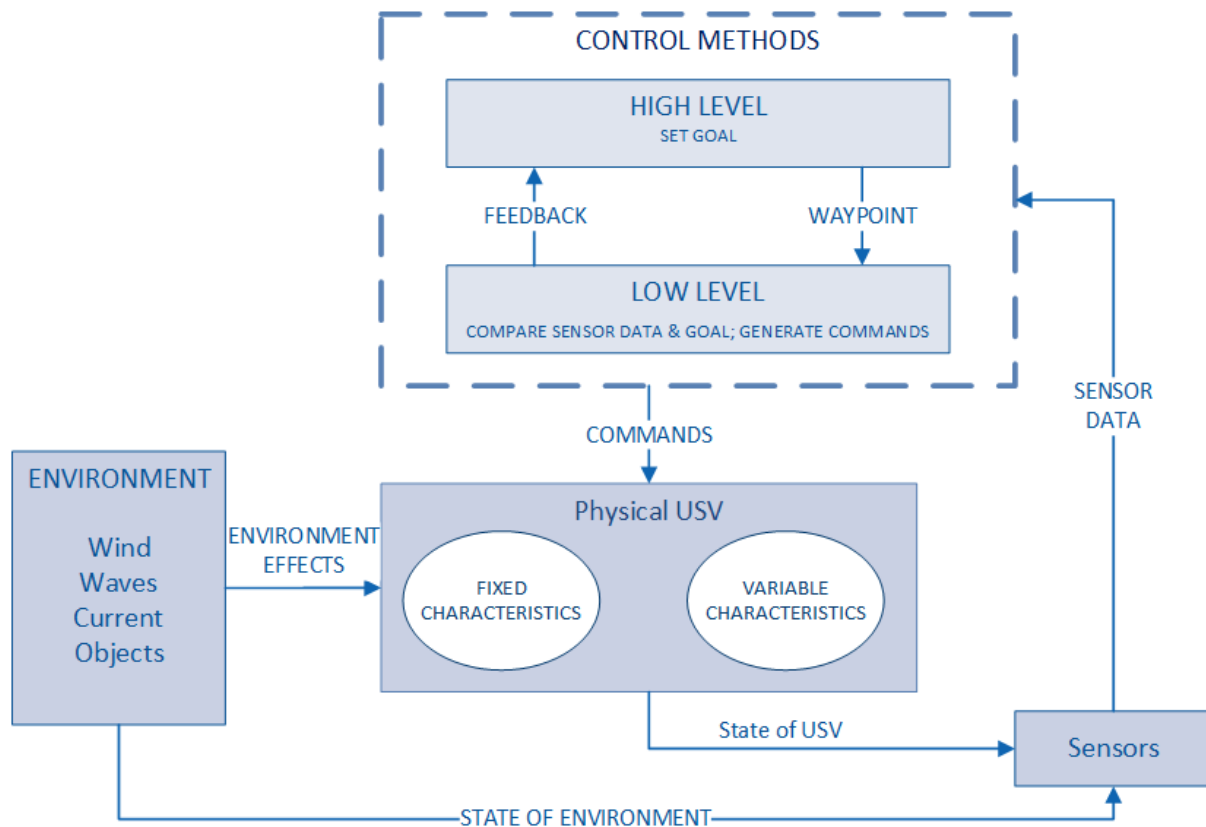


Figure 1. High Level USV Architecture

2.4 Simulation Software Architecture

The diagram shown in Figure 2 illustrates the software architecture of the Simulator. The Simulator is decomposed into components based on those in the high-level model, namely a Motion Model, USV control, environment, and sensors. Additionally, a Simulation Executive is included to provide execution control for the components and manage the advancement of simulation time. Converter components are used to compute and apply information required by the Motion Model from the information produced by the control and environmental components. The Observer controller component obtains simulation output data independently of other components in the Simulator. Arrows in the figure illustrate the flow of information. Solid lines indicate data flow between the components; dashed lines indicate a control flow.

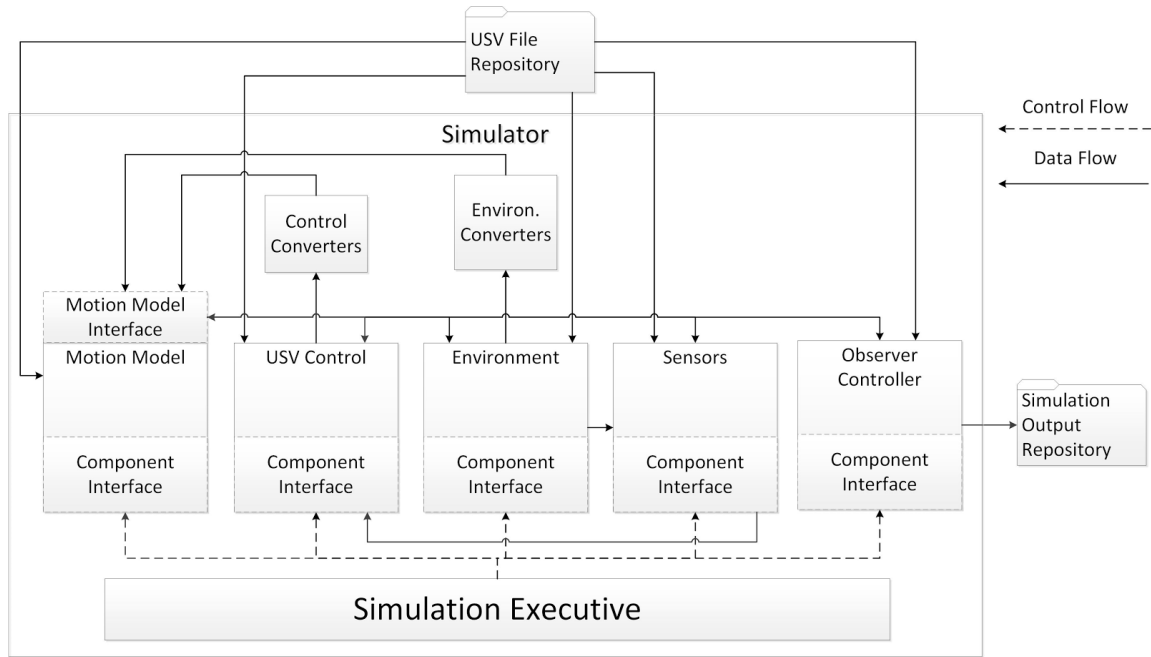


Figure 2. Simulator Software Architecture

3 TECHNICAL DETAILS

This section discusses each component of the Simulator in detail. Simulation input and output is also outlined and explained. Finally, a graphical user interface is described in terms of the information necessary to provide the Simulator before running the USV simulation.

3.1 Simulation Executive

The Simulation Executive acts as a controller for the Simulator. Its primary objectives are managing the advancement of time and determining the order of execution of different components in the system. The Simulation Executive maintains a set of components in the simulation and is provided a common interface with which to communicate with each component (shown in Figure 2). The Simulation Executive also maintains a current simulation time for the entire simulation.

3.2 Motion Model

The Motion Model component encapsulates the motion model functionality of the high-level USV architecture. A motion model interface is discussed that outlines the specific functionality provided to the rest of the Simulator. The component also maintains a mapping of USV identifiers with the internal representation of the USV in the motion model to facilitate isolating the motion model implementation. Lastly, Open Dynamics Engine (ODE) is used to implement the motion model for the prototype.

The Motion Model contains an interface to allow other components to interact with a dynamics engine in a common way. The interface also allows the Motion Model implementation to be replaced if necessary, such as with a different dynamics engine or a custom solution. The interface is given in the form of an API that can be used to gain access to state variables and apply force information to the underlying model. The API for the Motion Model includes the following methods:

- Add Force at Position – Applies the specified force vector and position to the force accumulator tracked in the Motion Model for a specified rigid body. The vector and position are assumed in body-frame coordinates.
- Get Position – Obtains the Cartesian position of the rigid body relative to the inertial

frame coordinate system.

- Get Orientation – Obtains the orientation angle(s) of the USV relative to its vehicle-carried coordinate system centered at the center-of-mass.
- Get Velocity – Obtains the translational velocity of the USV relative to its vehicle-carried coordinate system.
- Get Angular Velocity – Obtains the rotational velocity of the USV relative to its vehicle-carried coordinate system.
- Get Rotation Matrix – Obtains the transformation matrix for rotating from body coordinates to vehicle-carried coordinates.
- Get Translation Matrix – Obtains the transformation matrix for translating from vehicle-carried coordinates to inertial frame coordinates.

The Motion Model also keeps track of a mapping between its internal identifier for a rigid body and an external identifier unique to the USV. This mapping enables the underlying Motion Model to be replaced if necessary without having to rework code or the internal workings of other components of the system. Calls to the Motion Model interface require specifying this external USV identifier as an argument to reference the correct rigid body.

Open Dynamics Engine (ODE) is used for the implementation of the Motion Model. ODE is an open source library for simulated rigid body dynamics developed by Dr. Russel Smith (Smith 2007). The engine is used in several different video games and robotics simulators. It is portable between Windows, OS X, and Linux and does not require any installation besides including the dynamics libraries with the Simulator application.

3.3 Control Components

USV control is responsible for generating a set of commands that influence the Motion Model by producing force information. These commands are generalized as normalized values that are used in the calculations conducted by the force converters. Examples include motor angle and throttle percentage for a motor. The force converters are constructed to accept these variables and other information (such as the current orientation of the USV) in order to produce the appropriate force vector and position. To interface with the Simulation Executive, the USV control should implement the following interface methods:

- Initialization: Initialize control variables
- Pull Updates: Obtain updated variable information from the Motion Model or sensor data from available sensor components
- Compute Forces: Push control variables to converters to calculate force information to apply to the Motion Model

3.4 Environmental Components

An environment component of the Simulator represents the environmental effects that may impact the motion of the USV, such as wind, waves, and current. The actual model for the environment may vary in complexity, but the component must be able to interface with the Simulation Executive. This includes:

- Initialization: Initializes any internal state variables (e.g. wind direction)
- Pull Updates: Obtain state variables from the Motion Model
- Compute Forces: Push control variables to converters to calculate force information to apply to the Motion Model
- Move to Next State: Advances the environment model to its next state, updating any internal state variables

3.5 Sensor Components

The sensor components represent the sensor models for real sensors that may be used by the USV. The sensor model generally operates on actual state variable data from the Motion Model

and converts it into a format representative of the real sensor. Sensor models may also work in a separate coordinate system than the Motion Model and may advance time at different rates. To interface with the Simulation Executive, sensor components should implement the following interface methods:

- Initialization: Initializes any sensor internal variables
- Pull Updates: Obtain state variables from Motion Model for sensor operation
- Move to Next State: Perform operations on state variables to generate new sensor data

3.6 Converters

The Converters are components that act as adapters between the USV control, environmental state variables, and the Motion Model. The Motion Model interface is designed to accept force information that consists of a vector and a position to apply the force at a given point on the USV. The control and environment components are based on models that have their own set of state variables (e.g. throttle percentage, wind direction) that do not directly translate into force information. The converters create the mapping between the control and environment components and the Motion Model.

3.7 Simulation Input

A graphical user interface is included to input information into the Simulator. This information is not available from the Generation Tool and provides the Simulator information on how to manage the simulation runs. This includes:

- USV identifier: used to map a string to a Motion Model USV rigid body
- USV model file: an XML file defining USV characteristics, exported by the Generation Tool
- Simulation Output Directory: used by the Observer to output simulation run results.
- Ending Simulation Time: dictates when the simulation process flow will stop.
- Motion Model Time Step: used by the Motion Model to advance time.
- Sample Time Step: Interval between Observer outputs
- Number of Simulation Runs : How many iterations of simulation performed

3.8 Simulation Output

Simulation output will be handled by Observer components. An Observer is a component that is responsible for outputting simulation data. The Observer allows simulation results to be obtained without affecting the functionality of the other simulation components. An Observer is customized to accept data from several other components and output data in a certain format to a specific location. There can be multiple observers, if necessary, to output certain data to various places. An Observer Controller component is included to manage the observers. By default, at least one Observer must be in the system to obtain simulation results for a visualization tool.

4 SIMULATION EXECUTION

4.1 Section Overview

This section discussed the execution process of the Simulator and how the Simulation Executive manages state and time for the various components of the system. Some background information is first discussed about the issues of managing components that may advance time at different rates. The execution process is presented and a brief description is given of the activities performed at each step of the process.

4.2 Background Information

The simulation executive must be able to ensure that time and state are correctly viewed by the components of the Simulator. Figure 3 provides a view of the state transition for a single component from time 't' to the next time step 't+h'. The S_{t-} indicates the state of the component right before time 't' and S_{t+} indicates the state of the component just after time 't' due to the changes that occur from advancing time.

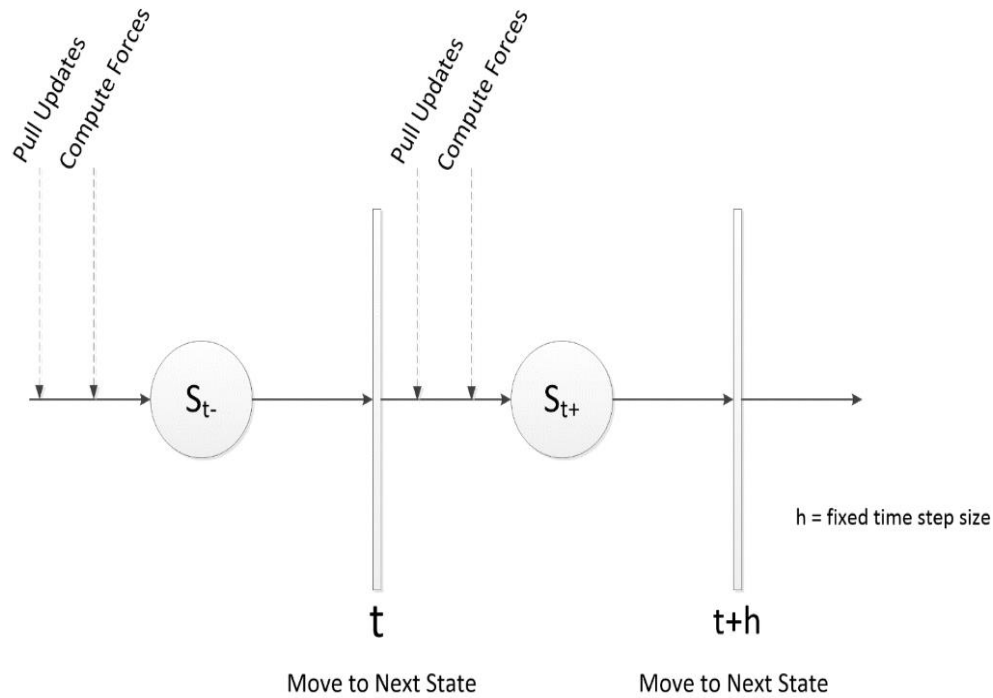


Figure 3. Component State Transition

When considering several components, this view of the state transitions becomes more complicated if the time step 'h' is different between components or if the order in which they update varies. In this case, some components may be using state information from other components that does not match the last time advance for the component. For example, the Motion Model component may be updated with a time step of 1 second before the USV control component which is updated with a time step of 5 seconds. If updated consecutively, the state of the simulation will be different for both components which internally are at different instances in time. The new forces passed to the Motion Model will be from a time stamp in the future and the state pulled from the motion will be at a state in the past for the USV control.

To address this problem, a process flow is constructed to organize and direct each component's actions. This process flow is given in Figure 4. The diagram consists of several steps that relate to actions taken by the components in the architecture. Each step is executed for every component in the simulation before moving to the next step in the process. The steps are comprised of initialization, pulling updated state variable information, computing required forces for the next time step, determining those components that may advance time, and advancing those components to their next state. In addition, the process flow maintains a loop that executes until a specified ending time is reached.

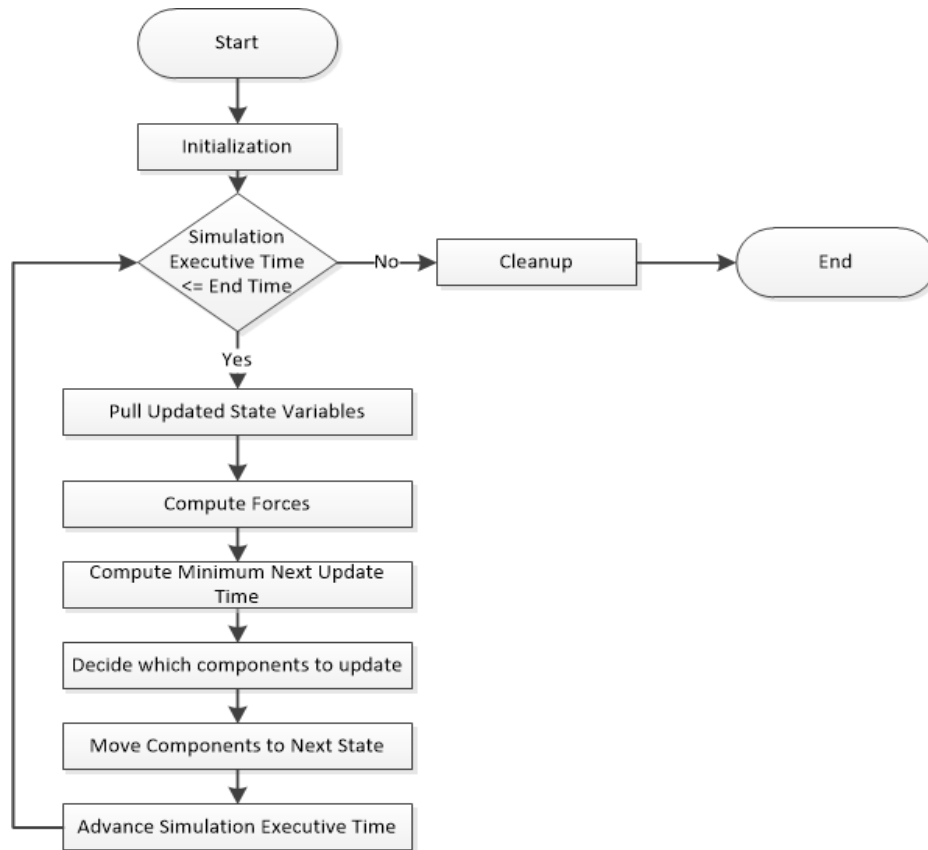


Figure 4. Simulator Process Flow

4.3 Initialization

The initialization step is the first step for the Simulation Executive. Any initialization necessary for the converters are done from the control or environmental components. This step is used to allow each component to set initial values for its variables and to allocate resources where necessary. The information given to each component is put in a common format for all components to view, even if not all of the information is used to perform initialization. The specific information is discussed in the Simulation Input section.

4.4 Pull Updated State Variables

The next step of the simulation is to pull updated state variable information from other components in the system. The information flow is pull-oriented. Each component is able to query for state variable information from the previous time step to use for computations in the current time step. This can include queries to the Motion Model in order to obtain the state variables for the USV. It also can include queries that are not as well-defined such as information about sensor data for the USV control component. In addition, observers in the simulation are able to write state information to the Simulation Output. It is assumed that these components are readily available in order to query state variable information.

4.5 Compute Forces

The Compute Forces step involves the conversion of USV control variables and environment variables into force information that can readily be applied to the Motion Model. Force information refers to a force vector and a position to apply the force that is relative to the USV (in the body-frame coordinate system). The converters use the state variables in applying a mathematical or logical operation to produce a force vector and position that are compatible with the Motion Model coordinate system. After computing the forces, the information is passed to the Motion Model via the given interface. The force information is stored internally in an accumulator or array to keep track of the forces applied to the USV.

4.6 Compute Minimum Next Time Step

The next two steps of the process flow involve determining the minimum time to advance the simulation and determining which components must advance internal state variables to that time. The Simulation Executive assumes that all of the components may have different rates at which time is advanced. The Simulation Executive also assumes that each component has knowledge of the simulation time it must update next. Each component is queried for this next update time and the minimum is computed from all components. With this knowledge, a subset of all components that must be updated can be constructed by comparing the next update time with the minimum time.

4.7 Advance State

The final step of the process flow involves advancing each component of the compiled subset to its next state and advancing simulation time for the executive. Each component will update internal state variables based on information provided to it during the simulation step. For example, the Motion Model will compute the next state of the USV based on the forces being applied to it. Once all components have completed the Advance State step, the Simulation Process will loop back to the Pull Updated State Variables step.

4.8 Cleanup

Once a simulation run is completed (after simulation time surpasses the specified ending time), the simulation resets its internal state. This is required if another simulation run is to be executed. The Cleanup step allows each component to perform a deallocation of resources and resets state variables before beginning another simulation run. After this is completed, the system should mirror its initial state. The cycle then starts again with the Initialization step.

5 ADVANTAGES AND FUTURE EXTENSIONS

5.1 Section Overview

This section describes several advantages of the Simulator design in relation to rapid prototyping. In addition, several extensions are described as areas of future research and development with the Simulator portion of RUMPS.

5.2 Ease of Use

The RUMPS Simulator is designed to be easy to use even for users who have limited to no programming knowledge. This is reflected by the architecture of several of the components. For example, the Simulator has a GUI which allows users to change several options for the simulation, such as USV parameter files, output directory, and simulation time step. This alleviates the user from needing to directly edit the code. There are limited options for initializing simulations as most functions are either standard among all simulations or automatically selected upon runtime. This further alleviates the burden on users for setting up the simulations. Simulation output is also automatically stored in a chosen directory. This output can then be parsed and visualized by the Visualizer to allow for easy comprehension of the simulation results. A Generation Tool will be able to provide USV parameter files. The RUMPS Simulator is intended to provide users rapid prototyping abilities and reduce the software literacy needed to use. These design decisions support this goal.

5.3 Modular Architecture

The Simulator is also modular in its design. Partitions in the responsibilities of the Simulator allow the Simulator to be connected to a variety of external components. Components can be used interchangeably as a result of this design decision, allowing the addition of new features and functionality to RUMPS. An additional advantage of this modularity is the RUMPS Simulator will be able to handle new iterations of external software, such as advances in

environmental simulations. This also allows the Simulator to act as an emulator when connected to real life hardware. The overall design of the Simulator also allows users to extend the capabilities of RUMPS.

5.4 Emulating Capabilities

As a future extension, the ability of the Simulator to connect to actual hardware can greatly extend the functionality of RUMPS. The interaction would allow users to perform tests on the physical hardware as well as the device's software using RUMPS. The ability to run simulations based on real-world, real-time weather conditions is also possible with RUMPS. In addition, the Simulator's emulating capabilities could allow USV engineers to train the use of control software with better correspondence to a completely physical system.

5.5 Higher Resolution

The modularity of the Simulator facilitates the ability to increase the resolution in the future. Higher degree-of-freedom motion models may be used to capture more complex environmental effects. Converter components may also be modified to implement more realistic control behavior such as input delay or smoothing. Additionally, geographic information could be used to replicate other aspects of the environment within the simulation. This includes other physical objects within the simulation and geographic barriers or obstacles to USV movement.

6 CONCLUSION

The RUMPS Simulator is the component of the larger RUMPS system responsible for running simulations of unmanned surface vehicles. The Simulator can accept a variety of USVs, environmental effects, and control commands as input and output data based on simulations ran using these inputs. The Motion Model moves the simulated USV based on rigid body dynamics and applied forces. Communication between components and time management is handled by a simulation executive. Converter components map environment and control inputs to forces to act on the simulated vessel. Observers record and save simulation data based on state variables in each time step of the simulations.

The design of the Simulator allows users to quickly and easily simulate various USVs and scenarios. This design also provides user extensibility, modularity, and emulation capabilities. In the future, RUMPS can be applied to simulations with higher resolution and fidelity. Higher resolution motion model and converter components are possible ways to increase the fidelity of the Simulator. A geographic overlay and geodetic coordinates are additional future extensions providing higher resolution. While there are many extensions that can be made to the Simulator, the architectural design of the Simulator is capable of integrating them.

REFERENCES

- Fossen, T. I. and T. Perez. 2004. "Marine Systems Simulator (MSS)". <http://www.marinecontrol.org> Accessed Jan. 21, 2017
- Gupta, Satyandra K. 2013. "Developing Autonomy for Unmanned Surface Vehicles." University of Maryland, Accessed 25 Jan 2017.
<<http://drum.lib.umd.edu/bitstream/handle/1903/14704/USV%20Report.pdf>>
- Smith, Russell. 2007. "Open Dynamics Engine - Home". *Ode.Org*. <http://www.ode.org> Accessed Jan. 10, 2017.

AUTHOR BIOGRAPHIES

BEAU H. BRANCH is currently a senior undergraduate student in Modeling and Simulation Engineering at Old Dominion University. He is also minoring in Applied Mathematics. His email is bbran016@odu.edu.

SAMANTHA C. COLLINS is currently enrolled in the linked Bachelor of Science in Modeling and Simulation Engineering to Master of Science in Modeling and Simulation degree program at Old Dominion University. Her email address is scoll003@odu.edu.

LEE C. DUMALIANG is a current senior at Old Dominion University pursuing a major in Modeling and Simulation Engineering and minors in Computer Science and Engineering Management. He has several years of experience with software development in the defense industry. His email address is lduma002@odu.edu.

NATHAN D. GONDA is an undergraduate student majoring in Modeling and Simulation Engineering with a minor in Computer Science at Old Dominion University. He has earned an Associate in Computer Science from Paul D. Camp Community College and has experience in computer programming for about 5 years. His research interests are in modeling and simulation visualization and computer game design. His email address is ngond002@odu.edu.

TIM P. LANE is a current senior at old dominion university pursuing a degree in modeling and simulation engineering he has had many years of experience in software development and design. His email is tlane006@odu.edu.

KARI R. MILES is currently enrolled at Old Dominion University pursuing a major in Modeling and Simulation. He has some experience with visualization programming. His email address is kmile006@odu.edu.

MELISSA R. PERIMAN is currently a senior undergraduate student majoring in Modeling and Simulation Engineering with a minor in Computer Engineering. Her email is mperi003@odu.edu.

DOMINIC A. SCERBO is a current senior at old dominion university pursuing a degree in Modeling and Simulation Engineering with a minor in Biomedical Engineering and Computer Science. He has experience with testing and evaluation in the defense industry. His email address is dscer001@odu.edu.

INSTABILITY AND PATTERNS OF ACTIVE SUSPENSIONS OF LIQUID CRYSTALS

R. Anthony Williams

Dr. Ruhai Zhou

Department of Mathematics and Statistics
Old Dominion University
5115 Hampton Boulevard
Norfolk, VA, USA
ra1willi@odu.edu

Department of Mathematics and Statistics
Old Dominion University
5115 Hampton Boulevard
Norfolk, VA, USA
rzhou@odu.edu

ABSTRACT

We first present kinetic model equations for active, anisotropic fluids arising from biological and materials science applications. The equations include transport equations for the polarity vector and the local concentration, as well as the Navier-Stokes equations. Then we examine the dynamics of the polar active liquid crystals near equilibrium in a uni-axial system. This is done by conducting a linear stability analysis about constant steady states in order to explore near equilibrium dynamics near the steady states. Numerical simulation results for different model parameters, including the parameters being perturbed, are provided.

Keywords: liquid crystals, stability analysis, numerical simulations.

1 INTRODUCTION

An active suspension can be described as a large-scale collection of self-propelled particles that interact with each other. Active suspensions have been the focus of many researchers throughout the years due to the complex dynamics they entail as well as the numerous applications. Some of the many applications of active suspensions include those in biology such as (Pedley and Kessler 1992), as well as those in technology (Dreyfus et al. 2005) and (Darnton et al. 2004) and many many others. Two of the most prominent applications of liquid crystals are liquid crystal displays (Gray, Harrison, and Nash 2007) and Kevlar (Zhang and Kumar 2008) which is used to produce many things including personal armor. The goal of this paper is to study the dynamics of polar active liquid crystals near equilibrium. This topic has been studied by many including (Xiao-Gang, Forest, and Wang 2014) as well as others, but these studies have mainly included a constant local concentration c . However, the local concentration is not constant generally and it depends on many things including the polarity vector, the potential for spatial inhomogeneity, and the polarity strength. Thus, in this paper we will examine the effects of allowing the local concentration to remain non-constant on the stability of the system near equilibrium. First we will present the mathematical model we will be studying, and then complete stability analysis on a linearized system. Within this analysis we will examine further a specific case of the wave number in the perturbation, and finally we present some numerical results from simulations that were conducted.

2 MATHEMATICAL MODEL

The mathematical model used here involves a system of equations including a transport equation for the polarity vector and the local concentration as well as the Navier-Stokes equations, all of which are derived from the kinetic equations. Using the notations $c = \langle 1 \rangle$, $\mathbf{p} = \langle \mathbf{m} \rangle$, $\mathbf{M} = \langle \mathbf{m}\mathbf{m} \rangle$, and $\mathbf{M}_3 = \langle \mathbf{m}\mathbf{m}\mathbf{m} \rangle$ where \mathbf{m} is a unit vector representing the axis of symmetry of the rod-like molecule and $\langle \cdot \rangle$ is the molecular average, the equations for the polarity vector \mathbf{p} and the local concentration c are

$$\begin{aligned} & \frac{\partial \mathbf{p}}{\partial t} + \mathbf{v} \cdot \nabla \mathbf{p} + U_0 \nabla \cdot \mathbf{M} \\ = & D_s \left(\nabla^2 \mathbf{p} + \nabla \cdot (N_1 \nabla c \mathbf{p} - \alpha \nabla \mathbf{p} \cdot \mathbf{M} - \frac{3N}{2} \nabla \mathbf{M} : \mathbf{M}_3) \right) \\ & - D_r (2\mathbf{p} + \alpha(\mathbf{M} - c\mathbf{I}) \cdot \mathbf{p} + 3N(\mathbf{M}_3 : \mathbf{M} - \mathbf{M} \cdot \mathbf{p})) \\ & + \boldsymbol{\Omega} \cdot \mathbf{p} + a(\mathbf{D} \cdot \mathbf{p} - \mathbf{D} : \mathbf{M}_3) \end{aligned} \quad (1)$$

$$\begin{aligned} & \frac{\partial c}{\partial t} + \mathbf{v} \cdot \nabla c + U_0 \nabla \cdot \mathbf{p} \\ = & D_s \left(\nabla^2 c + \nabla \cdot (N_1 c \nabla c - \alpha \nabla \mathbf{p} \cdot \mathbf{p} - \frac{3N}{2} \nabla \mathbf{M} : \mathbf{M}) \right) \end{aligned} \quad (2)$$

where U_0 is the rod self-propulsion speed, α is the polarity strength, D_s is the translational diffusion coefficient, D_r is the rotational diffusion coefficient, N_1 is the strength of the potential for spatial inhomogeneity, and a is the geometric particle parameter for rods. At equilibrium in a uni-axial system, the orientation tensor \mathbf{Q} (symmetric and traceless) may be written as $\mathbf{Q} = s(\mathbf{m}\mathbf{m} - \frac{1}{3}\mathbf{I})$, where s is the orientational order parameter (Doi and Edwards 1986). Also, the vector \mathbf{m} is parallel to the direction of average orientation so that $\mathbf{m} = \frac{\mathbf{p}}{\|\mathbf{p}\|}$, and then let the order parameter be $s = \|\mathbf{p}\|^2$. Thus it is true that $\mathbf{Q} = \mathbf{p}\mathbf{p} - \frac{1}{3}(\mathbf{p} \cdot \mathbf{p})\mathbf{I}$. Since $\mathbf{M} = \mathbf{Q} + \frac{1}{3}c\mathbf{I}$ (Yang et al. 2010) and \mathbf{Q} is traceless, we note that $\text{tr}(\mathbf{M}) = c$. This relationship gives rise to the closure rule $\mathbf{M} = \mathbf{p}\mathbf{p} + \frac{1}{3}(c - \mathbf{p} \cdot \mathbf{p})\mathbf{I}$. Additionally, we neglect the nematic strength N , so taking equations (1) and (2) and applying the closure rule results in the following governing system of equations

$$\begin{aligned} & \frac{\partial \mathbf{p}}{\partial t} + \mathbf{v} \cdot \nabla \mathbf{p} + U_0 \left[\nabla \cdot \mathbf{p}\mathbf{p} + \frac{1}{3} \nabla (c - \mathbf{p} \cdot \mathbf{p}) \right] \\ = & D_s \left(\nabla^2 \mathbf{p} + \nabla \cdot (N_1 \nabla c \mathbf{p} - \alpha [\nabla \mathbf{p} \cdot \mathbf{p}\mathbf{p} + \frac{1}{3}(c - \mathbf{p} \cdot \mathbf{p}) \nabla \mathbf{p}]) \right) \\ & - 2D_r \left(1 - \frac{1}{3}\alpha c + \frac{1}{3}\alpha \mathbf{p} \cdot \mathbf{p} \right) \mathbf{p} + \boldsymbol{\Omega} \cdot \mathbf{p} + a \left(\mathbf{D} \cdot \mathbf{p} - \mathbf{D} : \mathbf{p}\mathbf{p}\mathbf{p} \right) \end{aligned} \quad (3)$$

$$\frac{\partial c}{\partial t} + \mathbf{v} \cdot \nabla c + U_0 \nabla \cdot \mathbf{p} = D_s (\nabla^2 c + \nabla \cdot (N_1 c \nabla c - \alpha \nabla \mathbf{p} \cdot \mathbf{p})) \quad (4)$$

$$\frac{d\mathbf{v}}{dt} + (\mathbf{v} \cdot \nabla) \mathbf{v} = \nabla \cdot (-\Pi \mathbf{I} + \boldsymbol{\tau}_p + \boldsymbol{\tau}_a), \quad \nabla \cdot \mathbf{v} = 0 \quad (5)$$

$$\begin{aligned} \boldsymbol{\tau}_p = & 2Re^{-1} \mathbf{D} + G \left[\mathbf{p}\mathbf{p} - \frac{1}{3}(\mathbf{p} \cdot \mathbf{p})\mathbf{I} \right] \\ & - \frac{1}{6}\alpha_0 G [2\mathbf{p}\mathbf{p} - (\mathbf{p}\mathbf{p}\mathbf{p} \cdot \mathbf{p} + \mathbf{p} \cdot \mathbf{p}\mathbf{p}\mathbf{p})] \\ & + Re_2^{-1} \left[\mathbf{D} \cdot \left(\mathbf{p}\mathbf{p} - \frac{1}{3}(\mathbf{p} \cdot \mathbf{p})\mathbf{I} \right) + \left(\mathbf{p}\mathbf{p} - \frac{1}{3}(\mathbf{p} \cdot \mathbf{p})\mathbf{I} \right) \cdot \mathbf{D} \right] \end{aligned} \quad (6)$$

$$\tau_a = G\zeta_a(\mathbf{p}\mathbf{p} - \frac{1}{3}(\mathbf{p} \cdot \mathbf{p})\mathbf{I}) \quad (7)$$

where G is the anisotropic stress coefficient, α_0 tunes the stored stress due to polarity, and ζ_a is the particle activation parameter. The tensors \mathbf{D} and $\mathbf{\Omega}$ are defined as the rate of strain tensor and vorticity tensor, respectively. They are the symmetric and antisymmetric parts of the velocity gradient tensor, or equivalently $\mathbf{D} = \frac{1}{2}(\nabla\mathbf{v} + \nabla\mathbf{v}^\top)$ and $\mathbf{\Omega} = \frac{1}{2}(\nabla\mathbf{v} - \nabla\mathbf{v}^\top)$.

3 LINEAR STABILITY ANALYSIS

3.1 Equilibrium

We seek to study the stability of this system near equilibrium. This occurs when the flow velocity $\mathbf{v} = 0$ and the polarity vector \mathbf{p} is homogeneous in space. If we examine equation (5), we can see that there will be two different equilibrium states depending on \mathbf{p} . The first equilibrium is when $\mathbf{p} = \mathbf{0}$ and α takes any value, and this is known as the isotropic state. The second equilibrium occurs when $1 - \frac{1}{3}\alpha c + \frac{1}{3}\alpha\mathbf{p} \cdot \mathbf{p} = 0$. At this equilibrium, $c = \langle 1 \rangle = \int_{||\mathbf{m}||=1} 1 \cdot f(\mathbf{m})d\mathbf{m} = 1$ since the nanorod number density function f is normalized, so the polar state will be when

$$||\mathbf{p}||^2 = \frac{\alpha - 3}{\alpha}.$$

We will study the stability of this system in the polar state because it is more interesting than the isotropic state.

3.2 Stability

Now a solution to the governing system of equations (5)-(9) can be found by perturbing the initial state using a plane wave perturbation as follows:

$$\mathbf{v} = \mathbf{v}_0 + \epsilon\mathbf{v}_1, \mathbf{p} = \mathbf{p}_0 + \epsilon\mathbf{p}_1, \Pi = \Pi_0 + \epsilon\Pi_1, c = c_0 + \epsilon c_1$$

where $\epsilon \ll 1$ and the perturbations are $\mathbf{v}_1 = e^{i\mathbf{k} \cdot \mathbf{x}}\mathbf{v}_1(t)$, $\mathbf{p}_1 = e^{i\mathbf{k} \cdot \mathbf{x}}\mathbf{p}_1(t)$, $\Pi_1 = e^{i\mathbf{k} \cdot \mathbf{x}}\Pi_1(t)$, and $c_1 = e^{i\mathbf{k} \cdot \mathbf{x}}c_1(t)$ (Xiao-Gang, Forest, and Wang 2014). Here $\mathbf{k} = (k_1, k_2, k_3)$ is the wave vector, Π_1 and c_1 are constant unknown scalars, and \mathbf{v}_1 and \mathbf{p}_1 are constant unknown vectors. If we truncate at the linear order of ϵ , then from equation (3) we can get an equation for the first order correction \mathbf{p}_1 which is given by

$$\begin{aligned} \mathbf{p}_1' + \left[S_1\mathbf{I} + iU_0\mathbf{p}_0\mathbf{k} - \frac{2i}{3}U_0\mathbf{k}\mathbf{p}_0 + \frac{4}{3}\alpha D_r\mathbf{p}_0\mathbf{p}_0 \right] \mathbf{p}_1 \\ + \left[S_2\mathbf{I} + \frac{i}{2}(1-a)\mathbf{k}\mathbf{p}_0 + ia(\mathbf{k} \cdot \mathbf{p}_0)\mathbf{p}_0\mathbf{p}_0 \right] \mathbf{v}_1 \\ + \left[\frac{i}{3}U_0\mathbf{k} + D_s N_1 ||\mathbf{k}||^2 \mathbf{p}_0 - \frac{2}{3}\alpha D_r \mathbf{p}_0 \right] c_1 = 0. \end{aligned} \quad (8)$$

In a similar fashion, equation (4) provides us with

$$c_1' + S_3 c_1 + (\mathbf{k} \cdot \mathbf{p}_1) S_4 = 0. \quad (9)$$

The Navier-Stokes equations (5) together with the passive nematic stress τ_p (6) and the active stress τ_a (7) will result in

$$\begin{aligned}
& \left[S_5 \mathbf{I} - \frac{2i}{3} \alpha_0 G \left((\|\mathbf{p}_0\|^2 - \frac{1}{2} \alpha_0) \mathbf{p}_0 \mathbf{k} + 2(\mathbf{p}_0 \cdot \mathbf{k}) \mathbf{p}_0 \mathbf{p}_0 \right) \right. \\
& \quad \left. + \frac{iG}{3} (1 + \zeta_a) (2\mathbf{k} \mathbf{p}_0 - \mathbf{p}_0 \mathbf{k}) \right] \mathbf{p}_1 + \mathbf{v}_1' \\
& + \left[S_6 \mathbf{I} + \frac{1}{2Re_2} \left(\|\mathbf{k}\|^2 \mathbf{p}_0 \mathbf{p}_0 + (\mathbf{p}_0 \cdot \mathbf{k}) (\mathbf{p}_0 \mathbf{k} + \mathbf{k} \mathbf{p}_0) \right) \right] \mathbf{v}_1 \\
& + i\mathbf{k} \Pi_1 = 0, \\
& \mathbf{k} \cdot \mathbf{v}_1 = 0.
\end{aligned} \tag{10}$$

All of the constants $S_i, i = 1, 2, 3, 4, 5, 6$ are defined in the appendix. From the linearized momentum equation above, we can see that Π_1 can be eliminated from the system. This is done by taking the dot product of both sides of the equation and solving for Π_1 and then substituting in for the result, which is

$$\Pi_1 = \frac{(C\mathbf{p}_1 + D\mathbf{v}_1) \cdot \mathbf{k}}{i\|\mathbf{k}\|^2}.$$

The linearized system of equations can now be written in the following way:

$$\begin{pmatrix} c_1' \\ \mathbf{p}_1' \\ \mathbf{v}_1' \end{pmatrix} = \begin{pmatrix} -S_3 & -S_4 \mathbf{k}^\top & 0 \\ \mathbf{z} & A & B \\ 0 & C' & D' \end{pmatrix} \begin{pmatrix} c_1 \\ \mathbf{p}_1 \\ \mathbf{v}_1 \end{pmatrix}$$

where $\mathbf{z} = \frac{i}{3} U_0 \mathbf{k} + (D_s N_1 \|\mathbf{k}\|^2 - \frac{2}{3} \alpha D_r) \mathbf{p}_0$, $C' = (\mathbf{I} - \frac{1}{\|\mathbf{k}\|^2} \mathbf{k} \mathbf{k}) \cdot C$ and $D' = (\mathbf{I} - \frac{1}{\|\mathbf{k}\|^2} \mathbf{k} \mathbf{k}) \cdot D$. The matrices A, B, C , and D are defined in the appendix. For simplicity, we will denote the above system by $\mathbf{y}' = F\mathbf{y}$ where \mathbf{y} is the vector of unknown quantities, \mathbf{y}' is the time derivative of the unknown quantities, and F is the 7x7 coefficient matrix operating on \mathbf{y} . In order to examine the stability of this system, we need to determine the eigenvalues of the matrix F .

As stated previously, we are interested in the dynamics of this system near equilibrium which corresponds to the liquid crystal state $\|\mathbf{p}_0\|^2 = \frac{\alpha-3}{\alpha}$. Without loss of generality we can take an initial polarity vector position to be $\mathbf{p}_0 = (p, 0, 0)$, where $p = \sqrt{\frac{\alpha-3}{\alpha}}$. An additional assumption made is that the initial velocity of the surrounding fluid $\mathbf{v}_0 = 0$. This creates five more zero entries in the coefficient matrix F , leaving 38 total non-zero entries.

3.3 Specific Case of Wave Number k

This section addresses the case for when $k_2 = k_3 = 0$ but $k_1 \neq 0$, or $\mathbf{k} = (k_1, 0, 0)$. Also, we will let the rod self-propulsion velocity $U_0 = 0$ for this analysis. Thus, this case leads to the coefficient matrix F having eigenvalues equal to

$$\begin{aligned}
\lambda_1 &= 0 \\
\lambda_{2,3} &= -\frac{1}{2} \left[\frac{4}{3} \alpha D_r p^2 + D_s k_1^2 (N_1 - \alpha p^2) \right. \\
& \quad \left. \mp \sqrt{\left(D_s k_1^2 (1 + \alpha p^2) - \frac{4}{3} \alpha D_r p^2 \right)^2 + 4 D_s p^2 k_1^2 \left(D_s N_1 k_1^2 - \frac{2}{3} \alpha D_r \right)} \right]
\end{aligned}$$

$$\lambda_{4,5} = \lambda_{6,7} = -\frac{1}{2} \left\{ k_1^2 \left[D_s(N_1 - \alpha p^2) + \frac{1}{Re} + \frac{p^2}{Re_2} \right] \right. \\ \left. \mp k_1 \sqrt{k_1^2 \left[\frac{1}{Re} + \frac{p^2}{Re_2} - D_s(N_1 - \alpha p^2) \right]^2 + Gp^2(1+a) \left[\frac{1}{3}\alpha_0 - \frac{2}{3}\alpha_0 p^2 - 1 - \zeta_a \right]} \right\}$$

The goal is to determine what values for the parameters α and ζ_a will result in a positive real part of the eigenvalue. If the real part of the eigenvalue is positive, then this leads to an unstable state in the system (Sanchez 1968). These states are very interesting dynamically, so we wish to find situations where this is true. In addition, we are interested in the case for $\zeta_a < 0$ which is referred to as the pusher case (Saintillan and Shelley 2008). The remaining parameters will be fixed and will be equal to $D_s = 0.02$, $D_r = 5$, $G = 4$, $\alpha_0 = 1$, $Re = 5$, and $Re_2 = 5$ and $a = 0.5$. Also, we will let $k_1 = 1$ because the magnitude of k_1 has no significant impact on whether the eigenvalue has positive real part or not. We discovered that the eigenvalue λ_4 listed above will be positive for $\alpha > 3$ provided that $\zeta_a < 0$ and $|\zeta_a|$ is large enough. In this case, $-\zeta_a$ dominates the eigenvalue λ_4 . The relationship between ζ_a and α which represents when $\lambda_4 = 0$ can be seen in Figure 1 below. When $\alpha > 3$, $\lambda_4 > 0$ whenever $\zeta_a < 0$ and smaller than a critical value $\zeta_{a,c}$ which will be visibly apparent in Figure 2 below. Thus, under these circumstances the system will have an unstable state.

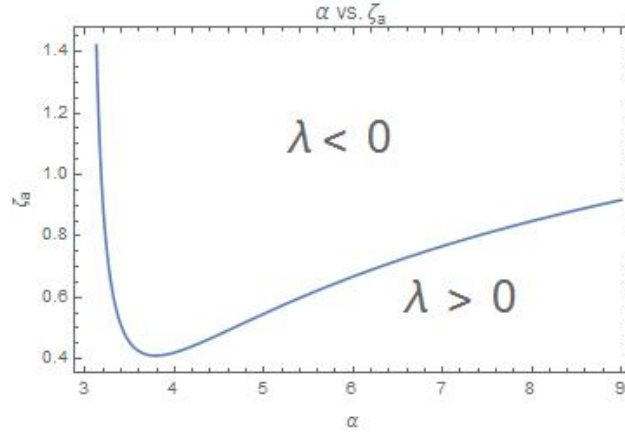


Figure 1: A depiction of the relationship for $\zeta_{a,c} = \zeta_{a,c}(\alpha)$. Along this curve, the eigenvalue $\lambda_4 = 0$.

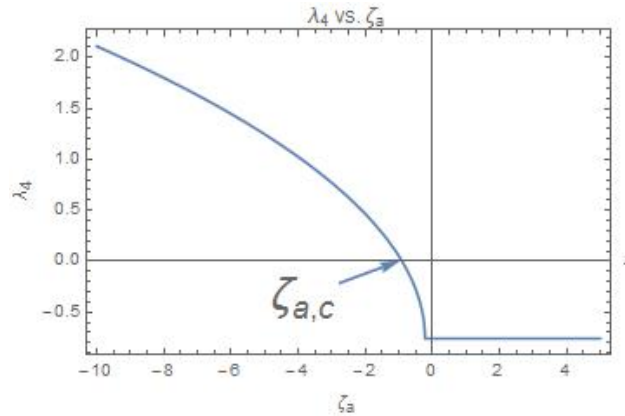


Figure 2: A special case of $\text{Re}(\lambda_4)$ vs. ζ_a for when $\alpha = 5$, showing the critical value $\zeta_{a,c}$.

4 NUMERICAL SIMULATIONS

4.1 Outline of Numerical Method

The Navier-Stokes equations were solved using the finite difference method on a staggered grid as in (Griebel, Dornseifer, and Neunhoeffer 1998). The horizontal velocity component v_1 is located at the midpoints of the vertical edges of each cell, and the vertical velocity component v_2 is located at the midpoints of the horizontal edge of each cell. Both the pressure Π and v_3 are located at the cell centers. The reason a staggered grid might be used is to prevent possible pressure oscillations.

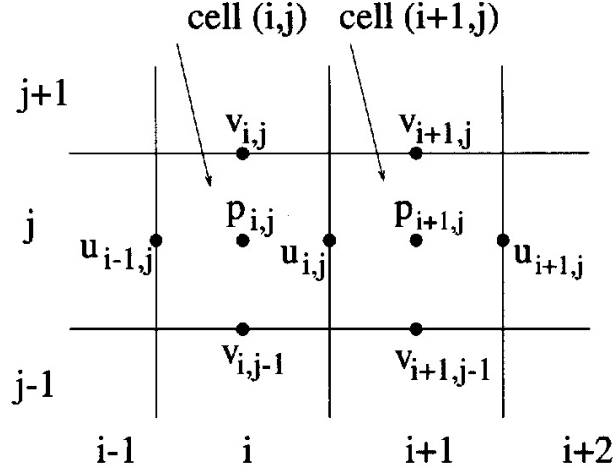


Figure 3: An illustration of the staggered grid deployed in (Griebel, Dornseifer, and Neunhoeffer 1998), where p is pressure in this context.

Once the solution to the Navier-Stokes equations is found numerically, which yields the velocity and pressure, these values are input into the transport equations for the polarity vector and local concentration. These equations are subsequently solved using a semi-implicit method where the time derivative is represented by a backward Euler scheme, resulting in the following two elliptic equations: $[1 - D_s \Delta t \nabla^2] \mathbf{p}^{n+1} = \mathbf{p}^n + \Delta t \mathbf{F}^n$ and $[1 - D_s \Delta t \nabla^2] c^{n+1} = c^n + \Delta t H^n$ where

$$\begin{aligned} \mathbf{F} = & D_s \left(\nabla^2 \mathbf{p} + \nabla \cdot (N_1 \nabla c \mathbf{p} - \alpha [\nabla \mathbf{p} \cdot \mathbf{p} \mathbf{p} + \frac{1}{3} (c - \mathbf{p} \cdot \mathbf{p}) \nabla \mathbf{p}]) \right) \\ & - 2D_r \left(1 - \frac{1}{3} \alpha c + \frac{1}{3} \alpha \mathbf{p} \cdot \mathbf{p} \right) \mathbf{p} + \boldsymbol{\Omega} \cdot \mathbf{p} + a \left(\mathbf{D} \cdot \mathbf{p} - \mathbf{D} : \mathbf{p} \mathbf{p} \mathbf{p} \right) \\ & - \mathbf{v} \cdot \nabla \mathbf{p} - U_0 \left[\nabla \cdot \mathbf{p} \mathbf{p} + \frac{1}{3} \nabla (c - \mathbf{p} \cdot \mathbf{p}) \right] \end{aligned}$$

$$H = D_s (\nabla^2 c + \nabla \cdot (N_1 c \nabla c - \alpha \nabla \mathbf{p} \cdot \mathbf{p})) - \mathbf{v} \cdot \nabla c + U_0 \nabla \cdot \mathbf{p}.$$

With this setup of the equations, there is now 4 equations (one for the local concentration and one for each component of the polarity vector) that need to be solved in order to get the solution for the polarity vector and the local concentration. These elliptic equations are solved using the MUDPACK Multigrid solver (Adams 1989) at each time iteration to update the values of the polarity vector and local concentration.

4.2 Results From Numerical Simulations

Numerical simulations were run using the same parameter values presented above, and the results confirm the findings presented here. For positive values of ζ_a or values that are negative but close to zero, the system will be stable and will converge to this stable state. To present this graphically, we have chosen to plots of the velocity vector dotted with the polarity vector after nearly two million iterations of the solver were run. At this point, if $\vec{v} \cdot \vec{p}$ remained at some constant value, then the system is seen as stable. However, if this is not the case, then the system is seen as unstable. In Figure 4, an example of a stable state is shown while in Figure 5 an unstable state is shown.

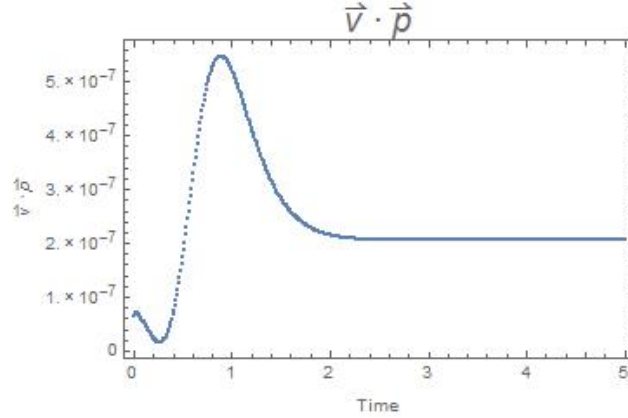


Figure 4: Stable state for the case where $\zeta_a = -1$.

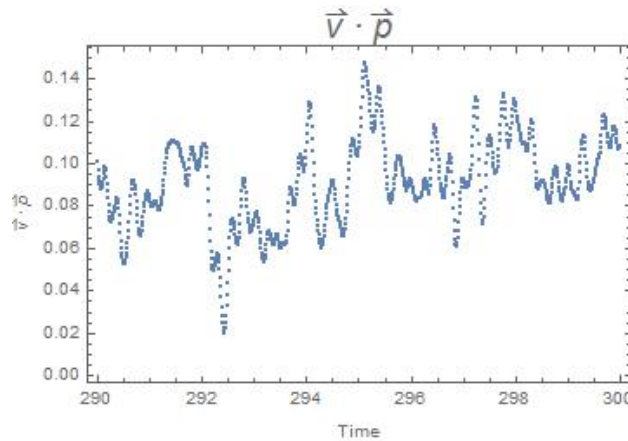


Figure 5: Unstable state for the case where $\zeta_a = -10$.

5 CONCLUSION

In this paper, we have presented our mathematical model for polar active liquid crystals which include the implementation of certain closure rules. Then we presented a method of linear stability analysis in order to examine the dynamics of the system near equilibrium. While doing this, we observed that for a specific case of the wave number k , there are stable and unstable states of the system based on the values of the particle activation parameter ζ_a and the polarity strength α . Next, we provided some results from the numerical simulations that were conducted which confirmed our results from the linear stability analysis. The goal of future work in this subject matter is to explore different values of the wave number k including

allowing $k_2, k_3 \neq 0$, as well as continue to conduct numerical simulations using different parameter values to determine the full effect of this (and numerical boundary conditions) on the stability of the liquid crystal system.

ACKNOWLEDGMENTS

This research is supported by the National Science Foundation, NSF-1517519. I would also like to acknowledge and thank the Modeling and Simulation Department at ODU for funding my work through a Modeling and Simulation assistantship.

A APPENDIX

The constants that were defined in 8, 9, and 10 are as follows:

$$\begin{aligned}
 S_1 &= i(\mathbf{v}_0 + U_0 \mathbf{p}_0) \cdot \mathbf{k} + D_s [|\mathbf{k}|^2 (1 - \frac{\alpha}{3}(c_0 - \|\mathbf{p}_0\|^2) + N_1 c_0) - \alpha(\mathbf{k} \cdot \mathbf{p}_0)^2] \\
 &\quad + 2D_r (1 - \frac{1}{3}\alpha c_0 + \frac{1}{3}\alpha \|\mathbf{p}_0\|^2) \\
 S_2 &= -\frac{i}{2}(1 + a)(\mathbf{k} \cdot \mathbf{p}_0) \\
 S_3 &= i(\mathbf{v}_0 \cdot \mathbf{k}) + D_s |\mathbf{k}|^2 (1 + N_1 c_0) \\
 S_4 &= U_0 - \alpha D_s (\mathbf{k} \cdot \mathbf{p}_0) \\
 S_5 &= iG(\mathbf{p}_0 \cdot \mathbf{k}) \left(\frac{1}{3}\alpha_0 - (1 + \zeta_a) - \frac{2}{3}\alpha_0 \|\mathbf{p}_0\|^2 \right) \\
 S_6 &= i(\mathbf{v}_0 \cdot \mathbf{k}) + \frac{1}{Re} |\mathbf{k}|^2 + \frac{1}{2Re_2} (\mathbf{p}_0 \cdot \mathbf{k})^2 - \frac{1}{3Re_2} \|\mathbf{p}_0\|^2 |\mathbf{k}|^2.
 \end{aligned}$$

The matrices A , B , C , and D from the system of equations are defined as:

$$\begin{aligned}
 A &= - \left(S_1 \mathbf{I} + iU_0 \mathbf{p}_0 \mathbf{k} - \frac{2i}{3} U_0 \mathbf{k} \mathbf{p}_0 + \frac{4}{3} \alpha D_r \mathbf{p}_0 \mathbf{p}_0 \right) \\
 B &= - \left(S_2 \mathbf{I} + \frac{i}{2} (1 - a) \mathbf{k} \mathbf{p}_0 + i a (\mathbf{k} \cdot \mathbf{p}_0) \mathbf{p}_0 \mathbf{p}_0 \right) \\
 C &= - \left[S_5 \mathbf{I} - \frac{2i}{3} \alpha_0 G \left((\|\mathbf{p}_0\|^2 - \frac{1}{2} \alpha_0) \mathbf{p}_0 \mathbf{k} + 2(\mathbf{p}_0 \cdot \mathbf{k}) \mathbf{p}_0 \mathbf{p}_0 \right) + \frac{iG}{3} (1 + \zeta_a) (2\mathbf{k} \mathbf{p}_0 - \mathbf{p}_0 \mathbf{k}) \right] \\
 D &= - \left[S_6 \mathbf{I} + \frac{1}{2Re_2} \left(|\mathbf{k}|^2 \mathbf{p}_0 \mathbf{p}_0 + (\mathbf{k} \cdot \mathbf{p}_0) (\mathbf{p}_0 \mathbf{k} + \mathbf{k} \mathbf{p}_0) \right) \right].
 \end{aligned}$$

REFERENCES

- Adams, J. C. 1989. “mudpack: Multigrid portable fortran software for the efficient solution of linear elliptic partial differential equations”. *Applied Mathematics and Computation* vol. 34, pp. 113–146.
- Darnton, N., L. Turner, K. Breuer, and H. C. Beng. 2004. “Moving Fluid with Bacterial Carpets”. *Biophysical Journal* vol. 86, pp. 1863–1870.

- Doi, M., and S. Edwards. 1986. *The Theory of Polymer Dynamics*. 1st ed. New York, New York, Oxford University Press.
- Dreyfus, R., J. Baudry, M. Roper, M. Ferminger, H. Stone, and J. Bibette. 2005. “Microscopic artificial swimmers”. *Nature* vol. 437, pp. 862–865.
- Gray, G., K. Harrison, and J. Nash. 2007. “New family of nematic liquid crystals for displays”. *Electronics Letters* vol. 9, pp. 130–131.
- Griebel, M., T. Dornsiefer, and W. Neunhoeffter. 1998. *Numerical Simulation in Fluid Dynamics*. 1st ed. Philadelphia, Society for Industrial and Applied Mathematics.
- Pedley, T. J., and J. O. Kessler. 1992. “Hydrodynamic Phenomena in suspensions of swimming microorganisms”. *Annual Review of Fluid Mechanics* vol. 24, pp. 313–358.
- Saintillan, D., and M. J. Shelley. 2008, Apr. “Instabilities and Pattern Formation in Active Particle Suspensions: Kinetic Theory and Continuum Simulations”. *Phys. Rev. Lett.* vol. 100, pp. 178103.
- Sanchez, D. A. 1968. *Ordinary Differential Equations and Stability Theory: An Introduction*. 1st ed. New York, New York, Dover Publications.
- Xiao-Gang, Y., M. G. Forest, and Q. Wang. 2014. “Near equilibrium dynamics and one-dimensional spatial-temporal structures of polar active liquid crystals”. *Chinese Physics B* vol. 23, pp. 1187011–1187037.
- Yang, X., M. G. Forest, W. Mullins, and Q. Wang. 2010. “2-D lid-driven cavity flow of nematic polymers: an unsteady sea of defects”. *Soft Matter* vol. 6, pp. 1138–1156.
- Zhang, S., and S. Kumar. 2008. “Carbon Nanotubes as Liquid Crystals”. *Small* vol. 4, pp. 1270–1283.

A SIMULATION AND VISUALIZATION APPROACH OF AIR POLLUTION IN CHINA

Sai Danganao and Yuzhong Shen
Modeling, Simulation & Visualization Engineering
Old Dominion University
ECS Building, Norfolk, VA 23529, USA
SDANG001@ODU.EDU, YSHEN@ODU.EDU

ABSTRACT

This paper proposes an approach for simulation and visualization of smog forming based on Weather Research and Forecasting Model with chemistry (WRF/Chem). WRF/Chem is one of the major forecast models with its online availability and consistency on the chemistry and meteorology processing. The numerical weather model was encapsulated to gain highly interactive capability that can be utilized by numerous applications such as educational software or local air quality forecast established on PC with augmented reality, and virtual reality. This extended abstract outlines the steps taken and displays initial achievements.

Keywords: Interactive, weather simulation, education.

1 INTRODUCTION

Since 2005, air pollution in major cities of China has become noticeable. The Chinese government enforced mandatory measures to maintain air quality for the 2008 Beijing Olympic Games. From 2012, smog in China's major areas caused panic and concern in press and public. In 2015, 80% of 366 cities in China failed to meet the national standard on air quality, even with the average PM 2.5 levels having fallen by 10%. In February 2015, "Under the Dome", a self-financed documentary film concerning air pollution in China, was viewed over 150 million times on Tencent video site within three days of its release.

It is common understanding that the certain sources, such as vehicles, industry, and biomass burning, are the causes of air pollution. However, it is not clear to the general public how these different sources contributed to smog formation. A tool that shows the pollution process will greatly facilitate understanding of air pollutions. Also, it is important to show people "what if scenarios" by changing the real smog process with the altered initial chemistry and meteorology conditions.

The interactive air pollution simulator presents an opportunity to vary different parameters for air pollution modeling and simulation.

2 METHOD

There are several build-in data sets, including weather observation data and emission inventory. An emission inventory is an estimation of air pollutant emissions by source in a certain area during a specific

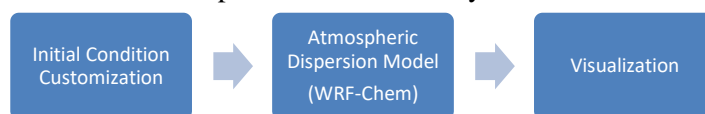


Figure 1 Workflow of the program

time. It is the most significant component of the initial input for the simulation. The build-in data are considered as the “real” data, which could be altered in simulations as the “imaginary” initial conditions.

In an interactive simulation, users are supposed to input certain variables, such as individual preferences, environmental factors, and different types of land use. Per the selection, a new emission inventory will be generated by multiplying scaling factors to the anthropogenic emissions grid.

Atmospheric dispersion model, in this case, WRF-Chem, is used to simulate distribution of air pollution. WRF/Chem model is ideally suited for air quality forecasting at regional to cloud resolving scales. It handles the chemistry and meteorology process simultaneously. The running environment of WRF-Chem is Unix. For software developing purpose, we developed WRF-Chem in Docker containers so it runs in windows (Hyper-V). It would take 5-30 minutes (or even more) for a personal computer to finish the computation, depending on the time and space resolution as well as the other configurations. The output file is in NETCDF format.

Eventually, based on the output of WRF-Chem, further visualization and visual analysis could be done. The pollution distribution could be visualized through various method, such as 3-D city street view with the effect of the smog, 2-D heat map and AQI-time graph etc. It is possible to save and load formatted .nc files. Comparison between different input and output are also possible for the educational purposes.

To improve user experience, it is reasonable to offer tutorials and case input with pre-computed output, to skip running the model and avoid waiting. Yet another application is to give users access to the real-time data, so the program will be running as the air quality forecast.

3 RESULTS AND CONCLUSION

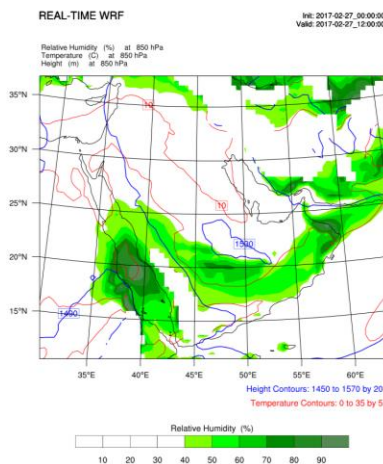


Figure 2 A test simulation (with emission inventory of sand storm but real time weather data) of middle-east asia

This abstract illustrates the initial implementation of a simulation and visualization approach of air quality. It could reveal basic relationships among multiple variables. The application is both scientific and intuitive, thus satisfy the demands of educational purposes. Further development of the software has the potential to become the air quality forecast running on PC.

For the project intentions, the Atmospheric Dispersion Model is done with WRF-Chem 3.8.1. Different emission inventories have been collected to establish the build-in data. A mid-way simulation can be made to test the parameter. Results could be visualized through NCAR Command Language (NCL) into .png format. Unity will be utilized to generate interactions and animations.

FAULT TOLERANCE FOR FINE-GRAINED PARALLEL ITERATIVE METHODS

Evan Coleman
Naval Surface Warfare Center - Dahlgren Division
ecole028@odu.edu

Masha Sosonkina
Old Dominion University
msosonki@odu.edu

ABSTRACT

This paper provides an introduction to fine-grained iterative methods and an overview of possible fault tolerance techniques for these methods. The computational model presented allows for either synchronous or asynchronous computation, and both the linear and non-linear cases are examined. Examples from recent research are provided and reviewed in the framework that is given.

Keywords: fault tolerance, fine-grained parallelism, asynchronous iteration, GPU/MIC acceleration

1 INTRODUCTION

Fine-grained parallel iterative methods are of great importance in the modern high-performance computing (HPC) environment. These algorithms can be run in either a synchronous or an asynchronous method and are naturally suited towards the memory/processing element style found in typical co-processor accelerators such as GPUs and MICs. As work using fine-grained parallel iterative methods continues to expand it is important to ensure that they are robust; i.e. that they are resilient to the occurrence of a computing fault and that they will be able to continue successfully to completion.

Fault tolerant methods are devised to increase both reliability and resiliency of HPC applications. There are reports (i.e. Asanovic, Bodik, Catanzaro, Gebis, Husbands, Keutzer, Patterson, Plishker, Shalf, Williams, et al. 2006, Cappello, Geist, Gropp, Kale, Kramer, and Snir 2014, Snir, Wisniewski, Abraham, Adve, Bagchi, Balaji, Belak, Bose, Cappello, Carlson, et al. 2014, Geist and Lucas 2009) that discuss the expected increase in the number of faults experienced by HPC environments. This is projected to be a more prevalent problem as HPC environments continue to evolve towards larger systems. As the landscape of HPC continues to grow into one where experiencing faults during computations is increasingly commonplace, the software used in HPC applications needs to continue changing alongside it in order to provide an additional measure of resilience against the increased number of faults experienced. As the rate of faults occurring in HPC environments continues to rise, it becomes increasingly important to ensure that these solvers are able to execute without suffering the negative consequences associated with a fault.

The next section of this paper provides an overview of related work and then Section 3 presents general mathematical models for asynchronous iteration in both the linear and non-linear cases including some general convergence theory. Section 4 gives results concerning fault tolerance of asynchronous iterative methods (with respect to both hard and soft faults), and details how the convergence results from Section 3 change. In Section 5, computational examples are presented for both the linear and non-linear case, and finally Section 6 concludes the paper.

2 RELATED WORK

Work on fine-grained (asynchronous) methods has been increasing as of late. Frommer and Szyld 2000 provides a survey of the theoretical results on asynchronous methods. A fine-grained parallel version of incomplete LU factorization was presented in Chow and Patel 2015, further analyzed in Chow, Anzt, and Dongarra 2015, and applied to model order reduction in Anzt, Chow, Saak, and Dongarra 2016. The general problem of asynchronously computing fixed points is discussed in Bertsekas 1983 and Bertsekas and Tsitsiklis 1989. Fault tolerance for asynchronous iterations has been explored in Anzt, Dongarra, and Quintana-Ortí 2015 in the linear case and Coleman, Sosonkina, and Chow 2017 in the non-linear case. Implementations of many asynchronous algorithms are provided for both MIC and GPU acceleration by the MAGMA library which is provided by Innovative Computing Lab 2015.

The expected increase in faults is detailed in Asanovic, Bodik, Catanzaro, Gebis, Husbands, Keutzer, Patterson, Plishker, Shalf, Williams, et al. 2006, Cappello, Geist, Gropp, Kale, Kramer, and Snir 2014, Snir, Wisniewski, Abraham, Adve, Bagchi, Balaji, Belak, Bose, Cappello, Carlson, et al. 2014, Geist and Lucas 2009, and many different approaches to fault tolerance have been investigated. Bronevetsky and de Supinski 2008 and Elliott, Mueller, Stoyanov, and Webster 2013 investigate traditional checkpoint-based fault tolerance in response to bit flips, while Elliott, Hoemmen, and Mueller 2014a, Bridges, Ferreira, Heroux, and Hoemmen 2012, Elliott, Hoemmen, and Mueller 2014b all investigate algorithmically based fault tolerance built upon the ideas of selective reliability. Alternative fault models – in order to generalize the negative impact that any fault may cause upon an algorithm – are studied and proposed in Elliott, Hoemmen, and Mueller 2015, Coleman and Sosonkina 2016a, Coleman and Sosonkina 2016b.

3 FINE-GRAINED (ASYNCHRONOUS) ITERATIVE METHODS

In fine-grained (parallel) computation, each component of the problem - i.e. a matrix or vector entry - is updated in a manner that does not require information from the computations involving other components while the update is being made. This allows for each computing element (i.e. a single processor or CUDA core) to act independently of all other computing elements while still being assigned multiple components to update. The generalized mathematical model that is used throughout this paper comes almost directly from Frommer and Szyld 2000, which in turn comes from both Baudet 1978 and Szyld 1998.

To keep this mathematical model as general as possible, consider a function, $G : D \rightarrow D$ where D is a product space $D = D_1 \times D_2 \times \dots \times D_m$. The goal is to find a fixed point of the function G inside of the domain D . To this end, iteration is performed such that $x^{k+1} = G(x^k)$ and a fixed point is declared if $x^{k+1} \approx x^k$. Note that the function G has internal component functions G_i for each sub-domain, D_i , in the product space, D . In particular,

$$x = (x_1, x_2, \dots, x_m) \in D \longrightarrow G(x) = G(x_1, x_2, \dots, x_m) = (G_1(x), G_2(x), \dots, G_m(x)) \in D \quad (1)$$

As a concrete example, let each $D_i = \mathbb{R}$. Forming the product space of each of these D_i 's gives that $D = \mathbb{R}^m$. This leads to the more formal functional mapping, $f : \mathbb{R}^m \rightarrow \mathbb{R}^m$. Additionally, let $f(\vec{x}) = 2\vec{x}$. In this case, each of the individual f_i component functions are defined by, $f_i(\vec{x}) = 2x_i$. Note that each of the component functions is said to operate on *all* of the vector \vec{x} even if the individual function definition does not require all of the components of \vec{x} to perform its specific update.

The assumption is also made that there is some finite amount of processing elements, P_1, P_2, \dots, P_p and that each of these is assigned to a block of components, B_1, B_2, \dots, B_m , to update. Note that the number of processing elements p will typically be significantly less than the number of blocks to update, m . The

restriction is made that $p < m$, however, often p will be significantly less than m . With these assumptions, the mathematical model can be stated in Algorithm 1.

Algorithm 1: Computational Model 1

```

1 for each processing element,  $P_k$  do
2   for  $i = 1, 2, \dots$  until convergence do
3     Read  $x$  from common memory
4     Compute  $x_j^{i+1} = G_j(x)$  for  $j \in B_k$ 
5     Update  $x_j$  in common memory with  $x_j^{i+1}$  for  $j \in B_k$ 

```

This computational model has each processing element read all pertinent data from global memory that is accessible by each of the processors, update the pieces of data specific to the component functions that it has been assigned, and update those components in the global memory. Note that the computational model presented in Algorithm 1 allows for either synchronous or asynchronous computation; it only specifies a method of fine-grained computing that allows each processing element to act independently of the other processors. If each processing element, P_k , were to wait for the other processors to finish each update, then the model would describe a parallel synchronous form of computation. However, the model is defined such that the processing elements do not have to wait for the other processors to finish which allows for an asynchronous form of computation.

To continue formalizing this computational model a few more definitions are necessary. First, set a global iteration counter, k , that increases *every* time any of the processor reads \vec{x} from common memory. At the end of the work done by any individual processor, p , the components associated with the block B_p will be updated. This results in a vector, $\vec{x} = (x_1^{s_1(k)}, x_2^{s_2(k)}, \dots, x_m^{s_m(k)})$ where the function $s_l(k)$ indicates how many times a specific component has been updated. Finally, a set of individual components can be grouped into a set, I^k , that contains all of the components that were updated on the k^{th} iteration. Given these basic definitions, the three following conditions (along with the model presented in Algorithm 1) provide a working mathematical framework for fine-grained asynchronous computation.

Definition 1. *If the following three conditions hold:*

1. $s_i(k) \leq k - 1$
Only components that have finished computing are used in the current approximation.
2. $\lim_{k \rightarrow \infty} s_i(k) = \infty$
The newest updates for each component are used.
3. $|k \in \mathbb{N} : i \in I^k| = \infty$
All components will continue to be updated.

Then given an initial $\vec{x}^0 \in D$, the iterative update process defined by,

$$x_i^k = \begin{cases} x_i^{k-1} & i \notin I_k \\ G_i(\vec{x}) & i \in I_k \end{cases}$$

where the function $G_i(\vec{x})$ uses the latest updates available is called an asynchronous iteration.

This basic computational model (i.e. Algorithm 1 and Definition 1) allows for many different results on fine-grained iterative methods that are both synchronous and asynchronous, though the three conditions given in Definition 1 are unnecessary in the synchronous case. In this work, the focus will be on fault tolerance for particular classes of linear and non-linear equations, with an emphasis on asynchronous methods. In certain (fine-grained) applications, the synchronous case is a special case of the asynchronous case.

3.1 Linear Methods

There is a significant amount of material in the current literature discussing fine-grained asynchronous linear methods (see for example: Frommer and Szyld 2000, El Tarazi 1982, Bertsekas and Tsitsiklis 1989, Baudet 1978), however the focus in this report will be on non-singular linear systems. These are described by the equation $Ax = b$ where $A \in \mathbb{R}^{n \times n}$ is non-singular. Further, consider stationary iteration methods that are defined by a splitting of A such that $A = M - N$ and the matrix M is also non-singular. The iteration operator associated with a given stationary iteration method is a mapping $H : \mathbb{R}^n \rightarrow \mathbb{R}^n$ such that $x^{k+1} = M^{-1}(Nx^k + b)$. The following result gives an indication of when the fine-grained asynchronous iteration will converge or diverge for a given problem.

Theorem 1. *Given an initial problem, $Ax = b$ and an iteration operator H as defined above the following statements hold:*

- *If $\rho(|H|) < 1$ then the asynchronous iteration defined in Section 3 converges to x^* where x^* is the solution of $Ax = b$.*
- *If $\rho(|H|) \geq 1$ then there exists some (asynchronous) order of updates to the components $x_i \in x$, and an initial guess x^0 such that the iterates produced by H do **not** converge to x^* .*

Two observations are worth making concerning this result. First, from Frommer and Szyld 2000, the set of matrices that have the property $\rho(|H|) < 1$ is the class of H -matrices (see Saad 2003); i.e. both M -matrices and matrices that are both strictly and irreducibly diagonally dominant. Second, computing the spectral radius of a given problem *may* be computationally prohibitive due to size and storage requirements, but does provide some means for ensuring that the given iteration is capable of converging in a fault free environment.

3.2 Non-linear Methods

A non-linear system is typically defined by a function $F(x) = 0, F : D \subseteq \mathbb{R}^n \rightarrow \mathbb{R}^n$ and the goal of an iterative procedure is to find a fixed point $x^* \in D$. The main result to be presented in this section is taken from Frommer and Szyld 2000. If H is an iterative version of the function F , then the following result holds.

Theorem 2. (El Tarazi) *Assume that $x^* \in D \setminus \bar{D}$ and that H is Frechét differentiable at x^* . If $\rho(|H'(x^*)|) < 1$ then there exists a neighborhood N of x^* such that the (fine-grained) asynchronous iterates (as defined by the assumptions in Section 3) converge to x^* provided that the initial guess $x^0 \in N$.*

Note that this result only guarantees convergence if the initial guess is inside of an (unknown) neighborhood of the solution x^* . In order to determine if the mapping will converge from its current location in the domain of the mapping H , it is necessary to first define what it means for a mapping to be a contraction.

Definition 2. *The function $H : D \subseteq \mathbb{R}^n \rightarrow \mathbb{R}^n$ is a contraction on D if there exists a constant $\alpha < 1$ such that:*

$$\|H(x) - H(y)\| \leq \alpha \|x - y\|$$

for some $x, y \in D$.

The form of the Jacobian $J(x) = H'(x)$, is determined by the ordering of the elements inside of x , but the norm of the Jacobian is associated only with the value of the elements in x . In particular, the spectral radius of the Jacobian is determined by the (partial) ordering imposed upon the mapping H when the Jacobian is computed, but the norm of the Jacobian changes as the iterative process progresses. The following result helps identify when the fixed-point iteration is a contraction:

Theorem 3. *The function $H : D \subseteq \mathbb{R}^n \rightarrow \mathbb{R}^n$ is a contraction at the location of the current iterate $x^k \in D$ if $\|H'(x^k)\| < 1$ for some matrix norm $\|\cdot\|$ and the domain $D \subseteq \mathbb{R}^n$ is convex.*

In the case where domain D is convex, this suggests an intuitive method for ensuring that the iterative process defined by H is progressing towards the solution x^* that will be explored further in Section 4. Note that $\rho(|H'(x^*)|) < 1$ only guarantees that there is a possible sequence of asynchronous updates that will cause an initial guess that is sufficiently close to the fixed point to converge; the condition $\|H'(x^k)\| < 1$ shows whether or not the current iterate will continue to progress towards the fixed point. These results can be combined to develop convergence statements for a given problem.

4 FAULT TOLERANCE FOR FINE-GRAINED ITERATIVE METHODS

This section deals with methods for developing fault tolerant variants of fine-grained iterative methods. Typically, faults are divided into two categories: hard and soft (e.g., Bridges, Ferreira, Heroux, and Hoemmen 2012, Elliott, Hoemmen, and Mueller 2014a). Hard faults cause immediate program interruption and typically come from negative effects on the physical hardware components of the computing platform or on the operating system itself. Soft faults represent all faults that do not cause the executing program to stop; they are the focus of this work. Most often, these faults refer to some form of data corruption that is occurring either directly inside of, or as a result of, the algorithm that is being executed. Currently, they often manifest as bit-flips. There are several fault mitigation that will be discussed in more detail below: migration, checkpointing, partial checkpointing, Algorithm-Based Fault Tolerance, self-stabilization.

Migration If the processing element, P_k , assigned to compute the updates for block B_k fails, the elements in B_k can be *migrated* to another block/processing element pair, B_l/P_l . Successful migration requires either: a flexible component assignment structure (i.e. block B_l can absorb all of the components in block B_k) or holding extra processing elements in reserve to protect against the occurrence of a fault.

Checkpointing During the course of the iterative update, the current iterate, x_k , is saved to memory. If a fault is detected, the corresponding iterate, x_F , is reset to the last known good state, x_k . This method is robust to the effects of a fault, but may be slow computational (especially at exascale: see Cappello, Geist, Gropp, Kale, Kramer, and Snir 2014 and Geist and Lucas 2009), and requires global communication which disrupts the fine-grained style of the iterative methods being discussed.

Partial Checkpointing As in the case of checkpointing, the current iterate, x_k , is saved to memory with some regularity. If a fault is detected, then some subset of the current iterate, x_F , is reset to the last known good state. Only the components that are determined to be affected by the fault need to be rolled back. This method requires a component level check on whether or not a fault has occurred. Note that this model has more synergy with fine-grained iterative methods; it is possible for individual components to detect faults and act accordingly. Further, the computational model of fine-grained (asynchronous) iteration given in Section 3 allows for the components to become out-of-sync with one another. In particular, (eventual) convergence of the fine-grained iterative is not affected negatively if the components in a given block, B_k , are reset to a state multiple iterations behind the other components that are being updated.

Algorithm-Based Fault Tolerance (ABFT) This category describes a class of methods that are based upon algorithmically modifying the fine-grained iterative method in question in order to mitigate the impact of a fault. These fault tolerance methods tend to be very application specific as they generally seek to avoid any sort of checkpointing, and instead encourage the iterative procedure to converge through the occurrence of a fault instead of retreating to a previously known good state.

Self-stabilization Self-stabilizing methods are a type of ABFT that are generally defined as methods that return a system to a valid state within some finite number of steps. As pointed out in Sao and Vuduc 2013,

the self-stabilizing property provides a means for fault tolerance; if a non-persistent fault occurs, the self-stabilizing method should be able to correct any impact of the fault in such a way that the algorithm will still converge. This technique tends to be application specific as it relies on correcting the values in the current iterate *without* saving a state to memory, and (if possible) avoiding an explicit fault detection mechanism.

4.1 Fault Recovery

Hard Faults Successful fault tolerance for asynchronous iterative methods with respect to hard faults relies upon two key components. First, successful detection of the fault itself. This will most likely be handled by the HPC platform, however the successful recovery of the fine-grained iterative method requires the algorithm to have the knowledge that a hard fault has occurred. This could be achieved internally in the algorithm by declaring a hard fault if components belonging to block B_F corresponding to processing element P_F fail to be updated within some stated bound. Second, the availability of extra processing elements, or the means to recover the failed hardware component quickly. Of the fault tolerance techniques described above, migration oriented methods are most well suited towards recovering from hard faults. Of the other methods, it is most easily possible to use the partial checkpointing method. In the case of migration, the components B_F associated with the failed processing element P_F will be reset to iteration 0, whereas in the case of partial checkpointing the elements in B_F will get reset to the last known good state; however there will also be an additional delay while the problem that affected P_F is remedied.

Soft Faults Unlike in the case of hard faults it is easier for recovery from soft faults to take a more algorithmically based approach, but harder to detect the fault directly. Each of the 5 results discussed above is capable of being utilized for the purposes of recovering from a soft fault, each with its own set of advantages and disadvantages. Similar to the case of a hard fault, the most important aspect to recovering from a soft fault is successful fault detection. However, this is often more difficult in the case of a soft fault since – though it corrupts data – it does not cause interruption to the flow of the iterative process. Many detection techniques rely on choosing an appropriate tolerance to check a property of the algorithm that has predictable behavior (e.g. a residual that is monotonically decreasing, a known property concerning a vector/matrix norm, etc); a tolerance that is too loose will allow potentially harmful errors to go undetected, while a tolerance that is too strict may report a fault when none actually occurred (“false positive”) and cause the program to do extra work to recover from a non-existent problem. The balance in choosing the correct fault tolerance method to recover from soft faults is typically application dependent and best done through the use of empirically derived information.

4.2 Convergence

Generally, the methods for recovering from *any* fault in the case of a fine-grained (asynchronous) iterative method are concerned with making the program in question robust to any negative numerical or fault-induced effects. Given the computational model presented in Section 3, it is important to note the convergence results presented in Section 3.1 and Section 3.2 will still hold naturally – assuming that the recovery process is executed in a reliable manner – for the following methods: migration, checkpointing, partial checkpointing. This is due to the asynchronous nature of the computational model presented in Section 3 allowing for certain elements to be updated significantly after others. If the components associated with a failed block B_F are not updated for some finite amount of time, the asynchronous fine-grained iterative algorithm will still converge so long as the components are *eventually* updated. This eventual update is guaranteed by each of the three methods listed. Convergence for both the general ABFT method and the specific self-stabilizing methodology tend to be both application and problem specific. Examples of such results are

presented in Anzt, Dongarra, and Quintana-Ortí 2015 (ABFT) and Coleman, Sosonkina, and Chow 2017 (self-stabilizing) for fine-grained iterative methods.

5 EXAMPLES

Examples of the fault tolerance methods discussed in Section 4 will be given for both the linear and non-linear case. This will show a possible instantiation of the techniques discussed above as an example of how to utilize the theory presented here.

5.1 Linear Method - Jacobi Iteration

The Jacobi iteration is an iterative method for determining a solution to the system of equations given by $Ax = b$. A brief introduction to the method is presented here before examining the modifications made to the algorithm for the purposes of fault tolerance. For further details on the Jacobi method and its derivation, see Saad 2003 or Bertsekas and Tsitsiklis 1989.

For a matrix $A \in \mathbb{R}^{n \times n}$, if each equation in the matrix-vector product is written individually, the system of equations becomes $\sum_{j=1}^n a_{ij}x_j = b_i$, where the values $a_{ij} \in A, x_j \in x$, and $b_i \in b$. If the assumption is made that $a_{ii} \neq 0$ then the a_{ii} terms can be removed from the sum above to obtain:

$$a_{ii}x_i + \sum_{j \neq i} a_{ij}x_j = b_i \longrightarrow x_i = \frac{-1}{a_{ii}} \left[\sum_{j \neq i} a_{ij}x_j - b_i \right] \quad (2)$$

for any given $i \in 1, 2, \dots, n$. The expression on the right lends itself towards iteration; with any initial guess for all of the $x_i \in x$ the individual x_i components can be updated repeatedly to arrive at a solution vector. Iteratively updating this expression represents the Jacobi algorithm. Note that the fine-grained nature of this method arises from the fact that each of the x_i components can be computed individually.

5.1.1 Fault Tolerance

Modifying the Jacobi algorithm to become fault tolerant can be done in two main ways – both of which rely on the fact that the residual at the k^{th} iteration, $r^k = b - Ax^k$, is less than the residual at the $(k-1)^{th}$ iteration, r^{k-1} (Saad 2003). This monotonic residual decay allows for both of the fault tolerance schemes that will be discussed here. The first variant is a simple checkpointing method. A copy of the vector x is stored with some regularity, and the residual r^k is computed at every iteration. If an increase in residual is detected than x is rolled back to the last known good state. Pseudo-code for this procedure is given in Algorithm 2.

Note that x_s is the version of x that is saved to rollback to and is computed every S iterations. This approach is very robust in that it will detect any fault or numerical instability that is large enough to cause an impact to the monotonicity of the residual (faults that are small enough to have no impact on the residual can typically be safely ignored), however computing the residual is computationally very expensive relative to the cost of a typical iteration. Additionally, since the computation of the residual itself has a chance to be affected negatively by a fault, this algorithm will not make much progress towards a solution when the computing environment has a high fault rate (Anzt, Dongarra, and Quintana-Ortí 2015).

The second approach to fault tolerance for the Jacobi method comes from Anzt, Dongarra, and Quintana-Ortí 2015 and functions on an individual component level. For each component, x_i , there is a constant, ϕ_i ($0 < \phi_i < 1$), such that, $|x_i^k - x_i^{k-1}| \leq \phi_i |x_i^{k-1} - x_i^{k-2}| \leq \phi_i^2 |x_i^{k-2} - x_i^{k-3}| \leq \dots$. Let $z_i^k = |x_i^k - x_i^{k-1}|$.

Algorithm 2: Checkpointing Jacobi Algorithm**Input:** Initial guess for $x_i \in x$, an input matrix A , and right-hand side b **Output:** The solution vector x to the equation $Ax = b$

```

1 for  $iter = 1, 2, \dots$  until convergence do
2   if  $iter == S$  then
3      $x_s = x$ 
4   for  $i = 1, 2, \dots, n$  do
5      $x_i = \frac{-1}{a_{ii}} [\sum_{j \neq i} a_{ij} x_j - b_i]$ 
6   Compute  $r^{iter} = b - Ax^{iter}$ 
7   if  $r^{iter} > r^{iter-1}$  then
8      $x = x_s$ 

```

If the problem in question has a linear convergence rate, than the component specific convergence ratio, $c_i = \frac{z_i^{k-1}}{z_i^k}$, can be used as a fault detector since the value for c_i should remain constant. In particular, one should compute component wise convergence ratio values for every element and use them in to detect faults throughout the algorithm. In particular, given a valid estimate of c_i , the following fault detector can be used,

$$\left| \frac{z_i^{k-1}}{z_i^k} - c_i \right| \leq c_i \cdot \delta \quad (3)$$

where δ is a user-defined threshold parameter. This leads an algorithm very similar to Algorithm 2 where the fault detection is replaced by Eq. (3) and instead of rolling back the entire vector x to a previous good state, the updates to individual components can be rejected on a case-by-case basis.¹

An important distinction between the two methods is that the second method preserves the fine-grained nature of the algorithm. This allows the algorithm to be executed in an asynchronous manner if desired, as there is no global communication between the individual components required to compute the traditional residual, $r^k = b - Ax^k$.

5.2 Non-linear Method - Fine Grained Parallel Incomplete LU Factorization

The non-linear example that will be considered is the Fine-Grained Parallel Incomplete LU (FGPILU) factorization. The FGPILU factorization approximates the true LU factorization and expresses a matrix A as the product of two factors L and U where, $A \approx LU$. Normally, the individual components of both L and U are computed in a manner that does not allow easy use of parallelization. The recent FGPILU algorithm proposed in Chow and Patel 2015 allows each element of both the L and U factors to be computed independently in a *fine-grained* manner. The algorithm progresses towards the incomplete LU factors that would be found by a traditional algorithm in an iterative process. To do this, the FGPILU algorithm makes use of the property $(LU)_{ij} = a_{ij}$ for all (i, j) in the sparsity pattern S of the matrix A , where $(LU)_{ij}$ represents the (i, j) entry of the product of the current iterate of the factors L and U . This leads to the observation that the FGPILU algorithm (given in Algorithm 3) is defined by the following two non-linear equations:

$$l_{ij} = \frac{1}{u_{jj}} \left(a_{ij} - \sum_{k=1}^{j-1} l_{ik} u_{kj} \right) \quad u_{ij} = a_{ij} - \sum_{k=1}^{i-1} l_{ik} u_{kj} \quad (4)$$

¹Due to space constraints, the pseudo-code for this algorithm is not provided

Following the analysis presented in Chow and Patel 2015, it is possible to collect all of the unknowns l_{ij} and u_{ij} into a single vector x , then express these equations as a fixed-point iteration $x^{(p+1)} = G(x^{(p)})$, where the function G implements the two non-linear equations described above. In a fault-free environment, it can be proven that the FGPILU algorithm is locally convergent in both the synchronous and asynchronous cases (see Section 3 in Chow and Patel 2015). The FGPILU algorithm is given in Algorithm 3.

Algorithm 3: FGPILU algorithm as given in Chow and Patel 2015

Input: Initial guesses for $l_{ij} \in L$ and $u_{ij} \in U$, an input matrix A

Output: Factors L and U such that $A \approx LU$

```

1 for  $sweep = 1, 2, \dots, m$  do
2   for  $(i, j) \in S$  do in parallel
3     if  $i > j$  then
4        $l_{ij} = (a_{ij} - \sum_{k=1}^{j-1} l_{ik}u_{kj})/u_{jj}$ 
5     else
6        $u_{ij} = a_{ij} - \sum_{k=1}^{i-1} l_{ik}u_{kj}$ 

```

5.2.1 Fault Tolerance

Similarly to Section 5.1.1, two different approaches towards fault tolerance for the FGPILU algorithm will be presented. Both of these approaches come from Coleman, Sosonkina, and Chow 2017. To track the progression of the FGPILU algorithm, it was proposed in Chow and Patel 2015 to monitor the non-linear residual norm. This is a value $\tau = \sum_{(i,j) \in S} \left| a_{ij} - \sum_{k=1}^{\min(i,j)} l_{ik}u_{kj} \right|$, which decreases as the number of sweeps progresses the algorithm closer to the conventional ILU factorization. If a fault occurs then one or both non-linear equations from the FGPILU algorithm will have some amount of error. In particular, if a fault occurs then the non-linear residual norm τ will be affected. In order to ensure that a fault does not negatively affect the outcome of the algorithm, a simple monitoring of the non-linear residual norm can be utilized. In principle, the non-linear residual norm will be at a minimum when the algorithm converges. Further, since there is a contribution from every non-zero (i, j) , the individual non-linear residual norms for each (i, j) , denoted here by τ_{ij} , can be defined as $\tau_{ij} = \left| a_{ij} - \sum_{k=1}^{\min(i,j)} l_{ik}u_{kj} \right|$. The total non-linear residual norm can always be recovered by taking the sum of all the individual non-linear residual norms over all non-zero (i, j) pairs. To establish a baseline for fault tolerance, define individual non-linear residual norms τ_{ij} for each (i, j) based on the initial guess that is used to seed the FGPILU algorithm. In particular, if L^* and U^* are the initial guesses for the incomplete L and U factors, then take $l_{ij}^* \in L$ and $u_{ij}^* \in U$ and define baseline individual non-linear residual norms τ_{ij}^* using the original values τ_{ij} and the values $l_{ij}^* \in L$ and $u_{ij}^* \in U$.

Since for each sweep of the FGPILU algorithm, the components $l_{ij} \in L$ and $u_{ij} \in U$ can be computed, by testing the individual non-linear residual norms it is possible to determine if a large fault occurred. Specifically, it is of interest to determine if a fault occurred that was large enough to cause a potential divergence of the algorithm. To do this, first a tolerance t is set and then a fault is signaled if $\tau_{ij} > t$. Since the individual non-linear residual norms are generally decreasing as the FGPILU algorithm progresses, set $t = \max(\tau_{ij}^*)$ initially (see Line 3 of Algorithm 4), and then update t during the course of the algorithm if desired. Note that if a fault is signaled by any of the individual non-linear residual norms, it is only known that a fault occurred somewhere in the current row of the factor L or the current column of the factor U . As such, the conservative approach would require the rollback of both the current row of L and the current column of U to their values at the previous checkpoint (e.g., Lines 5 to 6 of Algorithm 4). Further, it is possible for the individual non-linear residuals as defined to increase by a small amount, especially at early

iterations. To counteract the potential for reporting false positives on fault detection, the derivative of the global non-linear residual can be checked to ensure that it is also increasing before switching the current row and/or column (see Line 13 of Algorithm 4). This algorithm is detailed in Algorithm 4. Note that if

Algorithm 4: Partial Checkpoint-Based Fault Tolerant FGPIU

Input: Initial guesses for $l_{ij} \in L$ and $u_{ij} \in U$

Output: Factors L and U such that $A \approx LU$

```

1 for  $(i, j) \in S$  do in parallel
2    $\tau_{ij} = \left| a_{ij} - \sum_{k=1}^{\min(i,j)} l_{ik} u_{kj} \right|$ 
3  $t = \max(\tau_{ij})$ 
4 for  $sweep = 1, 2, \dots, m$  do
5   if Fault then
6      $i = \max_{i,j}(k_{ij}^1), j = \max_{i,j}(k_{ij}^2), Fault = FALSE, sweep = sweep - 1$ 
7     Rollback  $\{l_{ik}\}_{k=1}^{i-1}$  &  $\{u_{kj}\}_{k=1}^{j-1}$ 
8   else
9     for  $(i, j) \in S$  do in parallel
10      if  $i > j$  then  $l_{ij} = (a_{ij} - \sum_{k=1}^{j-1} l_{ik} u_{kj}) / u_{jj}$ 
11      else  $u_{ij} = a_{ij} - \sum_{k=1}^{i-1} l_{ik} u_{kj}$ 
12       $\tau_{ij} = \left| a_{ij} - \sum_{k=1}^{\min(i,j)} l_{ik} u_{kj} \right|$ 
13      if  $\tau_{ij} > t$  and  $\tau' > 0$  then
14        Set  $k_{ij}^1 = i, k_{ij}^2 = j$ , and Fault = TRUE

```

a fault is detected, the algorithm only restores the affected row of L and column of U . While no global communication is required to check for the presence of a fault, if a fault is detected there will be some communication required between processes to fix the effects of the fault.

The second method for fault tolerance for the FGPIU algorithm to be discussed here is a traditional “full” checkpointing method². In this case, a fault is declared if the currently computed global non-linear residual norm τ is some factor α greater than the previously computed non-linear residual norm τ_{i-1} . Note that, due to a combination of the asynchronous nature of the the FGPIU algorithm, the non-linear residual norm will not be strictly monotonically decreasing, especially as the algorithm proceeds closer to convergence. Therefore using the factor $\alpha = 1$, i.e., expecting a strict monotonic decrease, may cause the algorithm to report false positives, especially when nearing convergence.

The partial checkpointing method for fault tolerance of the FGPIU algorithm has the advantage of a fine-grained approach to fault detection, however it still requires global communication between the individual components in order to remedy the affect of a fault. It also does not require the entire matrix to be rolled back to a previous state; this may allow convergence to occur quicker, and/or may be able to be tuned to be more computationally efficient. In Coleman, Sosonkina, and Chow 2017 both methods were shown to be robust to the occurrence of faults in initial experiments.

6 SUMMARY & FUTURE WORK

This report has presented a survey of basic results concerning (asynchronous) fine-grained iterative methods, attempted to develop some general analytical statements about how fault tolerance methods can be used for this class of iterative method, and provided examples of these techniques from the recent research literature.

²Due to space constraints, the pseudo-code for this algorithm is not provided

Moving forward, it will be important to continue developing the general theory of fault tolerance for fine-grained iterative methods and applying it to specific applications. As these methods become increasingly more popular in the HPC environment due to their facility of use on both GPUs and MICs it is imperative to continue working towards applicable theoretical results that ensure they can execute successfully despite the possible negative effects associated with a computing fault.

ACKNOWLEDGMENTS

This work was supported in part by the Air Force Office of Scientific Research under the AFOSR award FA9550-12-1-0476, by the U.S. Department of Energy (DOE) Office of Advanced Scientific Computing Research under the grant DE-SC-0016564 and through the Ames Laboratory, operated by Iowa State University under contract No. DE-AC00-07CH11358, by the U.S. Department of Defense High Performance Computing Modernization Program, through a HASI grant, and the ILIR program at NSWC Dahlgren.

REFERENCES

- Anzt, H., E. Chow, J. Saak, and J. Dongarra. 2016. “Updating incomplete factorization preconditioners for model order reduction”. *Numerical Algorithms*, pp. 1–20.
- Anzt, H., J. Dongarra, and E. S. Quintana-Ortí. 2015. “Tuning stationary iterative solvers for fault resilience”. In *Proceedings of the 6th Workshop on Latest Advances in Scalable Algorithms for Large-Scale Systems*, pp. 1. ACM.
- Asanovic, K., R. Bodik, B. Catanzaro, J. Gebis, P. Husbands, K. Keutzer, D. Patterson, W. Plishker, J. Shalf, S. Williams et al. 2006. “The landscape of parallel computing research: A view from Berkeley”. Technical report, Technical Report UCB/EECS-2006-183, EECS Department, University of California, Berkeley.
- Baudet, G. M. 1978. “Asynchronous iterative methods for multiprocessors”. *Journal of the ACM (JACM)* vol. 25 (2), pp. 226–244.
- Bertsekas, D. P. 1983. “Distributed asynchronous computation of fixed points”. *Mathematical Programming* vol. 27 (1), pp. 107–120.
- Bertsekas, D. P., and J. N. Tsitsiklis. 1989. *Parallel and distributed computation: numerical methods*, Volume 23. Prentice hall Englewood Cliffs, NJ.
- Bridges, P., K. Ferreira, M. Heroux, and M. Hoemmen. 2012. “Fault-tolerant linear solvers via selective reliability”. *arXiv preprint arXiv:1206.1390*.
- Bronevetsky, G., and B. de Supinski. 2008. “Soft error vulnerability of iterative linear algebra methods”. In *Proceedings of the 22nd annual international conference on Supercomputing*, pp. 155–164. ACM.
- Cappello, F., A. Geist, W. Gropp, S. Kale, B. Kramer, and M. Snir. 2014. “Toward exascale resilience: 2014 update”. *Supercomputing frontiers and innovations* vol. 1 (1).
- Chow, E., H. Anzt, and J. Dongarra. 2015. “Asynchronous iterative algorithm for computing incomplete factorizations on GPUs”. In *International Conference on High Performance Computing*, pp. 1–16. Springer.
- Chow, E., and A. Patel. 2015. “Fine-grained parallel incomplete LU factorization”. *SIAM Journal on Scientific Computing* vol. 37 (2), pp. C169–C193.
- Coleman, E., and M. Sosonkina. 2016a. “A Comparison and Analysis of Soft-Fault Error Models using FGMRES”. In *Proceedings of the 6th annual Virginia Modeling, Simulation, and Analysis Center Capstone Conference*. Virginia Modeling, Simulation, and Analysis Center.

- Coleman, E., and M. Sosonkina. 2016b. "Evaluating a Persistent Soft Fault Model on Preconditioned Iterative Methods". In *Proceedings of the 22nd annual International Conference on Parallel and Distributed Processing Techniques and Applications*.
- Coleman, E., M. Sosonkina, and E. Chow. 2017. "Fault Tolerant Variants of the Fine-Grained Parallel Incomplete LU Factorization". In *Proceedings of the 2017 Spring Simulation Multiconference*. Society for Computer Simulation International.
- El Tarazi, M. N. 1982. "Some convergence results for asynchronous algorithms". *Numerische Mathematik* vol. 39 (3), pp. 325–340.
- Elliott, J., M. Hoemmen, and F. Mueller. 2014a. "Evaluating the impact of SDC on the GMRES iterative solver". In *Parallel and Distributed Processing Symposium, 2014 IEEE 28th International*, pp. 1193–1202. IEEE.
- Elliott, J., M. Hoemmen, and F. Mueller. 2014b. "Tolerating Silent Data Corruption in Opaque Preconditioners". *arXiv preprint arXiv:1404.5552*.
- Elliott, J., M. Hoemmen, and F. Mueller. 2015. "A Numerical Soft Fault Model for Iterative Linear Solvers". In *Proceedings of the 24th International Symposium on High-Performance Parallel and Distributed Computing*.
- Elliott, J., F. Mueller, M. Stoyanov, and C. Webster. 2013. "Quantifying the impact of single bit flips on floating point arithmetic". *preprint*.
- Frommer, A., and D. Szyld. 2000. "On asynchronous iterations". *Journal of computational and applied mathematics* vol. 123 (1), pp. 201–216.
- Geist, A., and R. Lucas. 2009. "Major computer science challenges at exascale". *International Journal of High Performance Computing Applications*.
- Innovative Computing Lab 2015. "Software distribution of MAGMA". <http://icl.cs.utk.edu/magma/>.
- Saad, Y. 2003. *Iterative methods for sparse linear systems*. Siam.
- Sao, P., and R. Vuduc. 2013. "Self-stabilizing iterative solvers". In *Proceedings of the Workshop on Latest Advances in Scalable Algorithms for Large-Scale Systems*, pp. 4. ACM.
- Snir, M., R. Wisniewski, J. Abraham, S. Adve, S. Bagchi, P. Balaji, J. Belak, P. Bose, F. Cappelto, B. Carlson et al. 2014. "Addressing failures in exascale computing". *International Journal of High Performance Computing Applications*.
- Szyld, D. B. 1998. "Different models of parallel asynchronous iterations with overlapping blocks". *Computational and applied mathematics* vol. 17, pp. 101–115.

AUTHOR BIOGRAPHIES

EVAN COLEMAN is a scientist with the Naval Surface Warfare Center Dahlgren Division. He holds an MS in Mathematics from Syracuse University and is working on a PhD in Modeling and Simulation from Old Dominion University. His email address is ecole028@odu.edu.

MASHA SOSONKINA is a professor of Modeling, Simulation and Visualization Engineering at Old Dominion University. Her research interests include high-performance computing, large-scale simulations, parallel numerical algorithms, and performance analysis. Her email address is msosonki@odu.edu.

A VECTOR INTRINSIC POINT-IMPLICIT LINEAR SOLVER FOR UNSTRUCTURED GRID APPLICATIONS ON INTEL XEON PHI KNIGHTS LANDING

Aaron Walden
Department of Computer Science
Old Dominion University
5115 Hampton Blvd, Norfolk, VA, USA
awalden@cs.odu.edu

ABSTRACT

We present a hand-optimized linear solver for unstructured fluid dynamics applications written in AVX512 vector intrinsics which replaces a legacy FORTRAN code optimized for multi-core CPUs. We discuss our strategies for vectorizing the given data layout, present our results, and discuss future optimization ideas.

Keywords: Computational Fluid Dynamics, Intel Xeon Phi Knights Landing, optimization, unstructured linear solvers, vectorization, code generation

1 INTRODUCTION

FUN3D is a collection of Computational Fluid Dynamics (CFD) codes developed and maintained by the NASA Langley Research Center which models fluid flow by solving the Navier-Stokes equations using an unstructured grid. Dominating the run time of most FUN3D simulations is a point-implicit linear solver routine which operates on a block-sparse matrix. FUN3D's point solver has been optimized for multi-core CPU performance (e.g. Intel Xeon) and it performs well “out of the box” on Xeon Phi Knights Landing (KNL) using only compiler-generated optimizations. It remains to be seen, however, if the machine’s full potential is being harnessed by FUN3D. In this work, we describe the implementation, optimization, and results of a point solver algorithm tailored for KNL written using *intrinsics*, a collection of C-style functions designed by Intel to give programmers register-level control of KNL's vector processing units (VPUs) .

2 METHODS

The inputs to the point solve routine are a block-sparse matrix with n block-rows of $nb \times nb$ dense blocks of off-diagonal and diagonal elements, **A_off** and **A_diag**, respectively, and vectors **b** and **r** (the residual) of size $n \times nb$. Sizes of n typically number in the millions. In this work we restrict the value of nb to 5, which is a common case. **A_off** is stored column-wise using the block CSR format. **A_diag** is stored linearly and column-wise in double precision. **A_off**, **b**, and **r** are stored in single precision.

The point solver routine solves $\mathbf{Ax} = \mathbf{b}$ first by computing the sparse matrix-vector product of a single input row of **A_off** and appropriate elements of **b**. The residual **r** minus the product becomes the right-hand side of forward and backward triangular solvers computed using the diagonal block of the input row, which is decomposed in-place into lower and upper triangular matrices.

A) *Matrix-vector product:* KNL's floating point performance is reliant on each core's dual SIMD VPUs which consist of 32 512-bit registers. We compute the matrix-vector product of **A_off** and **b** for row i by

looping over the k non-zero elements of $\mathbf{A_off}_i$ and computing $\mathbf{A_off}_{ij}\mathbf{b}_j$. We vectorize our code by loading 3 columns of $\mathbf{A_off}$, which is stored column-wise with unit stride, into a vector register, v_a . The final vector entry will be set to zero. Because $\mathbf{A_off}$ is not necessarily aligned to a 64-byte boundary, there may be a penalty associated with loading data which crosses cache lines. Aligning $\mathbf{A_off}$ by padding is unwise, however, because the code has a very low arithmetic intensity (apx. 0.5) and is certainly memory-bound on KNL. After loading v_a , corresponding values of \mathbf{b}_j are loaded into register v_b by copying them 5 times so that the correct column of $\mathbf{A_off}_{ij}$ will be multiplied by the same value of \mathbf{b}_j . We use a single vector register as an accumulator for the matrix-vector product, v_p , and compute $v_p = v_p + v_a v_b$ using a fused multiply-add to compute the product and sum it in a single vector instruction. We repeat this process for the last 2 columns of $\mathbf{A_off}_{ij}$. Following the loop, we are left with a single register of 15 values, which are partial sums of the matrix-vector product. We complete the product by adding the second and third 5-element lanes of v_p to the first.

B) Triangular solvers: Computation of this portion of the solver is performed in double precision, to which the matrix-vector product is upconverted before computation begins. Vectorization of 5x5 triangular solvers is less efficient due to data dependencies. We lack the space for full details, but briefly, the requirements are such that the first output, call it \mathbf{b}_1 , is required by the remaining 4, the second output is required by the remaining 3, and so on. Operations involving \mathbf{b}_i cannot be done until every operation involving \mathbf{b}_{i-1} has been completed. We can, however, vectorize over all operations involving \mathbf{b}_i once it has been computed. This allows us 4 SIMD operations in the first row, 3 in the second, and so on. To avoid undue cache pressure in a highly parallel environment, we load $\mathbf{A_diag}$ only once, hold it in 2 vector registers, and perform our computations using *blends* and *permutations* of $\mathbf{A_diag}$, which are instructions that allow us to shuffle and mix the contents of two registers. The final computed value of \mathbf{b}_i is downconverted to single precision before being stored.

C) Prefetching: We first fetch all data needed by the current row into L1 then fetch data for the next row into L2. This strategy gives ample time for latencies to be hidden behind computation and other prefetching.

3 RESULTS

Testing of our intrinsic solver conducted on the NASA Pleiades supercomputing system shows a total speedup of 1.7 over the original kernel for both small and large matrices (up to 5 million rows). A speedup of 1.3 is observed without prefetching. Prefetching increases the total speedup to 1.7. Combining L2 and L1 prefetching increases speedup by roughly 5% over L1 or L2 alone. Unrolling the matrix-vector product loop such that 3 blocks are done at once, giving nearly full vectorization and decreasing the number of instructions executed by nearly 20%, yields a paltry 2% speedup. Holding values of $\mathbf{A_diag}$ in registers and using permutations to compute gave no speedup over simply reloading data from L1.

4 CONCLUSION AND DISCUSSION

We achieved a significant speedup with our intrinsic solver despite using only basic optimizations. The fact that such large speedup is gained through prefetching exposes deficiency in KNL's hardware prefetchers, though this may be no surprise given the unstructured nature of the memory access pattern. The failure of unrolled matrix-vector loop to give speedup proportional to its drop in instruction count is a testament to how memory-bound the computation must be on KNL. In our future work, we can more aggressively optimize the solver using more sophisticated memory layouts, such as ELL or by organizing the diagonal blocks by column, which will triple vectorization in the triangular solvers. We also plan to expand intrinsic solvers beyond restriction to 5x5 blocks. A single kernel which can dynamically handle any block size is waylaid by integer bookkeeping, but we can maintain the speedup achieved by our intrinsic kernel by using code generation to write unrolled loops for a given block size which will compute vector position and appropriate masks at compile time.

MODELING TASK MIGRATION FOR FAULT TOLERANCE IN MATRIX-MATRIX MULTIPLICATION

Erik Jensen

Department of Modeling, Simulation, and
Visualization Engineering
Old Dominion University
4600 Elkhorn Avenue
Norfolk, Va, USA
Ejens005@odu.edu

Masha Sosonkina

Department of Modeling, Simulation, and
Visualization Engineering
Old Dominion University
4600 Elkhorn Avenue
Norfolk, Va, USA
Msosonki@odu.edu

ABSTRACT

Matrix-matrix multiplication kernels are subject to premature termination due to hardware failure. Using ULFM, a parallel MMM algorithm can recover from a fault and continue to completion with new processing units, or it can terminate. A mathematical model considers multiple factors to advise migration versus termination plus restarting. Results are forthcoming.

Keywords: High-Performance Computing, HPC, Fault Tolerance, Task, Migration, MPI, ULFM, Matrix-Matrix Multiplication, MMM

1 INTRODUCTION

Matrix-matrix multiplication (MMM) is one of the most common linear algebra kernels used in scientific and engineering applications. Parallel MMM algorithms such as those developed by Strassen and Cannon employ blocking to decompose the calculation into many smaller work units that can be distributed to multiple processing units, e.g. nodes within a supercomputer. Using message passing interface (MPI), a master processing unit can distribute the work units to worker processing units, which perform calculations and return the result to the master. A variant of OpenMPI, User Level Failure Mitigation (ULFM), permits an MPI application to continue beyond a fault and message delivery failure, at which point the application can recover from the fault. As hardware faults are projected to increase, a model is needed to determine the best course of action to take when a fault interrupts a MMM calculation: (1) migrate the failed task(s) from the faulty node to a new node, or (2) terminate the calculation and restart on new nodes. Several factors determine the more desirable route, including the problem size, proximity to completion, and the number of tasks pending on the failed node that would have to be migrated to reserve ranks on other nodes.

2 PRELIMINARY TESTING

A framework has been developed for testing task-based computations on a shared or distributed memory system. A naïve block matrix-matrix multiply algorithm is implemented within the framework to test the efficacy of the framework and the characteristics of the MMM algorithm, prior to developing a migration model. Preliminary testing experiments with several parameters: (1) matrix size, (2) block size, (3) number of processing units. Figures 1-4 depict results of various testing configurations. Time per task is defined as the total time divided by the number of tasks and does not reflect the actual time required to complete a task in the parallel computation. Tasks are completed in sets, which are distributed by the master. Task-set size is equal to the number of workers or remaining tasks. Another metric, time per task-set, may be defined as

$$T_{ts} = \frac{T_c}{\lceil t_n/w \rceil},$$

where T_c is the time to complete the MMM calculation, t_n is the number of tasks, and w is the number of workers. T_{ts} describes the amount of time to complete a set of tasks, and for large block sizes, is comparable to the time required for completion of a single task.

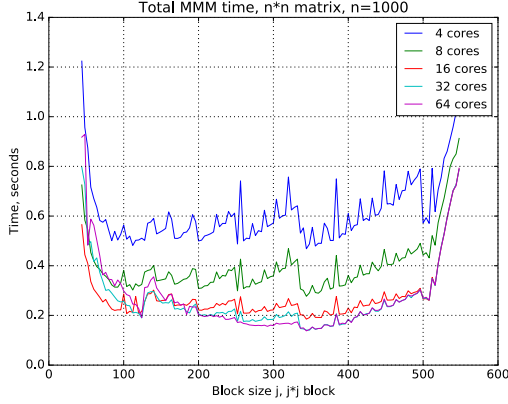


Figure 1: The framework and algorithm demonstrate scaling with increased cores.

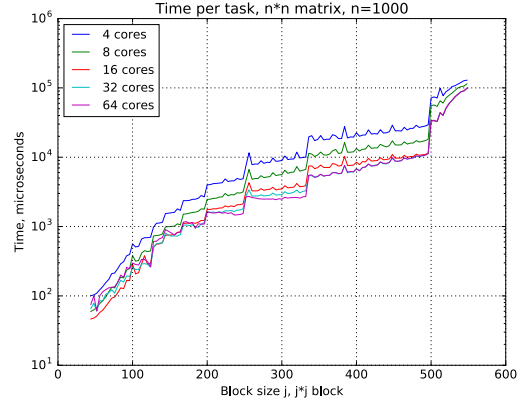


Figure 2: Wall time per task increases with block size and decreases with cores.

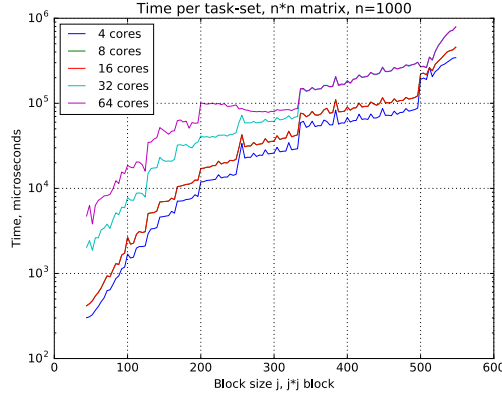


Figure 3: Larger sets increase latency.

3 FUTURE WORK

The framework and algorithm are subject to inspection and improvement. More efficient algorithms may be substituted for the naïve MMM algorithm. When the framework and MMM algorithm are considered acceptable, they will be tested in a fault-simulated environment. From that testing and analysis, a migration model will be developed.

4 ACKNOWLEDGEMENTS

This research was supported by the Turing High Performance Computing cluster at Old Dominion University.

PROTEIN SECONDARY STRUCTURE DETECTION USING PATTERN RECOGNITION AND MODELING

Tunazzina Islam
Ph.D. Student
Department of Computer Science
Old Dominion University
tislam@cs.odu.edu

ABSTRACT

Electron cryo-microscopy (Cryo-EM) technique produces density maps that are 3-dimensional (3D) images of molecules. In order to derive atomic structure of molecules, molecular features need to be identified from 3D images. Some molecular features of a protein show characteristic patterns in the image, and others show weak patterns or no pattern. We describe an approach that uses a combination of pattern recognition and geometrical modeling to recognize protein secondary features including α -helices and β -strands. We show the principle of modeling to distinguish the orientation of β -strands that are not visible in 3D images at medium resolution.

Keywords: Protein, β -strand, secondary structure, image, pattern, twist.

1 INTRODUCTION

Pattern recognition techniques have been successful in analysis of density maps obtained from cryo-electron microscopy (cryo-EM) technique. Cryo-EM is a biophysical technique to determine 3-dimensional structures of molecules [1, 2]. This technique is particularly suitable for large molecular assemblies that are often challenging for traditional techniques such as X-ray Crystallography and Nuclear Magnetic Resonance (NMR). A density map of molecules is a 3-dimensional image. When the resolution of the density map is higher than 4 Å, the quality of the image is sufficient to distinguish the protein chain and hence the molecular structure can be derived. However, for density maps at medium resolution, such as 5-10Å, the quality of the 3D image is not sufficient to distinguish the backbone of a protein. It is challenging to derive atomic structure from such images.

Though detailed molecular features are not visible for medium-resolution images, rough features such as secondary structures of a protein is visible. Secondary structure of a protein such as α -helices and β -sheets can be computational identified. An α -helix often appears as a cylinder and can be identified using image processing methods [3-6]. A β -sheet may appear as a thin layer of density and can be identified computationally [5-8].

Although β -sheets can be identified from cryo-EM density images at 5-10Å, it is almost impossible to detect the β -strands, the components of a β -sheet. The spacing between two neighboring β -strands is between 4.5 and 5Å, and therefore they are not visible when the resolution is at 5-10 Å. The detection of β -strands from β -sheets in such images has been a challenging problem since it was first attempted in 2004 [11].

A helix identified from the medium resolution cryo-EM image is often represented as a line referred as an α -trace that corresponds to the central axis of a helix. Location of major β -sheets can be identified from

image in low threshold value. But it is not possible to identify β -strands because they don't have any fixed pattern. In low threshold they may be visible but in high threshold they are not.

2 METHOD

2.1 Patterns of helices and β -sheets and detection

In a medium-resolution density map, a helix appears as a cylinder, and various methods exist to detect the location of these helices. We applied SSETracer [10] to detect the location of helices in a density map. SSETracer detects helices (and β -sheets) based on a characterization of local density features such as local structure tensor, local thickness, continuity of the skeleton, and density value. A detected helix is represented by a set of points located along the central axis of the helix. The current implementation of SSETracer contains a modified step in the axis extension to enhance the geometric characterization of the helix

2.2 Modeling of β -strands

A β -sheet is composed of multiple β -strands those can be parallel, antiparallel or mix of both. The sheet twist is defined as the angle between the backbone vectors of the two residues in the pair [12]. We have measured the twist angle from 3D atom coordinates extracted from PDB file generated from image having different orientations including best case that is similar to true orientation and bad case that is far from true orientation. .

3 RESULT

We used six proteins, for which the atomic structures were downloaded from the PDB, and their corresponding 3D density maps were simulated at 10 Å resolution using Chimera. For good orientation we have larger twist (3rd column of Table 1) than bad orientation (4th column of Table 1) because good orientation matches more with the atomic structure.

PDB ID_SHEET ID	No. of β -Strands in a β -Sheet	Maximum twist (good case)	Maximum twist (bad case)
1A12_A	4	13.526	8.954
1AKY_A	5	15.285	13.068
1ATZ_A	6	13.071	9.674
1CHD_SH1	7	10.93	5.736
1DTD_A	8	12.524	6.664
1QNA_C	9	13.415	3.501

4 REFERENCES

- [1] W. Chiu, M. L. Baker, W. Jiang, M. Dougherty, and M. F. Schmid, "Electron cryomicroscopy of biological machines at subnanometer resolution," *Structure*, vol. 13, pp. 363-72, Mar 2005.
- [2] Z. H. Zhou, "Atomic resolution cryo electron microscopy of macromolecular complexes," *Adv Protein Chem Struct Biol*, vol. 82, pp. 1-35, 2011.
- [3] W. Jiang, M. L. Baker, S. J. Ludtke, and W. Chiu, "Bridging the information gap: computational tools for intermediate resolution structure interpretation," *J Mol Biol*, vol. 308, pp. 1033-44, May 2001.
- [4] A. Del Palu, J. He, E. Pontelli, and Y. Lu, "Identification of Alpha-Helices from Low Resolution Protein Density Maps," presented at the Proceeding of Computational Systems Bioinformatics Conference(CSB), 2006.
- [5] M. L. Baker, T. Ju, and W. Chiu, "Identification of secondary structure elements in intermediate-resolution density maps," *Structure*, vol. 15, pp. 7-19, Jan 2007.
- [6] D. Si, S. Ji, K. A. Nasr, and J. He, "A machine learning approach for the identification of protein secondary structure elements from electron cryo-microscopy density maps," *Biopolymers*, vol. 97, pp. 698-708, Sep 2012.
- [7] Y. Kong and J. Ma, "A structural-informatics approach for mining beta-sheets: locating sheets in intermediate-resolution density maps," *J Mol Biol*, vol. 332, pp. 399-413, Sep 12 2003.
- [8] S. Dong and H. Jing, "Beta-sheet Detection and Representation from Medium Resolution Cryo-EM Density Maps," in *Proceedings of the International Conference on Bioinformatics, Computational Biology and Biomedical Informatics %@ 978-1-4503-2434-2*, ed. Wshington DC, USA: ACM, 2013, pp. 764-770.
- [9] D. Si and J. He, "Beta-sheet Detection and Representation from Medium Resolution Cryo-EM Density Maps.," *BCB'13: Proceedings of ACM Conference on Bioinformatics, Computational Biology and Biomedical Informatics*, vol. Washington, D.C., pp. 764-70, 2013.
- [10] K. Al Nasr, C. Liu, M. Rwebangira, L. Burge, and J. He, "Intensity-based skeletonization of CryoEM gray-scale images using a true segmentation-free algorithm," *IEEE/ACM Trans Comput Biol Bioinform*, vol. 10, pp. 1289-98, Sep-Oct 2013.
- [11] M. L. Baker, M. R. Baker, C. F. Hryc, T. Ju, and W. Chiu, "Gorgon and pathwalking: macromolecular modeling tools for subnanometer resolution density maps," *Biopolymers*, vol. 97, pp. 655-68, Sep 2012.
- [12] B. K. Ho and P. M. G. Curmi, "Twist and Shear in β -Sheets and β -Ribbons," DOI= 10.1006/jmbi.2001.5385

SIMULATION OF SOUND ABSORPTION BY SCATTERING BODIES TREATED WITH ACOUSTIC LINERS AND THE ASSESSMENT OF ITS HIGH-PERFORMANCE PARALLEL COMPUTING CAPABILITIES

Michelle E. Pizzo and Fang Q. Hu

Old Dominion University
Department of Mathematics and Statistics
2300 Engineering and Computational Sciences Building
Norfolk, VA, 23529, USA
mpizzo@odu.edu, fhu@odu.edu

ABSTRACT

Aircraft noise reduction is a major goal within the field of aviation research. When designing next generation quiet aircraft, it is important to be able to accurately and efficiently predict the acoustic scattering by an aircraft body from a given noise source. Acoustically treated liners are an effective tool for aircraft noise reduction. The objective of this work is to study the modeling and simulation of acoustic wave scattering by prototype geometric bodies treated with acoustic liners. Moreover, fast algorithms and high performance computing are considering to both assess and reduce the computational cost of the simulation. Carried out by large scale parallel computing, the scalability, speedup, and efficiency of the acoustic simulation will be studied. Preliminary results are included for modeling the acoustic wave scattering by non-treated geometries from a given point source. Geometries include a flat-plate and sphere.

Keywords: aeroacoustics, wave scattering, high performance computing, acoustic liners, admittance

1 INTRODUCTION

Aircraft noise reduction is a major goal within the field of aviation research. When designing next generation quiet aircraft, it is important to be able to accurately and efficiently predict the acoustic scattering by an aircraft body, both rigid as well as lined, from a given noise source (Hu 2013, Hu 2014, Hu et al. 2016). Acoustic scattering problems can be modeled using boundary element methods (BEMs) by reformulating the linear convective wave equation as a boundary integral equation (BIE), both in the frequency domain and the time domain (Chappell et al. 2006, Ergin et al. 1999, Jones and Hu 2007, Marburg 2015, Marburg and Schneider 2001).

Although frequency domain solvers are the most commonly used and researched within literature, there are several distinct advantages (Hu 2013) to using a time domain solver. For example, a time domain solution allows for the simulation and study of broadband sources and time-dependent transient signals, scattering solutions at all frequencies to be obtained within a single computation, and the avoidance of needing to invert a large dense linear system (as is required in the frequency domain). Additionally, a time domain solution is more naturally coupled with a nonlinear computational fluid dynamics simulation of noise sources.

Time domain boundary integral equations (TDBIE) have an intrinsic numerical instability and carry a high computational cost. In recent years, numerical techniques for modeling acoustic wave scattering by complex

geometries using the TDBIE have been under development (Hu 2013, Hu 2014, Hu et al. 2016). It has been shown that stability can be realized using a Burton-Miller type reformulation of the BIE. Moreover, the computational cost can be reduced using fast algorithms and high performance computing (HPC). This work uses a time domain BEM (TDBEM) approach, in which the scattering solution is obtained using temporal and surface basis functions and a March-On-in-Time scheme in which a sparse matrix is solved iteratively. The March-On-in-Time scheme reduces computational complexity.

The numerical formulation of the stable TDBEM is detailed in Section 2. The formulation is valid for both rigid and lined acoustic surfaces with a mean flow \mathbf{U} . The incorporation of an impedance boundary condition, valid for surfaces with acoustically treated liners, is detailed in Section 3. Preliminary results are included in Section 4 for modeling the acoustic wave scattering by non-treated geometries, i.e. geometries with a rigid surface, from a given point source. Concluding remarks are given in Section 5, and future work is detailed in Section 6.

2 FORMULATION OF THE STABLE TIME DOMAIN BOUNDARY ELEMENT METHOD

The convective wave equation with constant mean flow \mathbf{U} is given by:

$$\left(\frac{\partial}{\partial t} + \mathbf{U} \cdot \nabla\right)^2 p - c^2 \nabla^2 p = q(\mathbf{r}, t) \quad (1)$$

where $q(\mathbf{r}, t)$ is the source term, $p(\mathbf{r}, t)$ is the acoustic pressure, and c is the speed of sound. Consider a homogeneous initial condition $p(\mathbf{r}, 0) = \partial p / \partial t(\mathbf{r}, 0) = 0$. By introducing free-space adjoint Green's function:

$$\begin{aligned} \tilde{G}(\mathbf{r}, t; \mathbf{r}', t') &= \frac{1}{4\pi c^2 \bar{R}} \delta\left(t' - t + \beta \cdot (\mathbf{r}' - \mathbf{r}) - \frac{\bar{R}}{c\alpha^2}\right) \\ \alpha &= \sqrt{1 - M^2}, \quad \beta = \frac{\mathbf{M}}{1 - M^2}, \quad \mathbf{M} = \frac{\mathbf{U}}{c} \\ \bar{R} &= \sqrt{|\mathbf{M} \cdot (\mathbf{r} - \mathbf{r}')|^2 + \alpha^2 |\mathbf{r} - \mathbf{r}'|^2} \end{aligned} \quad (2)$$

where $M = |\mathbf{M}|$ and \mathbf{r} is a point on scattering body surface, the wave propagation problem can be reformulated into a TDBIE.

To reformulate the wave propagation problem into a TDBIE, we consider the Kirchhoff integral representation of the acoustic field in the presence of a mean flow:

$$p(\mathbf{r}', t') = \int_V \frac{q(\mathbf{r}, t'_R)}{4\pi c^2 \bar{R}} d\mathbf{r} + \int_0^{t'^+} \int_S \left[c^2 \left(\tilde{G} \frac{\partial p}{\partial \bar{n}} - p \frac{\partial \tilde{G}}{\partial \bar{n}} \right) - U_n \left(\tilde{G} \frac{\partial p}{\partial t} - p \frac{\partial \tilde{G}}{\partial t} \right) \right] d\mathbf{r}_S dt \quad (3)$$

where $\partial / \partial \bar{n} = \partial / \partial n - M_n(\mathbf{M} \cdot \nabla) = (\mathbf{n} - M_n \mathbf{M}) \cdot \nabla$ denotes a modified normal derivative, $\bar{\mathbf{n}} = \mathbf{n} - M_n \mathbf{M}$ and $M_n = \mathbf{M} \cdot \mathbf{n}$ such that \mathbf{n} is the inward normal vector on the scattering body, $t'_R = t' + \beta \cdot (\mathbf{r}' - \mathbf{r}) - \bar{R}/(c\alpha^2)$, and S denotes the surfaces of both the scattering body and acoustic source. Figure 1 illustrates the relationship between the mean flow, the surface of the scattering body, and the surface of the acoustic source.

In Eq. (3), the left integral is representative of direct sound and the right integral is representative of the sound contribution from the surface of both the scattering body and acoustic source. The TDBIE results by taking the limit of Eq. (3) as $\mathbf{r}' \rightarrow \mathbf{r}'_S$ where \mathbf{r}'_S is a point on the boundary and \mathbf{r}_S is some arbitrary point:

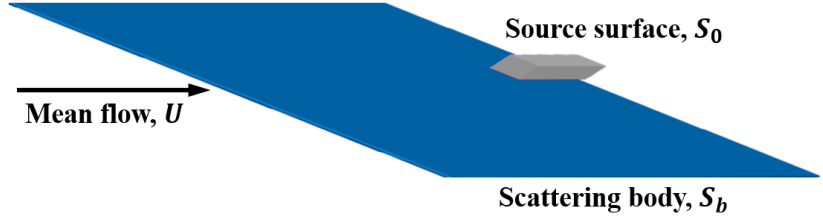


Figure 1: Schematic diagram illustrating the relationship between the mean flow, the surface of the scattering body, and the surface of the acoustic source.

$$C_S p(\mathbf{r}'_S, t') = \int_V \frac{q(\mathbf{r}, t'_R)}{4\pi c^2 R} d\mathbf{r} + \int_0^{t'^+} \int_S \left[c^2 \left(\tilde{G} \frac{\partial p}{\partial \bar{n}} - p \frac{\partial \tilde{G}}{\partial \bar{n}} \right) - U_n \left(\tilde{G} \frac{\partial p}{\partial t} - p \frac{\partial \tilde{G}}{\partial t} \right) \right] d\mathbf{r}_S dt \quad (4)$$

where $C_S = 1$ if \mathbf{r}'_S is on the exterior of S and $C_S = 1/2$ if \mathbf{r}'_S is on S . Using the relation that $\partial/\partial \bar{n} = \partial/\partial n - M_n(\mathbf{M} \cdot \nabla) = (\mathbf{n} - M_n \mathbf{M}) \cdot \nabla$ in Eq. (4), $\partial p/\partial \bar{n}$ can be expressed with the term $\partial p/\partial n$. For rigid acoustic surfaces, $\partial p/\partial n = 0$.

The resulting TDBIE has intrinsic numerical instabilities due to resonant frequencies resulting from nontrivial solutions in the interior domain. Using a Burton-Miller type reformulation of Eq. (4), resonant frequencies can be eliminated and stability achieved (Hu 2013, Hu 2014, Hu et al. 2016). The stable Burton-Miller type reformulation of Eq. (4) is discretized by dividing S into boundary elements using surface element basis functions $\phi_i(\mathbf{r}_S)$ at node i and temporal basis functions $\psi_j(t)$ at time j :

$$p(\mathbf{r}_S, t) = \sum_{i=1}^{N_e} \sum_{j=0}^{N_t} u_i^j \psi_j(t) \phi_i(\mathbf{r}_S) \quad (5)$$

In Eq. (5), N_e denotes the total number of surface nodes, N_t denotes the number of time-steps, and u_i^j denotes the solution at node i and time-step j . Equation (5) yields the collocation March-On-in-Time scheme:

$$\mathbf{B}_0 \mathbf{u}^j = \mathbf{q}^j - \sum_{m=1}^{j_M} \mathbf{B}_m \mathbf{u}^{j-m} \implies \mathbf{u}^j = \lambda^j \mathbf{e}_0 \quad (6)$$

which reduces computational complexity. In Eq. (6), \mathbf{u}^j denotes the vector of all coefficients u at time j . The solution is obtained by iteratively solving the system of equations. In the system (6), \mathbf{B}_m , $m = 0, \dots, j_M$ is sparse. Stability is realized if and only if the magnitude of the largest eigenvalue of \mathbf{B}_m is less than or equal to 1, i.e., $\max|\lambda| \leq 1$.

3 INCORPORATION OF AN IMPEDANCE BOUNDARY CONDITION

Consider a geometric body with a surface that is treated with an acoustic liner. In this case, $\partial p/\partial n \neq 0$ unlike rigid acoustic surfaces where $\partial p/\partial n = 0$. Without the presence of mean flow, $\partial p/\partial n$ can be represented (Marburg and Schneider 2001) in the frequency domain by:

$$P_n(\mathbf{x}, \omega) = \frac{\partial p}{\partial n}(\mathbf{x}, \omega) = i\omega \rho_0 Y(\omega) p(\mathbf{x}, \omega) \quad (7)$$

where ρ_0 is the density, $Y(\omega) = 1/Z(\omega)$ is the surface admittance, and $Z(\omega)$ is the surface impedance. It can be shown that (refer to Appendix A):

$$P_n(\mathbf{x}, t) = \frac{\partial p}{\partial n}(\mathbf{x}, t) = \frac{\rho_0}{2\pi} \int_{-\infty}^t \frac{\partial y}{\partial t}(t - \tau) p(\mathbf{x}, \tau) d\tau \quad (8)$$

where $\partial y / \partial t$ is the time derivative of $y(t)$, the inverse Fourier transform of $Y(\omega)$.

The ability to translate the acoustic liner admittance from the frequency domain to the time domain is heavily motivated by the work of Rienstra (Rienstra 2006), in which a frequency domain impedance boundary condition is translated to the time domain using Fourier transforms. Using experimentally determined admittance data and methodology similar to Rienstra (Rienstra 2006), a time domain solution of $y(t)$ can be calculated. Then, Eq. (8) can be coupled with the BIE that provides the framework for accurately calculating sound scattering from acoustically large bodies. This coupling will allow for the study of sound absorption by acoustic liners regarding both the spatial resolution of the TDBEM with respect to the surface element basis functions and the far field acoustic impedance behaviors of treated vs. non-treated surfaces. Additionally, this coupling will allow for the study of the scalability and computational performance by assessing the algorithmic speedup and efficiency. Refer to Section 6 for further details.

4 PRELIMINARY RESULTS FOR THE ASSESSMENT OF RIGID SURFACES

Preliminary work has been completed to assess the spatial resolution of the TDBEM with respect to the surface element basis functions $\phi_i(\mathbf{r}_s)$ in Eq. (5) by considering the scattering of an acoustic point source by prototype geometric bodies with rigid surfaces (Hu et al. 2016). Geometries considered include a flat plate and sphere. The flat plate has dimension $[-0.5, 0.5] \times [-0.5, 0.5] \times [-0.1, 0.1]$ and the point source is located at $\mathbf{x} = (x, y, z) = (0, 0, 1)$. The sphere is centered at $\mathbf{x} = (0, 0, 0)$ with radius 0.5 and the point source is located at $\mathbf{x} = (0, 0, 1)$. For the spherical geometry, we consider both a standard and rotated orientation as shown in Fig. 2a and Fig. 2b. For both the flat plate and sphere, we limit the study to the observed far field solution.

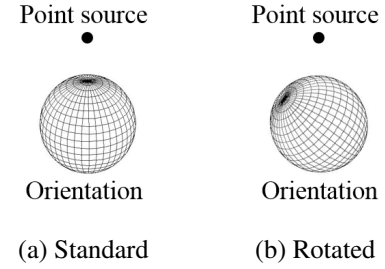


Figure 2: Orientation of the rigid body spherical elements with relation to the point source.

4.1 Flat Plate with Rigid Body

To assess the spatial resolution of the TDBEM with respect to the spatial basis functions, we consider the scattering of an acoustic point source by the flat plate with rigid body. The contour plots of the frequency domain solutions converted from the time domain solution at $\omega = 5\pi$ and $\omega = 15\pi$ are given in Fig. 3. The solution in Fig. 3 was computed using 80 elements in both the x - and y - directions and 16 elements in the z -direction.

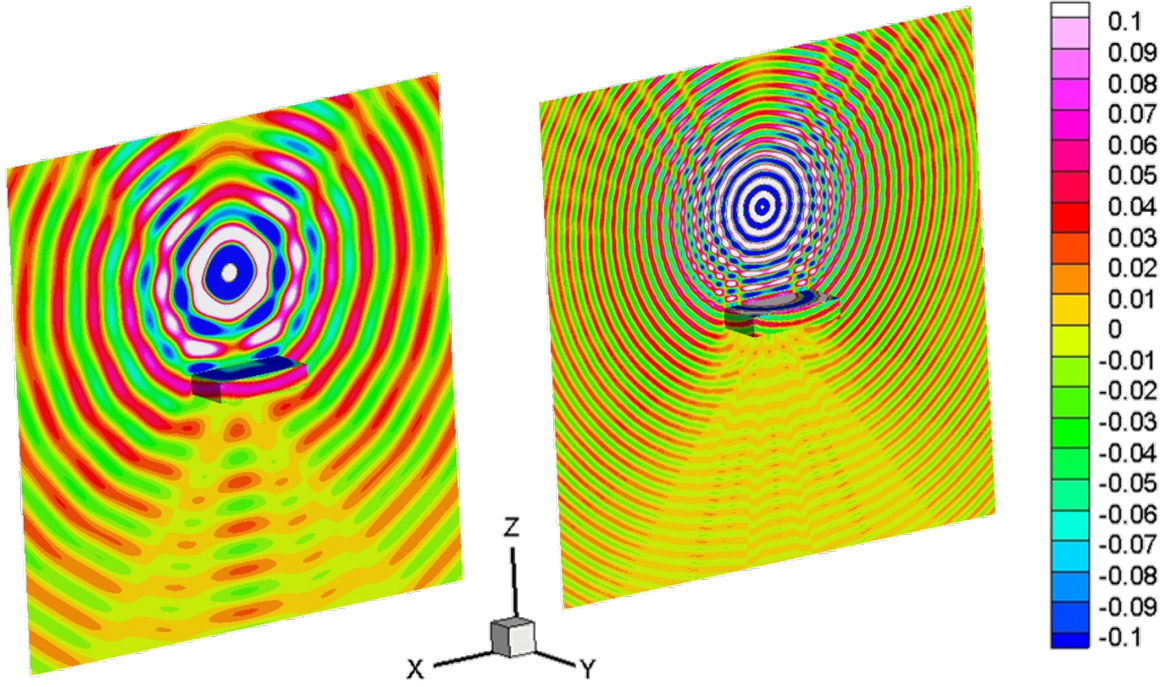


Figure 3: Rigid body flat plate contour plots of the frequency domain solution converted from the time domain solution at $\omega = 5\pi$ [left] and $\omega = 15\pi$ [right].

By discretizing the surface of the plate in the x , y , and z directions with N_x , N_y , and N_z elements respectively, a series of computations were carried out by increasing the number of elements used from $N_x \times N_y \times N_z = 20 \times 20 \times 4$ to $80 \times 80 \times 16$. Figure 4 depicts the converted frequency domain solution at $\omega = 15\pi$ along a field line of coordinates where $-2.5 \leq x \leq 2.5$, $y = 0$, and $z = -2.5$, as the number of surface elements increased from $20 \times 20 \times 4$ to $80 \times 80 \times 16$. The spatial resolution is measured along one direction on the surface using the metric of points-per-wavelength (PPW). PPW is computed by:

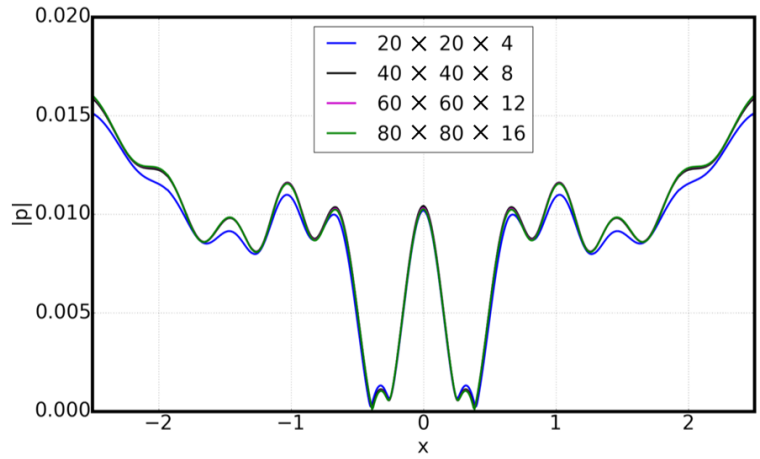


Figure 4: Rigid body flat plate frequency domain solution converted from the time domain solution at $\omega = 15\pi$ along a field line of coordinates $-2.5 \leq x \leq 2.5$, $y = 0$, and $z = -2.5$.

$$PPW = \frac{2\pi(p+1)N_x}{kL_x} \quad (9)$$

where p is the order of the basis function, $k = \omega/c$ is the wavenumber, L_x is the plate length along the x -direction, and N_x is the number of elements along the x -direction (Hu et al. 2016). Using the solution

computed by $100 \times 100 \times 20$ as the reference, the relative error in the L2 norm is plotted as a function of PPW. The results of the far field solution with constant basis functions are given in Fig. 5. For constant basis functions, $p = 0$. As indicated in the figure, the relative error measured in the L2 norm becomes as small as 2% with only 5 PPW. The excellent spatial resolution is likely due to the integration over a closed, hence periodic, domain. Although the basis functions being used are of zeroth order, the integrations over each element are computed by high-order Gauss quadrature on a 6×6 grid.

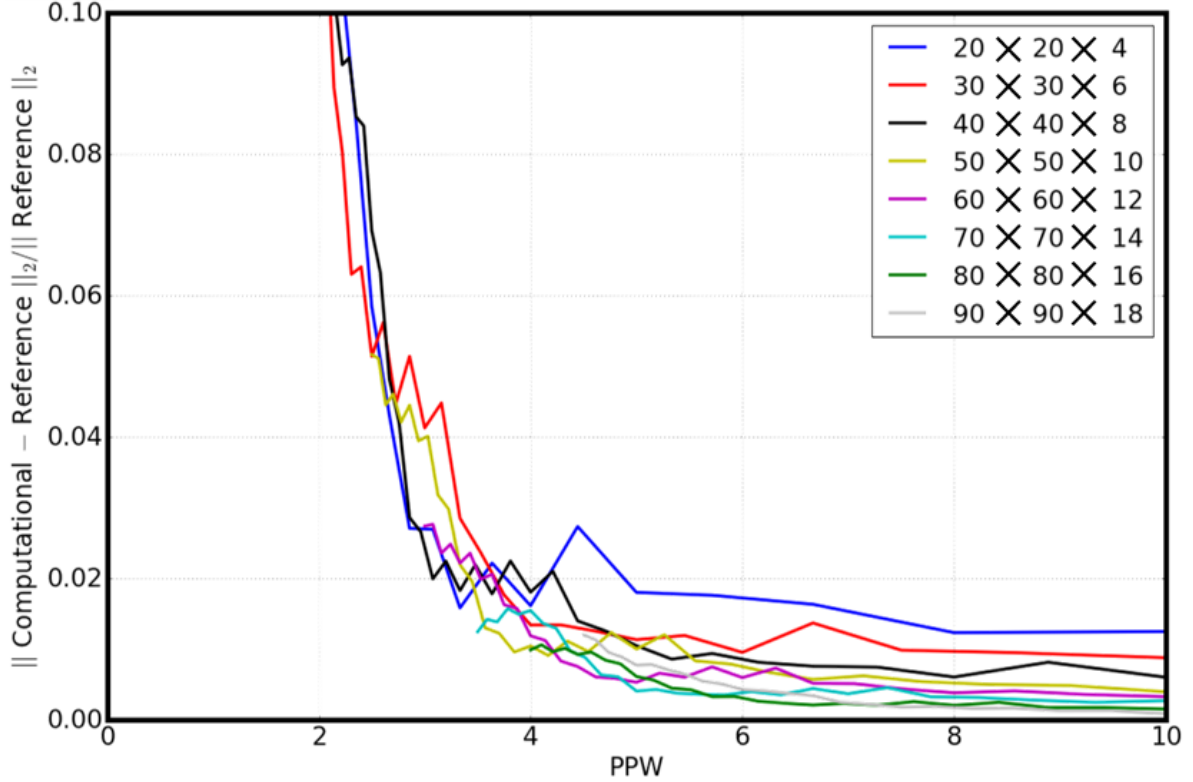


Figure 5: PPW results of the flat plate with rigid body far field scattering solution obtained with constant basis functions.

Moreover, the spatial resolution is measured along the entire surface using the metric of points-per-wavelength-squared (PPW2). PPW2 is computed by:

$$PPW2 = \frac{4\pi^2(p+1)^2[2N_xN_y + 2(N_x + N_y)N_z]}{k^2[2L_xL_y + 2(L_x + L_y)L_z]} \quad (10)$$

where L_y and L_z are the plate lengths along the y - and z - directions, and N_y and N_z are the number of elements along the y - and z - directions (Hu et al. 2016). Equation (10) is equivalent to $(4\pi^2 \times \text{degrees of freedom})$ divided by $(k^2 \times \text{surface area})$. Using the solution computed by $100 \times 100 \times 20$ as the reference, the relative error in the L2 norm is plotted as a function of PPW2. The results of the far field solution with constant basis functions are given in Fig. 6. As indicated in the figure, the relative error measured in the L2 norm of the far field solution becomes as small as 2% when PPW2 is only 25 with constant basis functions being used. This results are similar to the PPW results of the flat plate.

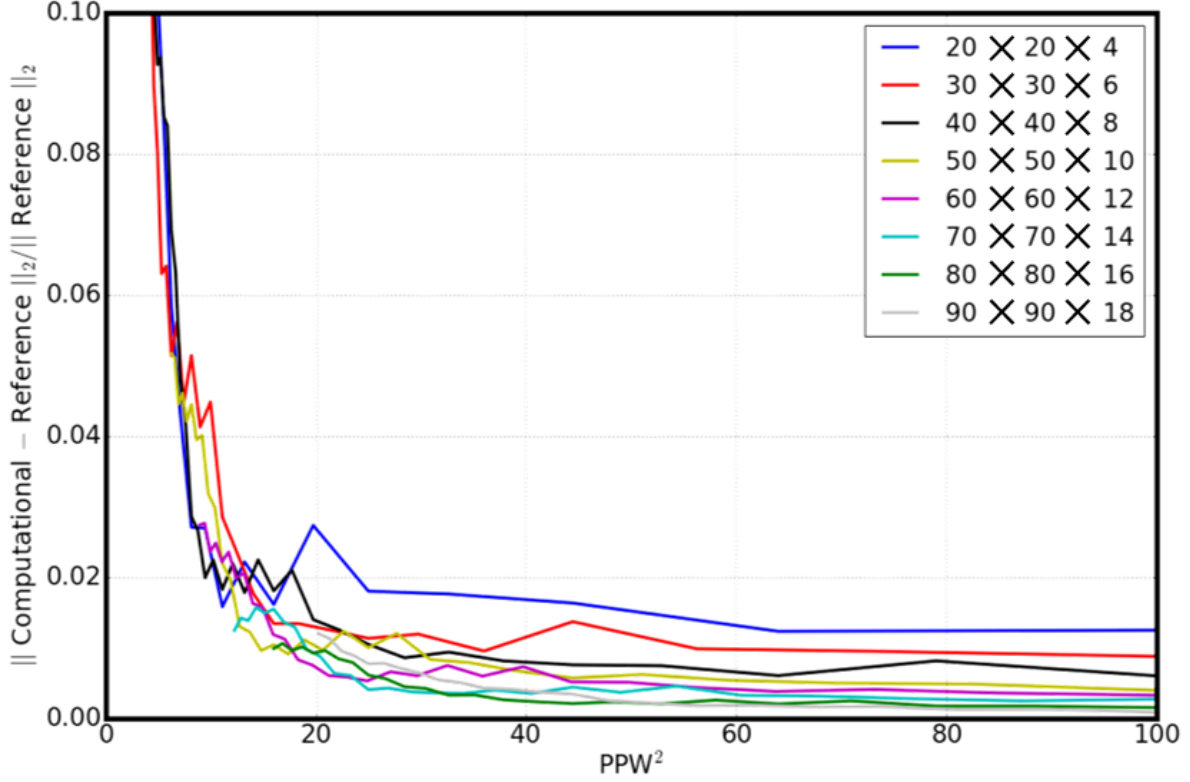


Figure 6: PPW2 results of the flat plate with rigid body far field scattering solution obtained with constant basis functions.

4.2 Sphere with Rigid Body

To further assess the spatial resolution of the TDBEM with respect to the spatial basis functions, we consider the scattering of an acoustic point source by the sphere with rigid body. Both the standard and rotated orientations were considered as illustrated in Figs. 2a and 2b, respectively. A series of computations were carried out by increasing the number of surface elements used from 729 to 32,389. The spatial resolution is measured along the entire surface using the metric of PPW2.

Using the exact solution, the relative error in the L2 norm is plotted for both the standard and rotated orientations as a function of PPW2. The results of the far field solution with constant basis functions are given in Fig. 7 for the standard orientation sphere and Fig. 8 for the rotated orientation sphere. As indicated in the figures, the relative error measured in the L2 norm of the far field solution becomes as small as 4% when PPW2 is only 25 with constant basis functions being used, for both orientated spheres. It appears that, at least for the far field problem, the use of constant elements can keep the overall problem size small while the high-order integration helps to maintain accuracy.

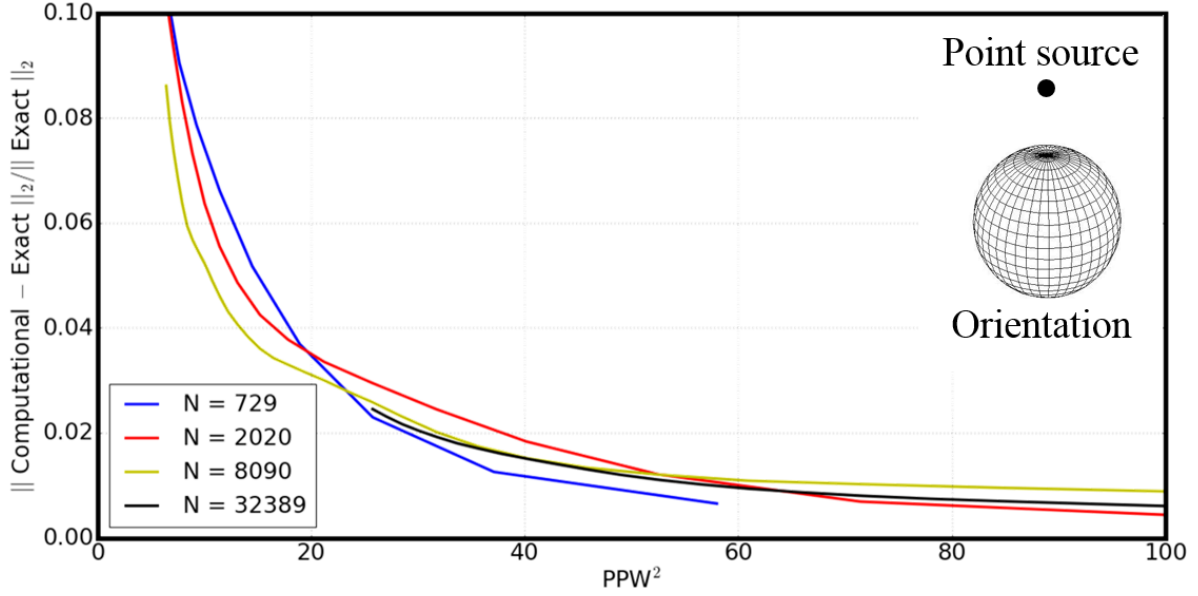


Figure 7: PPW2 results of the standard orientation sphere with rigid body far field scattering solution obtained with constant basis functions.

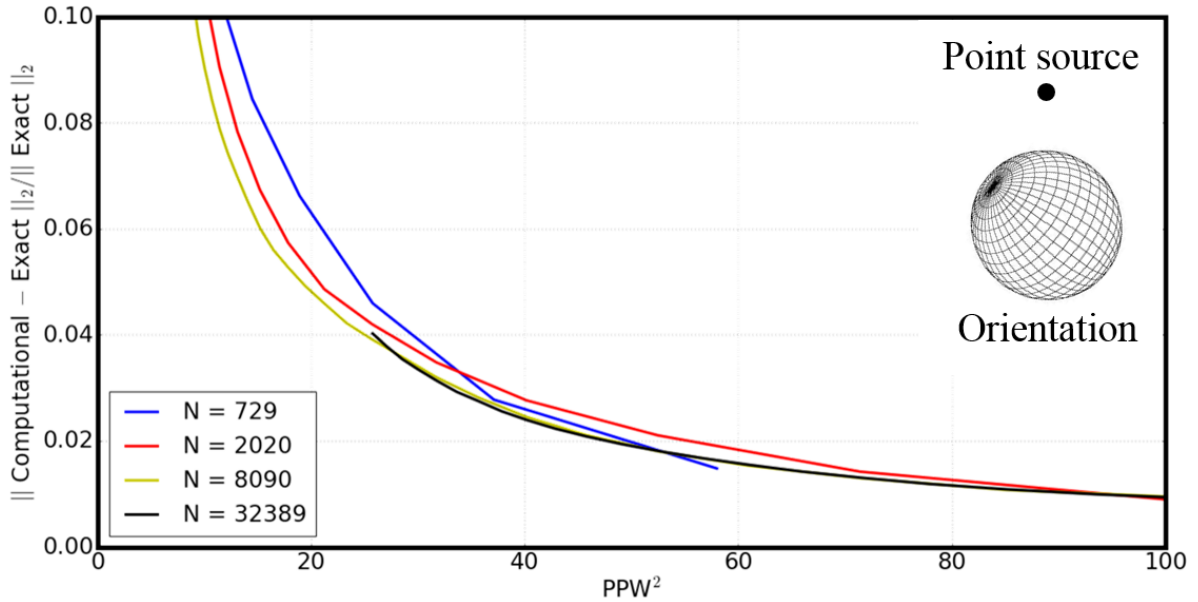


Figure 8: PPW2 results of the rotated orientation sphere with rigid body far field scattering solution obtained with constant basis functions.

4.3 High Performance Scalability Using Central Processing Units

Preliminary work has been completed to assess the scalability of the numerical algorithm by once more considering the scattering of an acoustic point source by the sphere with rigid body. The computations were run using standard compute nodes available through the Old Dominion University Turing Cluster (J. Pratt

2017). Each standard compute node carries 128 GB of memory and contains between 16 and 32 central processing unit (CPU) cores. Using the Turing Cluster, a series of computations using 8,090 and 32,389 surface elements were carried out by doubling the processing power from 8 CPU cores to 16, 32, 64, etc. With each simulation, the average time per iteration was calculated. The expectation was to obtain an inverse relationship between the time per iteration and the number of CPU cores.

The scalability results of using 8,090 surface elements are shown in Fig. 9a and using 32,389 surface elements are shown in Fig. 9b. As indicated in the figures, the numerical results closely match the expected outcome that the time per iteration would be inversely proportional to the number of CPU cores. It should be noted that in both assessments, the numerical results tapered off. This is likely due to an increase in parallel overhead, i.e. an increase in the time associated with performing parallel communications.

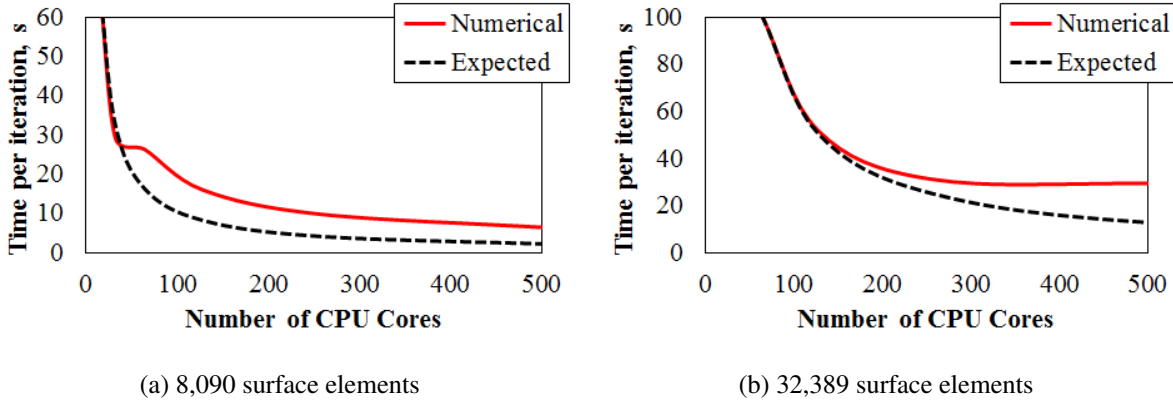


Figure 9: CPU scalability results obtained from modeling the scattering of an acoustic point source by a sphere with rigid body.

4.4 Algebraic Eigenvalue Stability Analysis

Preliminary work has also been completed to assess the numerical stability of the TDBEM by again considering the scattering of an acoustic point source by the sphere with rigid body. In this assessment, we consider a rotated orientation as illustrated in Fig. 2b. With Mach numbers $M = 0.0, 0.3$, and 0.6 , the magnitudes of the largest eigenvalue λ in Eq. (6) were calculated to be, respectively: 1.000000, 0.999994, and 0.999938 using 729 surface elements and 1.000000, 1.000000, and 1.000000 using 2,020 surface elements. These results demonstrate numerical stability since the condition $\max|\lambda| \leq 1$ is satisfied in each case.

5 CONCLUDING REMARKS

When designing next generation quiet aircraft, it is important to be able to accurately and efficiently predict the acoustic scattering by an aircraft body, both rigid as well as lined, from a given noise source. Acoustic scattering problems can be modeled using BEMs by reformulating the linear convective wave equation as a BIE. A TDBIE formulation, as well as its Burton-Miller type reformulation was discussed for solving acoustic wave scattering problems in the time domain using a BEM. The TDBEM formulation is valid for both rigid and lined acoustic surfaces with a mean flow U . Moreover, the incorporation of an impedance boundary condition in the formulation of the TDBEM was introduced. Impedance boundary conditions are assumed when assessing the scattering of an acoustic point source by a geometric body whose surface is treated with an acoustic liner. Further work will be completed to study the implementation of an impedance boundary condition as outlined in Section 6.

Preliminary work has been completed to assess the spatial resolution of the TDBEM with respect to the surface element basis functions by considering the scattering of an acoustic point source by prototype geometric bodies with rigid surfaces. Geometries considered include a flat plate and sphere. The study was limited to the observed far field solution and spatial resolution was measured using the metrics of PPW and PPW2. Results demonstrate that with constant basis functions, the relative errors measured in the L2 norm becomes as small as 2% when PPW2 is only 25 (PPW is only 5) for the flat plate and as small as 4% for the sphere. The scalability of the numerical algorithm was additionally assessed. A series of computations were carried out using standard compute nodes that carry 128 GB of memory and contain between 16 and 32 CPU cores. The scalability results using 8,090 and 32,389 surface elements both demonstrate the expected inverse relationship between the time per iteration and the number of CPU cores. Finally, numerical stability was demonstrated by considering a spherical geometry with rotated orientation. With Mach numbers $M = 0.0$, 0.3 , and 0.6 , the magnitudes of the largest eigenvalue λ in Eq. (6) all satisfied the condition $\max|\lambda| \leq 1$ for both 729 and 2,020 surface elements.

6 FUTURE WORK

Future work includes further studying the:

1. Analytical representation of experimentally determined admittance data and its connection to time domain implementations
2. Coupling of the admittance boundary condition with the TDBIE, including the stable reformulation of the integral equation
3. Analytical stability of the coupled system with respect to the (a) non-existence of resonant frequencies and (b) algebraic eigenvalue stability analysis
4. Simulation of acoustic wave scattering by prototype geometric bodies treated with liners and comparison with analytical and experimental results whenever possible
5. Scalability and computational performance by assessing the algorithmic speedup and efficiency

ACKNOWLEDGMENTS

F. Q. Hu and M. E. Pizzo are supported by a NASA Cooperative Agreement, NNX11AI63A. M. E. Pizzo is also supported in part by an Old Dominion University Modeling and Simulation graduate fellowship. This work used the computational resources at the Old Dominion University ITS Turing cluster and the Extreme Science and Engineering Discovery Environment (XSEDE), which is supported by the National Science Foundation grant number OCI-1053575.

A APPENDIX

Let $p(\mathbf{x}, \omega) = Z(\omega)v(\mathbf{x}, \omega)$ where $p(\mathbf{x}, \omega)$ is the acoustic velocity, $v(\mathbf{x}, \omega) = \mathbf{v} \cdot \mathbf{n}$, \mathbf{v} is the acoustic velocity vector, \mathbf{n} is the inward normal vector on the scattering body, and $Z(\omega)$ is the surface impedance. Then,

$$v(\mathbf{x}, \omega) = Y(\omega)p(\mathbf{x}, \omega) \quad (11)$$

where $Y(\omega) = 1/Z(\omega)$ is the surface admittance (Rienstra 2006). In the frequency domain, v can be represented (Marburg and Schneider 2001) by:

$$v(\mathbf{x}, \omega) = \frac{1}{i\omega\rho_0} \frac{\partial p}{\partial n}(\mathbf{x}, \omega) \quad (12)$$

A.1 Method

By setting Eq. (11) equal to Eq. (12), we obtain a relation for $\partial p / \partial n(\mathbf{x}, \omega)$. Using the inverse Fourier transform convolution property and causality condition which states that $y(t - \tau) = 0$ for all $t - \tau > 0$, i.e. $y(t - \tau) = 0$ for all $t > \tau$, we see that:

$$\begin{aligned}
 Y(\omega)p(\mathbf{x}, \omega) &= \frac{1}{i\omega\rho_0} \frac{\partial p}{\partial n}(\mathbf{x}, \omega) \implies \frac{\partial p}{\partial n}(\mathbf{x}, \omega) = i\omega\rho_0 Y(\omega)p(\mathbf{x}, \omega) \\
 \mathcal{F}^{-1} \left\{ \frac{\partial p}{\partial n}(\mathbf{x}, \omega) \right\} &= \mathcal{F}^{-1} \{ i\omega\rho_0 Y(\omega)p(\mathbf{x}, \omega) \} \\
 P_n(\mathbf{x}, t) = \frac{\partial p}{\partial n}(\mathbf{x}, t) &= \mathcal{F}^{-1} \{ i\omega\rho_0 Y(\omega) \} * \mathcal{F}^{-1} \{ p(\mathbf{x}, \omega) \} = \rho_0 \frac{\partial y}{\partial t}(t) * p(\mathbf{x}, t) \\
 &= \frac{\rho_0}{2\pi} \int_{-\infty}^{\infty} \frac{\partial y}{\partial t}(t - \tau) p(\mathbf{x}, \tau) d\tau = \frac{\rho_0}{2\pi} \int_{-\infty}^t \frac{\partial y}{\partial t}(t - \tau) p(\mathbf{x}, \tau) d\tau \\
 \implies P_n(\mathbf{x}, t) &= \frac{\partial p}{\partial n}(\mathbf{x}, t) = \frac{\rho_0}{2\pi} \int_{-\infty}^t \frac{\partial y}{\partial t}(t - \tau) p(\mathbf{x}, \tau) d\tau
 \end{aligned}$$

A.2 Method

Using the inverse Fourier transform convolution property and causality condition which states that $y(t - \tau) = 0$ for all $t - \tau > 0$, i.e. $y(t - \tau) = 0$ for all $t > \tau$, we see that from Eq. (11):

$$\begin{aligned}
 v(\mathbf{x}, \omega) &= Y(\omega)p(\mathbf{x}, \omega) \\
 \mathcal{F}^{-1} \{ v(\mathbf{x}, \omega) \} &= \mathcal{F}^{-1} \{ Y(\omega)p(\mathbf{x}, \omega) \} = \mathcal{F}^{-1} \{ Y(\omega) \} * \mathcal{F}^{-1} \{ p(\mathbf{x}, \omega) \} \\
 v(\mathbf{x}, t) &= \frac{1}{2\pi} \int_{-\infty}^{\infty} y(t - \tau) p(\mathbf{x}, \tau) d\tau = \frac{1}{2\pi} \int_{-\infty}^t y(t - \tau) p(\mathbf{x}, \tau) d\tau
 \end{aligned}$$

Hence by invoking the differentiation property and commutativity of convolution, we see that from Eq. (12):

$$\begin{aligned}
 v(\mathbf{x}, \omega) &= \frac{1}{i\omega\rho_0} \frac{\partial p}{\partial n}(\mathbf{x}, \omega) \implies \frac{1}{\rho_0} \frac{\partial p}{\partial n}(\mathbf{x}, \omega) = i\omega v(\mathbf{x}, \omega) \\
 \mathcal{F}^{-1} \left\{ \frac{1}{\rho_0} \frac{\partial p}{\partial n}(\mathbf{x}, \omega) \right\} &= \mathcal{F}^{-1} \{ i\omega v(\mathbf{x}, \omega) \} \\
 \frac{1}{\rho_0} P_n(\mathbf{x}, t) = \frac{1}{\rho_0} \frac{\partial p}{\partial n}(\mathbf{x}, t) &= \frac{\partial v}{\partial t}(\mathbf{x}, t) = \frac{\partial}{\partial t} \left[\frac{1}{2\pi} \int_{-\infty}^t y(t - \tau) p(\mathbf{x}, \tau) d\tau \right] \\
 &= \frac{1}{2\pi} \int_{-\infty}^t \frac{\partial}{\partial t} [y(t - \tau) p(\mathbf{x}, \tau)] d\tau = \frac{1}{2\pi} \int_{-\infty}^t \frac{\partial y}{\partial t}(t - \tau) p(\mathbf{x}, \tau) d\tau \\
 \implies P_n(\mathbf{x}, t) &= \frac{\partial p}{\partial n}(\mathbf{x}, t) = \frac{\rho_0}{2\pi} \int_{-\infty}^t \frac{\partial y}{\partial t}(t - \tau) p(\mathbf{x}, \tau) d\tau
 \end{aligned}$$

REFERENCES

- Chappell, D. J., P. J. Harris, D. Henwood, and R. Chakrabarti. 2006. "A Stable Boundary Element Method for Modeling Transient Acoustic Radiation". *Journal of Acoustical Society of America* vol. 120 (1), pp. 74–80.
- Ergin, A. A., B. Shankar, and E. Michielssen. 1999. "Analysis of Transient Wave Scattering from Rigid Bodies Using a Burton-Miller Approach". *Journal of Acoustical Society of America* vol. 106 (5), pp. 2396–2404.
- Hu, F. Q. 2013. "An Efficient Solution of Time Domain Boundary Integral Equations for Acoustic Scattering and Its Acceleration by Graphics Processing Units". *19th AIAA/CEAS Aeroacoustics Conference* (2013-2018).
- Hu, F. Q. 2014. "Further Development of a Time Domain Boundary Integral Equation Method for Aeroacoustic Scattering Computations". *20th AIAA/CEAS Aeroacoustics Conference* (2014-3194).
- Hu, F. Q., M. E. Pizzo, and D. M. Nark. 2016. "On the Assessment of Acoustic Scattering and Shielding by Time Domain Boundary Integral Equation Solutions". *22nd AIAA/CEAS Aeroacoustics Conference* (2016-2779).
- Jones, A. D., and F. Q. Hu. 2007. "A Three-Dimensional Time-Domain Boundary Element Method for the Computation of Exact Green's Functions in Acoustic Analogy". *13th AIAA/CEAS Aeroacoustics Conference* (2007-3479).
- Marburg, S. 2015. "The Burton and Miller Method: Unlocking Another Mystery of Its Coupling Parameter". *Journal of Computational Acoustics* vol. 23 (1550016).
- Marburg, S., and S. Schneider. 2001. "Influence of Element Types on Numeric Error for Acoustic Boundary Elements". *Journal of Computational Acoustics* vol. 11 (3), pp. 363–386.
- J. Pratt 2017. "Turing Community Cluster General Information". <http://www.odu.edu/hpc>. Accessed Apr. 14, 2017.
- Rienstra, S. W. 2006. "Impedance Models in Time Domain Including the Extended Helmholtz Resonator Model". *12th AIAA/CEAS Aeroacoustics Conference* (2006-2686).

AUTHOR BIOGRAPHIES

MICHELLE E. PIZZO is a Ph.D. Candidate in the Department of Mathematics and Statistics at Old Dominion University. She holds a M.S. in Computational and Applied Mathematics from Old Dominion University, M.S. in Mechanical Engineering from Embry-Riddle Aeronautical University, and B.D. in Aerospace Engineering from Embry-Riddle Aeronautical Engineering. Her research interests lie in modeling and simulation, high performance computing, aeronautics, astronautics, and aeroacoustics. Her email address is mpizzo@odu.edu.

FANG Q. HU is a Professor in the Department of Mathematics and Statistics at Old Dominion University. He holds a Ph.D. in Mathematics from Florida State University, M.S. in Engineering from Zhejiang University, and B.S. in Mathematics from Zhejiang University. His research interests include scientific computing, computational aeroacoustics, computational biology and geochemistry, and numerical simulation of fluids and sound. His email address is fhu@odu.edu.