## Multimodal Volume Visualization of Geophysical Data for Archaeological Analysis

Andrew Loomis, Margaret Watters

#### 09/30/2010

#### Abstract

We propose an interdisciplinary research project to investigate the usefulness of multimodal volume visualization techniques to integrate multiple geophysical data sets in an archaeological context. These techniques will be evaluated on their ability to enhance subsurface anomalies, as well as identify and classify archaeological features.

#### **Research Goals**

There are two distinct goals of this research proposal. The first is the use of multimodal volume visualization for the integration of multiple geophysical measures into a single visual representation. The second goal is to evaluate those visualizations based on how accurately they represent known archaeological features in the subsurface.

#### Background

Archaeologists use geophysical surveys as noninvasive methods to measure the physical properties of buried features. Surveys that make use of multiple geophysical measures offer greater insight into the composition and structure of the subsurface (Clay 2001). Recently, there has been significant work in creating integrated representations of these data (e.g. Kvamme 2006, Watters 2006). The work of Kvamme (2001) took data from several different surveys and integrated them in a variety of ways to produce a single image representation. A primary goal of this project is to extend those techniques to create a single volume representation. In many ways this is more natural because several of the survey methods themselves capture data in three dimensions. Previous work in integrating three-dimensional data has been carried out by Watters (2006). That research is focused of visualizing multiple three-dimensional datasets in the same space. While that is an effective integration technique, this project will focus on creating a single representation of the data and not multiple representations in the same space.

From a visualization perspective this problem is one of multimodal volume visualization. Multimodal visualization integrates data collected from multiple sources into a single visualization. The work by Cai (1999) set up a nice framework to explore several different approaches to volume visualization in this area. It describes three different techniques that are primarily distinguished by the level at which the data integration occurs in the rendering pipeline. This project is looking to extend that work by integrating datasets earlier in the rendering pipeline. More recent work in this area is focused on the integration of medical image data (e.g. Tang 2009, Robler 2006). While this work has proven to be effective, the difference in nature between the medical and geophysical data makes it less clear how effective those techniques will be for archaeological analysis.

#### Significance

The results of this project have direct benefit to the archaeological community. Improved techniques for the integrated visualization of multiple geophysical measures in a threedimensional environment can lead to the extraction of new information for archaeological analysis. In addition, the exploration and evaluation of multimodal volume visualization techniques within an archaeological context has benefits and applications beyond this specific field.

#### Methods

The archaeological site that we will be investigating is the Catholme Ceremonial Complex, which is a collection of ritual monuments located in England (Watters 2006). We will be using data that was obtained from two different sources, Ground Penetrating Radar (GPR) and Electrical Imaging. Both of these measures are collected in slices on a regular grid, and extensively processed to in order to remove noise and artifacts. See Watters (2006) for a more detailed treatment of the data processing. Once processed the slices can be assembled into three-dimensional raster volumes. Registration of these volumes to each other is assumed to have occurred in the processing step. Each voxel of those volumes contains a single scalar value representing the measurement taken by each geophysical survey at that point.

There are many different ways to integrate data for multimodal volume visualizations (e.g. Cai, 1999; Tang, 2009; Robler, 2006). The first approach that we will take is to reduce the dimensionality of the problem by integrating the data prior to rendering. This integration will take the form of some simple mathematical functions and builds directly on the work of Kvamme (2006) and Cai (1999). Effectively, it will produce new data in the form of a single volume that represents a combination of the geophysical measures of the site. It is unclear whether this simplification will be able to capture all of the required details for archaeological analysis. If it cannot, additional methods will be employed and expanded to tackle this issue. The works of Robler (2006) and Tang (2009) both describe techniques that may prove effective at this task.

Each of the above techniques will be compared and evaluated against the known archaeological features present at this site. These evaluations will occur with the help experts in the field of geophysical data visualization in archaeology. They will be judged on their ability to discern significant archaeological features, and to aid in the analysis of these features.

#### Work Plan

Week 1: Obtain the data and convert it to manageable/non-proprietary formats.

Week 2-3: Develop software for a basic volume renderer that will provide the framework for multimodal visualizations.

Week 3-4: Implement various techniques for multimodal volume visualization, and make modifications where necessary.

Week 5: Evaluate the effectiveness of these methods.

Week 6: Prepare final presentation and results.

#### References

Cai, W. and Sakas, G. "Data Intermixing and Multi-volume Rendering." *Computer Graphics Forum* 18 (1999): 359-369.

Clay, R. B. "Complementary Geophysical Survey Techniques: why two ways are always better than one." *Southeastern Archaeology*, 2001: 31-43.

Kvamme, K.L. "Integrating Multidimensional Geophysical Data." *Archaeological Prospection* (John Wiley & Sons) 13, no. 1 (2006): 57--72.

Nuzzo, L. and Leucci, G. and Negri, S. and Carrozzo, M.T. and Quarta, T. "Application of 3D Visualization Techniques in the Analysis of GPR Data for Archaeology." *Annals of Geophysics* 45, no. 2 (2002).

Robler, F. and Tejada, E. and Fangmeier, T. and Ertl, T. and Knauff, M. "GPU-based multi-volume rendering for the visualization of functional brain images." *Proceedings of SimVis.* 2006. 305--318.

Tang, H. and Dillenseger, J.L. and Bao, X.D. and Luo, L.M. "A Multi-Volume Visualization Framework for Spatial Aligned Volumes after 3D/3D Image Registration." *Information Science and Engineering (ICISE)*, Dec 2009: 3571--3574.

Watters, M.S. "Geovisualization; an example from the Catholme Ceremonial Complex." *Archaeological Prospection* (John Wiley & Sons) 13, no. 4 (2006): 282--290.

# #

# Representing Choice Reaching Tasks through the Visual Metaphor of Force Attraction

#### Bryan Tyler Parker1, Samuel Birch2 and Joo-Hyun Song3

PI; Department of Computer Science, Brown University, Providence, RI 02912 USA
 <sup>2</sup> Collaborator, Brown University, Providence, RI 02912 USA
 <sup>3</sup> Consultant, Department of Cognitive Science, Brown University, Providence, RI 02912 USA

#### Abstract

We propose an interdisciplinary research project to develop and investigate the effectiveness of methods for the visualization of user reach trajectories, with the goal of revealing hidden cognitive states in the choice reaching tasks and to develop a model for user reaching trajectories.

#### **1. Research Goals**

The first goal is to create a static visualization technique for 3D user hand movement data that effectively visualizes temporal information. The second goal is to create a multi-view technique of visualizing the correlation between 3D user hand movement and 2D eye-tracking data. The third goal is to develop a novel technique for the visualization of the influence of "attractor targets" on user reach trajectories that uses the visual metaphor of force attraction. The final goal is to evaluate the effectiveness of these techniques for developing a model for user reaching trajectories.

#### 2. Background

Perceptual and cognitive processes have largely been inferred based on reaction times and accuracies obtained from discrete responses. However, discrete responses are unlikely to capture dynamic internal processes, occurring in parallel, and unfolding over time. Recent studies measuring continuous hand movements during target choice reaching tasks reveal the temporal evolution of hidden internal events. For instance, the direction of curved reaching trajectories reflects attention, language representations and the spatial number line, in addition to interactions between the ventral and dorsal visual streams. This elucidates the flow of earlier cognitive states into motor outputs. Thus, this line of research provides new opportunities to integrate information across different disciplines such as perception, cognition and action, which have usually been studied in isolation.

Joo-Hyun Song's research attempts to create a model of continuous hand

movements during selection of a target when presented with distracter targets. However, understanding of the data collected has been hindered by ineffective visualization techniques.

#### 3. Methods

Song's experimental trials involve a user selecting a specified target on a screen, with other distracter targets on screen. The users hand motion is captured in 3d, and their eye gaze on the screen is tracked in 2D. Visualization of these various components has been limited to 2D plots.

For the 3D hand movement data, a visualization that allows for recognition of temporal information will be created. This will take the form of a trajectory that has points along it that represent discrete time units. Given a still 3D model of the trajectory, temporal information would be readily evident: points bunched together on the trajectory indicate the user slowing down or stopping while reaching; points spaced far apart indicate fast movement. An option of having multiple user trials of the same task visible at the same time would allow for visual pattern recognition.

This 3D hand movement data will also have to be related to the 2D eye-tracking data. For this, a multi-view solution will be implemented. A 2D window will contain trajectories of the eye-tracking data, with discrete time units much like those used in the 3D visualization. Alongside the eye-tracking trajectories, the 2D projection of the 3D data will also be displayed. This will allow for the analyzer of the data to see how the hand movement follows the eye-movement, as well as related the 3D movement data to the 2D plane of the target screen.

In addition, a novel technique for the visualization of the 'distracter targets'' influence on a task reaching trajectory will be developed. By knowing the 2D location of the different distracter targets, as well as the motion data, one can use an equation similar to the one of gravitational force to enumerate the distracter target's influence on the user. Visualizing this imaginary force can lead to real understanding. A possible visualization could be that each target has a visual circle of influence which scales up or down depending on how strongly it is affecting the user reach trajectory.

These techniques will be compared and evaluated against the current 2D graph visualization already in place. They will be judged by experts in the field of cognitive science for their ability to aid in their understanding of the cognitive states at work during the trials. If a decent model of human reaching trajectories is created with the use of these techniques, where with the current 2D graphs this was proving difficult and they were yet to be able to create such a model, then the new techniques would be considered effective.

#### 4. Impact

The techniques developed through this project will directly benefit the cognitive science community, for it will aid in the model of human target reach trajectories, which will lead to insight on the perception, cognition, and action behind target selection. This in turn would impact the HCI community in developing interfaces. The static visualization of 3D data technique would be applicable to any disciplines that analyze motion data. Most motion can be relatively defined as the summation of attraction/repulsion influences; therefore the novel force visualization technique is also applicable to a wide range of disciplines.

#### 6. Timeline

**Week 1:** Obtain the data and develop software to import it and roughly display it temporally.

Week 2-3: Implement technique to effectively display the 3D data temporarily, as well as view multiple user trials simultaneously.

Week 3: Implement multi-view capability for the viewing of 2D eyetracking data and visualizing the correlation between it and the 3D hand movement data.

**Week 4:** Implement "attractor target" influence visualization.

**Week 5:** Evaluate the effectiveness of these methods.

**Week 6:** Prepare final presentation and results.

#### 7. Facilities

Windows machines in the CIT as well as personal machines will be used to develop and deploy the software. Software will be developed with Cinder. No other facilities will be required.

#### 8. References

#### Statistical Assessment of Individual Peak Quantitation in Mass Spectrometric Data

PI: Justin A. DeBrabant<sup>1</sup> co-PI: Steven Gomez<sup>1</sup> Collaborators: Arthur Salomon<sup>2,3,4</sup>, Kebing Yu<sup>2</sup>

<sup>1</sup> Brown University, Department of Computer Science, Providence RI, USA
 <sup>2</sup> Brown University, Department of Chemistry, Providence RI, USA
 <sup>3</sup> Brown University, Department of Molecular Biology, Cell Biology, and Biochemistry, Providence RI, USA
 <sup>4</sup> Brown University, Center for Genomics and Proteomics, Providence RI, USA

Efficient analysis of large mass-spectrometric data sets has become essential in the field of MS-based proteomics. There are multiple methods for both collecting the data and quantitating it in a way such that the underlying protein and/or peptide structure can be determined with a high level of confidence. Although much work has been done on the quantitation process itself, to our knowledge there is no available metric to compare the relative qualities of different peak quantitations. We propose the development of such a metric, as well as a complimentary visualization tool that will allow quick and efficient quality comparisons among a group of peak quantitations. This work will allow researchers to not only collect the best possible data, but will also allow them to understand the possible error in the data they've collected in way comparable across experiments.

#### 1. SPECIFIC AIMS

There are several goals of this project. First we will develop a metric that will quantify the quality of an individual peak quantitation for mass spectrometric data. This metric will then be implemented in a visualization framework that will allow researchers to easily identify the quality of their data as well as better compare that data across multiple experiments. Within the visualization component will be the ability to filter and rank data based on their individual quality scores as well as simultaneously view data across several experiments.

Once the visualization component is developed, it will be implemented within the framework of the High Throughput Autonomous Proteomic Pipeline (HTAPP). HTAPP, developed by researchers at Brown University, is in an instance of a proteomic pipeline that completely automates the data collection, filtering and analysis as well as the transfer of data between all the distinct components [13]. Specifically, the proposed tool will be adjacent the "Quantitation" and "Relational Database Exploration Tool" components within the framework of the pipeline, with input coming from the quantitation component and output going to the exploration component.

Finally, after implementation within the proteomics pipeline is complete, we will conduct a small user study to evaluate the relative effectiveness of our metric and corresponding visualization component in the comparison of experimental results. As a basis for comparison, we will use the Xcalibur application as it is currently the state-of-the-art in data collection and analysis of MS data [1].

#### 2. BACKGROUND AND RELATED WORK

In recent years, proteomics has emerged as a technique for better understanding the metabolic pathways of the cell, to which proteins are essential. A large part of proteomics research is based around the technique of mass spectrometry (MS), a tool for measuring the molecular mass of molecules in a sample. In proteomics, peptides in a sample are identified by comparing the mass given from the mass spectrometer to known peptide masses, usually from a peptide database such as PeptideDepot [13]. By inducing different environmental conditions and measuring cellular response as a function of protein production, researchers can better understand intracellular functions.

Unfortunately, this process is inherently noisy, and there are several phases in the data acquisition and analysis pipeline that require quantitation of this noisy data, most notably to compare to known peptide masses. The efficient use and analysis of this data has been a continued challenge in the field of MS-based proteomics, as only a fraction of proteins in a sample can be identified, and even a smaller fraction of the identified proteins can be reliably quantified.

Much effort has been invested in the filtering and analysis phases of the pipeline and there are several techniques that can efficiently automate post-acquisition tasks such as peptide quantitation and protein identification [3,4,5,6,7]. Also there has been work in confidence scores of the protein identification phase, i.e. how good of a match is a given spectrogram to a known protein [8, 9, 10]. However, there is no technique to efficiently compare the individual peak quantitations produced by the pipeline. Basically what this means, is that there are methods measure how a given peak quantitation fits to a known protein, but no way to measure the quality of the given peak to begin with. The assumption here is that proteomic researcher has already intuitively categorized peaks as "good" or "bad" and is only using "good" peaks. This manual processing contrasts sharply with the automated nature of all other aspects of the proteomics pipeline, and will become all but impossible as the amount of proteomics data continues to grow.

Also lacking in the field of proteomics is a way to efficiently visualize changes over time and across different runs of the experiment with different conditions. This is the key to proteomics, i.e. given an outside stimulus what is the intracellular response as measured by protein synthesis, so it is surprising that no tool to effectively view these changes is currently available. The present software allows the visualization of an individual peak, with the ability to click through related peaks from other experiments [1]. This method has clear limitations on the ability to understand how the peaks are changing given a change in stimuli, which is the ultimate goal of the experiments.

This work proposes to define a set of variables that will accurately summarize the quality of MS peak data and allow meaningful inter-experiment comparisons. In conjunction with this metric, we will develop a visualization tool that will allow efficient comparison of the individual peaks and their associated quality score across a wide range of experiments. Also available will be the ability to filter results being visualized by the proposed quality score, i.e. only view results above a certain threshold.

#### 3. IMPACT

Already proteomics has emerged as our best method to both understand and subsequently model complex biological systems [11]. With many researcher turning to proteomics as the successor of genomics and the possible future of clinical research, the impact of the above contributions to proteomics are potentially profound. By better modeling and understanding cellular functions, scientists will be able to better understand how and why these functions sometimes break down, e.g. in cancer or autoimmune diseases. However, the difficulty of analyzing the mass of data available is only getting more challenging as the data grows, so better computational methods are needed. To this end, we propose this project.

On its own, the quality metric represents a huge step in the more efficient comparison of peak quantitations, as current methods are limited to labeling peaks "good" or "bad" based on expert user knowledge. With an explosion in the amount of proteomic data being produced in recent years, this manual process will only become less feasible. Using the quality metric to assess the peak quantitations would save the researches time and provide more consistent results by automating the filtering process. It would also give researchers a better understanding of the quality of the results they were viewing, which is important in many contexts.

The proposed quality metric, implemented within the framework of the proposed visualization, will greatly aid researchers in their quest to understand the different cellular responses to varying stimuli. The current software, Xcalibur [1], allows only the viewing of the data one experiment and one time point at a time, from which it is extremely difficult for researchers to effectively deduce cellular changes across varying stimuli. Again, as the amount of data available continues to grow, this will only become more difficult.

Overall, the more formal evaluation of peak quality will reduce the bottleneck in the automated proteomics pipeline currently caused by manual peak quality assessment. Also, a visualization of multiple results together with their quality score will simultaneously allow researchers to better understand their results as well as increase the confidence of their assessment, both of which are essential to the advancement of proteomics research.

#### 4. METHODS

This project can be broken down into 4 distinct phases, each depending on the previous one. A description of each phase follows.

#### 4.1 Phase 1: Data Collection

The first phase of the project will be the reformatting of the MS data from a proprietary binary file format into a relational database. The raw binary file produced by the instrument can only be read by the application, Xcalibur [1], that is provided by the instrument manufacturer. This software provides the quantitation of the data, but without access to the underlying data there is no way to quantify the quality of the quantitation. As such, we propose to create a layer in the pipeline that will take the binary data being produced by the mass spectrometer, translate it into text, and store it in a relational database. This way, the quantitation software can continue to use the binary file, while our quality metric will be able to access the data from a relational database. To make this possible, we have obtained an API from the manufacturer for extracting the binary data into text.

#### 4.2 Phase 2: Metric Implementation

The next phase of the project will be the calculation of the quantitation quality metric. This metric will rely on the characteristics of the underlying data acquisition component. We plan to take into account several error-inducing factors such as signal-to-noise ratio and continuity of the peak shape in order to estimate the selected ion chromatogram quality. We will then measure variance in the experimental results by applying a Gaussian curve to the individual peaks and aggregating variance from the curve.

#### 4.3 Phase 3: Visualization Component

The final phase will be the implementation of the visual component and its integration into the HTAPP. The most important characteristic of the visual component will be the simultaneous visualization of peaks from repetitions and variances from the same experiment. In this context, a repetition is when the same sample is run through the MS multiple times in order to increase the confidence of the results. The other type is when the sample is changed some minute way and run through the MS, looking for changes from the control run. For each experiment, there will be multiple condition runs, each with 2-5 repetitions. For each run, the individual peaks will have an associated quality score. The visualization will give the ability to view repetitious runs and their associated quality score in a 3 dimensional grid where each repetition occupies a unique plane in x-y grid parallel to the planes for the other repetitions. The quality score will be indicated by a color scheme on the graphs. A similar view will be available for each of the different conditions, with each condition graph representing an aggregate score from each of its repetitions and colored to represent the aggregate of the quality metrics of these aggregations. The key to this visualization will be the side by side comparison across both repetitions and experiments with varying conditions, thus visualizing, with confidence, intracellular response to a given stimuli.

#### 4.4 Phase 4: Evaluation

The final phase will be the evaluation of the quality metric and the visualization component of the project using a small user study. To evaluate the quality metric we will have expert users from the Salomon Lab at Brown University manually label peaks as "good" or "bad". We will then decide on a threshold that will separate the "good" and "bad" peaks and apply our quality metric and label and the peaks below the threshold "bad" and all those above as "good". We will then compare the intersection of the two sets to determine how accurately our metric models expert user intuition. For the evaluation of the visualization component, we will use the Xcalibur [1] system as a basis for comparison, as it is widely used to compare and analyze MS data.

#### 5. TIMELINE

Week 1: Transfer of raw binary data files to a relational database

- **Week 2:** Study of the possible underlying parameters that would best describe the quality of the quantitation
- Week 3: Quality metric for peak quantitations
- Week 4: Begin phase 3, implementation of the visual component
- Week 5: Finish visual component and insert into HTAPP
- Week 6: User study of quality metric and visualization component

#### 6. **REFERENCES**

#### [1] Xcalibur. http://thermoscientific.com

[2] M. Bantscheff, et al. (2007). 'Quantitative mass spectrometry in proteomics: a critical review'. Analytical and Bioanalytical Chemistry 389(4):1017-1031-1031.

[3] J. Cox & M. Mann (2008). 'MaxQuant enables high peptide identification rates, individualized p.p.b.range mass accuracies and proteome-wide protein quantification'. Nature biotechnology 26(12): 1367-1372.

[4] Craig, R., and Beavis, R. C. (2003). 'A method for reducing the time required to match protein sequences with tandem mass spectra'. Rapid Commun. Mass Spectrum. 17(20):2310-6

[5] Craig, R., and Beavis, R. C. (2004). 'TANDEM: matching proteins with tandem mass spectra'. Bioinformatics. 2004 Jun 12;20(9):1466-7. Epub 2004 Feb 19.

[6] J. K. Eng, et al. (1994). 'An approach to correlate tandem mass spectral data of peptides with amino acid sequences in a protein database'. Journal of the American Society for Mass Spectrometry 5(11): 976-989.

[7] J. K. Eng, Bernd Fischer, Jonas Grossmann, Michael J. MacCoss (2008). 'A Fast SEQUEST Cross Correlation Algorithm'. Journal of Protemics 7 (10), 4598-4602

[8] David Fenyö and, Ronald C. Beavis (2003). 'A Method for Assessing the Statistical Significance of Mass Spectrometry-Based Protein Identifications Using General Scoring Schemes'. Analytical Chemistry 2003 75 (4), 768-774

[9] Martinez Lundgren and Han Wright (2009). 'Protein identification using Sorcerer 2 and SEQUEST'. Curr Protocol Bioinformatics.13(13.3).

[10] L. C. McHugh & J. W. Arthur (2010). 'Harvest: an open-source tool for the validation and improvement of peptide identification metrics and fragmentation exploration'. BMC bioinformatics 11(1): 448+.

[11] S. Patterson and R. Aebersold (2003). 'Proteomics: the first decade and beyond'. Nature Genetics (33):311-323.

[12] K. Yu, et al. (2009). 'Integrated platform for manual and high-throughput statistical validation of tandem mass spectra'. Proteomics 9(11):3115-3125.

[13] K Yu, et al. (2009). 'PeptideDepot: flexible relational database for visual analysis of quantitative proteomic data and integration of existing protein information'. Proteomics 9(23):5350-8.

## CREATING MODEL OF HUMAN COGNITION IN GRAPHICAL QUANTITATIVE DATA EMBELLISHMENTS

PI: Diem Tran Consultant: Caroline Ziemkiewicz

Department of Computer Science, Brown University, Fall 2010

#### Abstract

Graphical data embellishment is reported as one of the top unsolved problem in visualization. Though several studies have been conducted to understand the problem, its solution is still in dispute due to opposite findings from the studies. We create a new model to predict human cognition in reading quantitative data displayed in graphs, in order to explain various outcomes of previous studies as well as establishing a unified approach in drawing graphs and designing visualization techniques, benefitting a large portion of the visualization community. A timeline of 6 week research includes a preliminary study to collect data for model design, a two-week period of developing the model and two weeks for validation and evaluation.

## Contents

1	INTRODUCTION	3
2	RELATED WORK	3
3	CONTRIBUTION & SIGNIFICANCE	4
	3.1 Contribution	4
	3.2 Significance	4
4	PROCESS	4
5	CONCLUSION	5

## **1 INTRODUCTION**

According to Chen [1], the problem of enhancing graphs is among top unsolved problems in the visualization field. This problem involves directly the understanding of how users perceive data from different kinds of visualization techniques, hence affects design of information visualization systems. It is further related to evaluation of effects that visualization tools have on users. Overall, whether or not and how to make graphs become appealing and easier to retrieve for people is an open problem to the field of visualization.

In this proposal, we will use cognitive architecture widely used by scientists to model human cognition on graphical data embellishment to understand more about how people read, interpret and memorize data, as the area has controversial results from different studies. Our model will be compared with real results obtained from user experiment to evaluate the validity of predictions we made in the model.

### 2 RELATED WORK

Several works have been done to model cognitive tasks in general. ACT-R [2] is a cognitive architecture to explore how human cognition works. It is used to physically model the process that occurs inside human's mind. Scientists can use it to replicate cognition when users perform different tasks, and test the validity of replication by comparing its results with results obtain from real users. GOMS [3] is also a similar architecture to ACT-R in terms of modeling human cognition through decomposing the cognitive process into procedural steps which can be imitated and performed by computers. Effectiveness of the two tools has been proved to be valid, and is used widely among computer scientists to understand what is happening in human's mind when interacting with computers [3, 4]. For the graph perception area, Lohse implemented a GOMS-based model to understand graphical retrieval of human in general [5]. However, no complex model has been built to understand human cognition in perceiving different kinds of enhancement in presenting data graphically.

On the other hand, there are empirical evidences of how graph embellishments affects capability to memorize data. Nevertheless, results are very controversial due to opposite results in different studies. Some results yield favor for enhancement of displaying data [6, 7] in various ways, from using picture [8] to applying emotional tones of the data [9] to make users remember them more accurately. Meanwhile, others follow a long-established trend of keeping graphs as plain and simple as possible [10, 11]. There are also neutral results

in which non of the two, enhanced of minimalist approach in graphing is effective [12, 13]. All of these outcomes put the visualization community in confusion where no direction is a better choice, since there is no general approach to identify the problem, as the key to it lies behind human cognition in perceiving data. To unify and explain all results, we need a general model to describe how a user retrieve data in a procedural way.

## **3 CONTRIBUTION & SIGNIFICANCE**

#### 3.1 Contribution

Understanding the unsolved problem of visualization, we create a new model using ACT-R and/or GOMS with the aim of predict what a user typically does when reading graphs. By modeling this cognitive process, we are able to have more insight in how graph embellishments affect users' memory. In particular, we can understand how specific portion of data display in a graph is retrieved into a reader's mind, and how the reader interpret, analyze and synthesize the data. Following this cognitive procedure, we are able to explain why there are controversial results in previous studies, and provide guidelines for future graph drawings as well as design of visualization systems.

#### 3.2 Significance

If successful, the model will contribute significant resources for the top unsolved problem of visualization. It will helps graph drawers in designing more appealing and effective graphs, transmitting data more persuasively in ways that drawers want readers to perceive. Furthermore, scientists interested in how displayed data affect users' perception can based on our model to design more effective techniques to display data. It also open paths for future development of models to predict human cognition in other complex data display techniques.

## 4 PROCESS

Week 1-2: Preliminary user study. We conduct a user study to grasp initial understanding of how a person read different kinds of graphs. In the study, questionnaires and interviews are used to collect data.

- Week 3-4: Developing the model. Based on results obtained from the user study, we will create the model to predict user's cognition process. We aim to use both GOMS and ACT-R. Small experiments may be conducted to verify and refine the model.
- Week 5: Model evaluation. We run an experiment to compare results from our model and real data collected from users. Comparisons are done between real data and either ACT-R based or GOMS based model.
- Week 6: Analysis of results and report writing. Using statistical testing units, we validate our results and write final report.

## 5 CONCLUSION

In this proposal, we address one of the top unsolved problem in visualization: embellishments on data displayed in graphs. At first, we discuss the current state of the area by citing different empirical studies which yield controversial results. Next, we point out the need of creating a new model to unify and explain those results. We also conduct a timeline to show what and how we are going to invest into the model. The success of our research will establish many advantages of the field visualization.

## References

- [1] C. Chen, "Top 10 unsolved information visualization problems," *IEEE Comput. Graph. Appl.*, vol. 25, no. 4, pp. 12–16, 2005.
- [2] J. R. Anderson, M. Matessa, and C. Lebiere, "Act-r: a theory of higher level cognition and its relation to visual attention," *Hum.-Comput. Interact.*, vol. 12, no. 4, pp. 439– 462, 1997.
- [3] B. E. John and D. E. Kieras, "The goms family of analysis techniques: Tools for design and evaluation," tech. rep., Pittsburgh, PA, USA, 1994.
- [4] J. R. Anderson, *How Can the Human Mind Occur in the Physical Universe?* Oxford University Press, USA, 2007.
- [5] G. L. Lohse, "A cognitive model for understanding graphical perception," *Hum.-Comput. Interact.*, vol. 8, no. 4, pp. 353–388, 1993.

- [6] S. Bateman, R. L. Mandryk, C. Gutwin, A. Genest, D. McDine, and C. Brooks, "Useful junk?: the effects of visual embellishment on comprehension and memorability of charts," in CHI '10: Proceedings of the 28th international conference on Human factors in computing systems, (New York, NY, USA), pp. 2573–2582, ACM, 2010.
- [7] N. Holmes, "Designer's guide to creating charts and diagrams," 1984.
- [8] W. Hockley, "The picture superiority effect in associative recognition," *Memory and Cognition 36*, pp. 1351–1359, 2009.
- [9] M. Mather and K. Nesmith, "Arousal-enhanced location memory for pictures," *Journal of Memory and Language 58*, pp. 449–464, 2008.
- [10] E. R. Tufte, *The visual display of quantitative information*. Cheshire, CT, USA: Graphics Press, 1986.
- [11] W. S. Cleveland, *The elements of graphing data*. Belmont, CA, USA: Wadsworth Publ. Co., 1985.
- [12] J. Kulla-Mader, "Graphs via ink: Understanding how the amount of non-data-ink in a graph affects perception and learning," 2007.
- [13] A. Blasio and A. Bisantz, "A comparison of the effects of data-ink ratio on performance with dynamic displays in a monitoring task," *International Journal of Industrial Ergonomics 30*, pp. 89–101, 2002.

## Do explanations promote confidence in uncertainty visualizations? A user-study in medical diagnosis

**Gideon Goldin** 

PI; Cognitive, Linguistic & Psychological Sciences, Brown University

#### **Steve Gomez**

Collaborator, Computer Science, Brown University

#### **Elizabeth Bird**

Consultant, Rhode Island Hospital

#### **Steven Sloman**

Consultant, Cognitive, Linguistic & Psychological Sciences, Brown University

#### Abstract

How to best represent uncertainty is an unsolved problem in visualization research. We propose a novel technique in which explanations for uncertainty, rather than just uncertainty itself, are visualized. We predict that the inclusion of explanatory information will lead to more confidence in decision-making under risk. A user-study is proposed.

#### Introduction

Risk in decision-making constitutes a type of uncertainty associated with the selection of a particular choice. The risk perception literature in the medical domain demonstrates that people behave as less than ideal Bayesian reasoners, succumbing to such fallacies as base-rate neglect (Hoffrage & Gigerenzer, 1998), especially when careful deliberation is inhibited (e.g., time-pressure; Maule & Edland, 1997). Medical doctors operate within environments where time-pressure, stress, and other inhibitory factors often prevent fully rational decision-making.

Despite progress in artificial intelligence on expert systems, Bayesian networks, and even human factors, physicians are still slow to adopt decision-support systems (Egea & Gonzalez, 2010). There are myriad hypotheses, some building off models of technology acceptance (e.g., T.A.M. (Davis, Bagozzi & Warshaw, 1989)) designed to try to explain why this may be. Reasons include but are far from limited to trust in information sources, and innumeracy (Lipkus & Peters, 2009).

We propose that a potential barrier for more widespread use of these types of software systems lies in a lack of desire to use these systems due to the format of the presented information. In particular, we hypothesize that the lack of causal and/or mechanistic explanations, even obvious ones, play a particularly detrimental role in diagnosis.

We are not necessarily attempting to create or discover methods which will foster better probabilistic reasoning in medical judgment and decision-making, though this is not an unlikely product of this research. Rather, we are trying to verify whether or not the provision of different types of explanation (including verbal and graphical ones) will foster *more confidence* in a given probabilistic assessment.

Thus, this project will aid decision-making under uncertainty by offering visualized explanations in-tandem with that uncertainty. These explanation visualizations may take different forms, ranging from a diagrammatic representation of a causal model in a graph format, to the simple inclusion of a verbal explanation (sometimes of something that is even already known). It is explanations, whether simple or complex, that contextualize data that may otherwise seem arbitrary or even myseterious. Indeed, Lombrozo T. (2006) says that explanations accommodate novel information in the context of prior beliefs. Without such explanations, we are not certain of what an expert system, for example, is computing. However, given explanations, we are drawn, automatically, to a conclusion, albeit a probabilistic one, that is based in reason. Whether or not this is a good thing, naturally, depends on its application, as simply including an explanation may promote trust in a faulty computation just as easily as in a safe one.

#### **Significant Contributions**

The contributions of this project are multi-fold. First, the unsolved problem of visualizing uncertainty will be addressed. Empirical data from the psychological literature will be combined with our user-studies to offer new insights and evidence on how to best integrate explanations into the representation of uncertainty. Second, contributions will be made across fields, as potential benefactors reside in computer science (e.g., human factors, human-computer interaction, visualization, visual analytics), psychology (e.g., probability judgment, decision-making), and even education (e.g., evidence-based approaches to medical school, patient/layperson self-diagnosis).

#### **Broader Impact**

The primary impact of this project is clear: improving medical diagnosis by way of increasing the adoption rate of decision-support systems. The project outlined in this proposal is comprised, at its most basic instantiation, as a simple visualization technique. This lends itself well to being implemented in examination rooms across the world, as it is inexpensive and would be only a matter of coupling graphical (or verbal) elements to otherwise textual and numeric data.

Furthermore, looking beyond our 6-week constraint, this work has the potential to be integrated into, or to begin defining a framework for visualizing explanations and uncertainty. The ultimate goal of a framework of this nature could be the automatic generation of visualizations. For example, one could imagine that given some input data, a system could automatically determine an explanation for some data (doing probabilistic inference), and could likewise present this information in the most appropriate manner (e.g., causal graphs for mechanistic explanations, time-lines for temporal explanations, examples cases for analogical explanations, etc.).

#### **Research Goals**

The research goals for this project include unifying literature in computer science and psychology, designing and running a user-study over the web to address our hypothesis, analyzing the data, and disseminating our results among the public.

#### **Related Work**

Related work in this field is far from unified. Relevant literature already exists in cognitive psychology (e.g., decision-making, probability judgment, knowledge representation, etc.), computer science (e.g., information visualization, artificial intelligence, human-computer interaction, etc.), medicine (e.g., medical education, medical practice, etc.). However, these literature bases are disparate and would benefit from generalization in a holistic account.

In visualization research, there is work addressing the uncertainty problem (for example, Riveiro (2007)). However, many of the concrete suggestions reside in 3D applications. In psychology, much research has been done to address information (especially probability) formats. Gigerenzer and Edwards (2003) review some of these findings, showing, for example, that often times

people reason better with natural frequencies (10 out of 100 people) than probabilities (10% of people). Others have demonstrated differences when using different graph types (Elting, Martin, Cantor & Rubenstein, 1999). Furthermore, presenting numeric vs. verbal descriptors can affect risk perception, as shown by Young & Oppenheimer (2009), among others. Our general motivation force is the notion that explanations offer a framework with which to understand data. If a computer can generate explanations, perhaps people will infer that it can understand its own data, and thus, make intelligent predictions.

#### Methods

In order to test our hypothesis that explanatory information promotes use of decision-support systems, we will run a user-study. This study will consist of multiple hypothetical patient cases in which the subject must indicate to some degree his or her desire to use the given system/information.

Consider one example, examining the influence of vacuous explanations, below:

Imagine you are presented with a case in which a particular patient is presenting with heavy sweating. You are wondering if this patient has the influenza virus. You have been given a computer which can purportedly make predictions about diagnoses. It suggests that:

(1) The probability of the influenza virus given heavy sweating is: 30%.

(2) The probability of the influenza virus given heavy sweating is: 30% because the influenza virus can cause a fever which raises body temperature, causing a person to sweat excessively.

(3) The probability of the influenza virus given heavy sweating is: 30% because the influenza virus can cause a fever which raises body temperature, causing a person to sweat excessively (a diagram is included):



"How confident are you in endorsing this assessment?"

In cases (2) and (3), we are not (arguably) given any more information than in (1) because supposedly, we already know about this explanation (thus the explanatory information is

vaccuous). However, it is hypothesized that when the computer seems to also "understand" this explanation, it seems much more credible.

A simple three-factor, between-subject condition (*no explanation vs. verbal explanation vs. verbal & graphical explanation*) will be compared for this task, and similar ones. We predict that the *explanation* conditions will result in significantly higher confidence levels than the *no explanation* condition, although no predictions are made to distinguish the two *explanation* conditions. All tasks in the study will be run on the web, using Mechanical Turk to recruit and pay participants.

In general, tasks will entail a comparison of one or more *explanation* conditions against a *no explanation* condition. For example, another task might compare simply providing the user with two relevant datum versus providing the user with two relevant datum linked by an explanation.

#### Work Plan

*Week 1 - Week 2*: During the first two weeks, the literature will be searched for all relevant articles. Two meetings will transpire between project participants to discuss study direction. *Week 3*: During this week, a study will be designed, agreed-upon, and run online with normal participants. This will provide us with quick feedback on a pilot experiment. This study will be run using Mechanical Turk.

*Week 4*: After the pilot, a modified paper-version of the study will be run in-person with actual physicians. This will be run in the Rhode Island Hospital under the supervision of Dr. Elizabeth Bird.

Week 5: Analysis of the study will be completed.

*Week 6*: The project will be written up.

#### **Progress Report Plan**

A private web site will be made available for disseminating project progress.

*Week 1 - Week 2*: In order to demonstrate progress of the literature search, an annotated bibliography will growing on the private web site.

*Week 3*: Details and results from the pilot experiment will be posted on the private web site for evaluation.

*Week 4*: Details and results from the final experiment will be posted on the private web site for evaluation.

*Week 5*: Analysis of the study will be completed, with results posted on the private web site. *Week 6*: A summary of the project and a discussion of results will be posted to the private web site. This will take the form of a final presentation (slideshow) and final report (paper).

#### References

Hoffrage, U., & Gigerenzer, G. (1998). Using natural frequencies to improve diagnostic inferences. *Academic Medicine*, 73(5), 538-540.

Maule, J. A., & Edland, A. C. (1997). The effects of time pressure on human judgment and decision making. In R. Ranyard, R. W. Crozier, & O. Svenson (Eds.), *Decision making: Cognitive models and explanations*. London and New York: Routledge.

Lipkus, I., & Peters, E. (2009). Understanding the role of numeracy in health: Proposed theoretical framework and practical insights. *Health Education & Behavior*, *36*(6), 1065-1081.

Davis, F.D., Bagozzi, R.P., & Warshaw, P.R. (1989). User acceptance of computer-technology: A comparison of two theoretical-models. *Management Science* 35(8), 982–1003.

Egea, J.M.O., & Gonzalez, M.V.R., (2010). Explaining physicians' acceptance of EHCR systems: An extension of TAM with trust and risk factors. *Computers in Human Behavior*, In Press.

Gigerenzer G., & Edwards A. (2003). Simple tools for understanding risks: from innumeracy to insight. *BMJ (Clinical research ed.)*. 327(7417), 741-744.

Elting L., Martin C., Cantor S., & Rubenstein E. (1999). Influence of data display formats on physician investigators' decisions to stop clinical trials: prospective trial with repeated measures. *British Medical Journal*. *318*(7197), 1527.

Young S., & Oppenheimer D.M. (2009). Effect of communication strategy on personal risk perception and treatment adherence intentions. *Psychology, health & medicine*. *14*(4), 430-42.

Lombrozo T. (2006). The structure and function of explanations. *Trends in cognitive sciences*. *10*(10), 464-470.

Riveiro M., (2007). Evaluation of uncertainty visualization techniques for information fusion. *10th International Conference on Information Fusion*. 1-8.

Curriculum Vitae for Gideon Goldin

#### **GIDEON GOLDIN**

Box 1978, Cognitive, Linguistic & Psychological Sciences, Brown University, Providence, RI 02906 (321) 720-5377 Gideon Goldin@Brown.edu

#### **EDUCATION**

**Doctor of Philosophy in Cognitive Science,** Fall 2008 – Present Department of Cognitive & Linguistic Sciences, Brown University Providence, Rhode Island

#### Bachelor of Science in Computer Software Engineering, Spring 2008 College of Computer Information Science and Engineering, University of Florida, Gainesville, Florida Minor: Linguistics, Spring 2008

#### **RELEVANT COURSEWORK**

#### **Independent Study in Human-Computer Interaction and Visualization** Professor Katherine Spoehr Department of Cognitive & Linguistic Sciences, Brown University

#### **Cognition, Human-Computer Interaction, and Visual Analysis**

Professor David Laidlaw Department of Computer Science, Brown University

#### **Human-Computer Interaction**

Professor Benjamin Lok College of Computer Information Science and Engineering, University of Florida

#### **Aesthetic Computing**

Professor Paul Fishwick College of Computer Information Science and Engineering, University of Florida

#### **RESEARCH EXPERIENCE**

**Research Assistant**, Fall 2010-present Sloman laboratory, Brown University

• Automated system for conducting online experiments.

#### Graduate Researcher, Fall 2008-present

Sloman Laboratory, Brown University

OCD Research, Butler Hospital

• Designed, conducted, and analyzed an experiment assessing the probability judgments of risk of healthy and obsessive-compulsive people.

#### Undergraduate Researcher, June 2007 - August 2007

Linguistics Laboratory, University of Florida

**Psycholinguistics Group** 

• Migrated an ERP study for usage in an eye-tracking environment.

#### Undergraduate Researcher, September 2006 - December 2006

High-Performance Computing & Simulation Research Laboratory, University of Florida Unified Parallel C Group

• Programmed small-scale simulation applications in Unified Parallel C.

#### SKILLS

**Computer Languages**: Java, Perl, LISP, C, XHTML, JavaScript, JQuery, CSS, ActionScript, MATLAB, PHP **Software Packages**: MATLAB, Net beans, Eclipse, SPSS, Adobe CS

#### ACHEIVEMENTS AND HONORS

Research Award, Brown Institute for Brain Sciences, Brown University (2009-2010) University Fellowship, Brown University (2008-2009) Dean's List: Fall 2004, Spring 2006, Spring 2007 Florida Academic Scholar

#### Curriculum Vitae for Steve Gomez

CONTACT INFORMATION Steven R. Gomez Department of Computer Science Brown University Box 1910 Providence, RI 02912 USA

office: CIT 423 cell: (802) 272-4941 e-mail: steveg@cs.brown.edu www.cs.brown.edu/ steveg

#### **RESEARCH INTERESTS**

Scientific visualization, visual computing and analytics, human-computer interaction, computer graphics, vision

#### **EDUCATION**

Brown University, Providence, Rhode Island USA

Second-year graduate student in doctoral program in Computer Science

• Courses completed: Computer Graphics, Computer Vision, Computational Photography, Distributed and Parallel Computing, Networking and Distributed Systems

#### Dartmouth College, Hanover, New Hampshire USA

B.A. magna cum laude, Computer Science, June 2007

#### INDUSTRY

M2S, Inc., Lebanon, NH - Software Engineer - 2008-09

Development for several tools, including virtual stent grafts for Preview aortic aneurysm visualization software.

#### AWARDS

Brown University Graduate Fellowship – 2009–10 Rufus Choate Scholar (top 5% of class) – 2006–07 John G. Kemeny Undergraduate Computing Prize – First Place, 2006, for BlitzChat Dartmouth Presidential Research Scholar – 2005–06 Class of 1928 Endowed Scholarship – 2003–07

#### PUBLICATIONS

Steven R. Gomez, Radu Jianu, and David H. Laidlaw. A Fiducial-Based Tangible User Interface for White Matter Tractography. In Proceedings of ISVC, 2010.

- Steven R. Gomez. Interacting with Live Preview Frames: In-Picture Cues for a Digital Camera Interface. In Proceedings of ACM UIST (Poster), 2010.
- Keller, R., Hunt, M., Jones, S., Morrison, D., Wolin, A., and Gomez, S. 2007. Blues for Gary: Design Abstractions for a Jazz Improvisation Assistant. Electron. Notes Theor. Comput. Sci. 193 (Nov. 2007), 47-60.

#### TECHNICAL SKILLS

Languages: Java, MATLAB, C, C++, Scheme, Perl, Tcl/tk, PHP, SQL, X/HTML, LATEX Applications: Eclipse, Netbeans, Xcode, Apache, CVS, Subversion, MS Office, Adobe CS Operating Systems: Mac OS X, Linux, Windows 98/2000/XP Letter of Support from Elizabeth Bird Support has been attained. A more formal email is coming soon. Letter of Support from Steven Sloman Support has been attained. A more formal email is coming soon.

## PRELIMINARY: Visualization of Hyperspectral Images Through Interactive Non-linear Sectionals

Ryan P. Cabeen

Department of Computer Science Brown University, Providence, RI, USA rpc@cs.brown.edu

September 30, 2010

#### Abstract

The proposed study presents a method for visualization of hyperspectral images where a user is able to interactively choose non-linear sections through an image, reducing the spatial dimension and resulting in several types of visualization of data along the chosen curve in spatial coordinates.

#### 1 Overview

In the proposed study, we intended to develop and evaluate a visualization method for understanding hyperspectral image data through interactive spatial dimension reduction. The research will result in an implementation of the method for use with remote sensing data. The efficacy of the method will be evaluated with expert users, and the integration of the tool with existing remote sensing data analysis environments will be explored.

In the field of remote sensing, there has been both an increase in the amount of images being taken, their spatial and spectral resolution, producing a need for new visualization techniques. Unlike traditional photography, where there are sensors for bands in red, green and blue bands, remote sensing systems sample a broad spectrum that includes the infrared at anywhere from a few to hundreds of bands. This type of imaging can give insight into the structure and mineral composition of the field of view; however, intepreting the shape of the spectra is a highly specialized skill that requires knowledge both of the physics of the imaging process and the subject's composition. Furthermore, the broad spectrum being recorded at each pixel is difficult to render as an image that can be interpreted by the human visual system's three bands of sensitivity. The combination of increases in data volume and resolution has produced a need for not a single visualization method but an array, to which the proposed study will contribute.

There are a variety of methods in use for the hyperspectral image visualization, and these have originated both as specific to the field and as an applications of more general visualization techniques. Typically, the visualization will involve dimension reduction, which can be affected in the spectral and spatial domains. Examples of spectral reduction can including taking band ratios, principle component analysis and many others. Examples of spatial reductions are plotting the spectrum at a single point, rendering an image cube and computing statistics of the spectra in a neighborhood.

The proposed study investigates the use of a spatial reduction method that uses interaction to benefit from the expertise of the user. The use will specify a simple curve in spatial coordinates. At each point along this curve, the spectra are sampled. The result can then be rendered as an image or a three dimensional elevation plot. Unlike spectral reduction techniques, the spectra along the curve are presented in full, which might be advantageous when the user is accustomed to interpreting the data based on fine details. Through the course of this study, the implementation, advantages and limitations of this method will be investigated from both the perspective of a tool maker and an expert user.



Figure 1: Point sample of a spectrum in a hyperspectral image

#### 2 Aims

There are several goals of the proposed study including the development of the visualization method, application to remote sensing data, evaluation by expert users and integration into existing research environments.

The proposed visualization method fills a niche in the array of visualization techniques by providing a spatial data reduction technique that is not available in the traditional remote sensing image analysis environment. Purely spatial data reduction techniques might be advantageous because they allow the user to understand specific characteristics of the spectra. Existing methods for this type of reduction include point sampling, image cube rendering and neighborhood statistics. Point sampling measures the spectrum at a single pixel, plotting the data in a typical graph, which is illustrated in Fig. 1. Image cube rendering samples the spectra along line segments that are parallel to the coordinate axes then renders those samples as the faces of a cube, which can be seen in Fig. 2. The proposed method is to extend this process to an arbitrary curve in spatial coordinates and visualize the spectra along the curve in several ways. One potential benefit for this is the ability to trace features in the image that are not linear, which is often the case in natural scenes. This non-linear sectional can be rendered in an image, where one dimension is length along the curve, another is wavelength and the intensity is data value of the image at that path length and band. Another possible rendering is as a 3D plot where the height is the data value at a certain path length and band. Part of the development stage will be to evaluate the curve representation, path parameterization and interpolation schemes. Additionally, it is not necessarily the case that the regions of interest are already quantified, and they may change with further visualization. Thus, interaction and real-time rendering are chosen to be a focus to the method.

In the development phase, the tractability of the tool will be investigated, and the aim is to create a full featured tool that allows the user to access real-world image data, interact and render in real time, and save and load intermediate measurements for later analysis. With the success of this stage, the strengths and limitations can be investigated. Real-world remote sensing images will be tested with the application to understand the performance with images of high spatial and spectral resolution. Then, a group of several expert geologists will evaluate the tool with respect to scientific problems that are part of their work. The aim is to compare similar parts of traditional image analysis workflows with the developed tool and gain a qualitative understanding of how they affect inferences. Finally, with the success of all of these parts, the application of the tool to existing image analysis environment will be investigated. Several tools offer a framework for third-party tool integration, and the utility of such an integration will be evaluated.



Figure 2: Image cube rendering of a hyperspectral image

#### 3 Significance

This project has numerous application in the scientific community including minerology, physics and agriculture. In particular, there have been a large number of data sets whose value has long before being exhausted. The missions collecting hyperspectral data include the terrestrial Airborne Visible/Infrared Imaging Spectrometer (AVIRIS), the Moon Minerology Mapper (M<sup>3</sup>), the lunar Observatoire pour la Minralogie, l'Eau, les Glaces, et l'Activit (OMEGA) and the Mars Thermal Emmission Spectrometer (TES). These projects focus on minerological and environmental analysis, which can potentially benefit from visualization methods that allow for better interactions with natural landmarks.

Outside of remote sensing, this method can be applied to any data with a planar domain that maps to data with a dimension of prohibitive size. For example, spatial maps of consumer preferences, political trends or environmental factors can generate high dimensional data that can be explored through the proposed method.

#### 4 Related Work

Existing visualization methods revolve around data reduction in the spectral or spatial domain. In spatial reduction, the spectrum a single location can be measured and plotted, and the spectra can also be measured at the boundary of a rectangle and rendered on the faces of a cube to give the impression of the spectral dimension mapping to depth[10], as in Fig. 1 and 2. Additionally, it is common to compute statistics of the spectra in a simple region of interest in the image. In the spectral domain, there are a number of methods for reducing the dimension of the data to three or fewer, to enable a color rendering on a standard display. Methods for spectral reduction include priciple component analysis, independent component analysis, band ratios and support vector machines, among many others[5, 11]. In addition to data dimension reduction, there have also been attempts to understand the variety of possible visualizations[3] and create experimental environments[1, 12, 6].

These methods are implemented in a variety of tools, which are part of a larger environment that includes many processing steps to calibrate, register and organize the image data. Some common packages include ENVI[7], ERDAS[2], Opticks[8] and SpecTIR[4]. The ENVI package includes a plugin framework, enabling thirdparty developers to integrate new tools into the workspace. This may offer a way to smoothly incorporate a new tool into the analysis chain.

#### 5 Research Plan

The schedule of the proposed work is to take place over the course of six weeks, with the following landmarks:

- Week 1: Work with collaborators to choose test data sets. Choose a development environment for the project. Learn data formats and find code to read it. Design and build a skeleton of the graphical user interface.
- Week 2: Design and develop curve representation and mechanism for interaction. Choose and develop interpolation scheme for image sampling.
- Week 3: Evaluate, design and develop image rendering of curve samplings. Develop code

to load and save curve descriptor and sample image data.

- Week 4: Evaluate, design and develop 3D height plot of sampled data. Test code with real data and investigate optimizations if necessary. Evaluate real-time performance.
- Week 5: Ask collaborators to prepare data for a comparision of methods. Meet with collaborators to demo and evaluate the tool qualitatively. Discuss efficacy and potential for integration into ENVI.

Week 6: Write up findings and present results.

#### References

- Shangshu Cai. Hyperspectral image visualization using double and multiple layers. PhD thesis, Mississippi State, MS, USA, 2008. Adviser-Moorhead, Robert J. and Adviser-Du, Qian.
- [2] ERDAS Inc. Erdas, inc. the earth to business company. http://www.erdas.com, September 2010.
- [3] Gupta MR Jacobson, NP. Design goals and solutions for display of hyperspectral images. Geoscience and Remote Sensing, *IEEE Transactions on*, 43(11):2684–2692, 2005.
- [4] SpecTIR LLC. Spectir: End to end hyperspectral solutions. http://www.spectir. com, September 2010.
- [5] Brian S. Penn. Using self-organizing maps to visualize high-dimensional data. *Comput. Geosci.*, 31(5):531–544, 2005.

- [6] Vito Roberto and Massimiliano Hofer. Theia: open environment for multispectral image analysis. In AVI '08: Proceedings of the working conference on Advanced visual interfaces, pages 462–465, New York, NY, USA, 2008. ACM.
- [7] ITT Visual Information Solutions. Envi software - image process and analysis solutions, September 2010.
- [8] Kip Streithorst. Opticks. http://opticks. org/confluence/display/opticks, September 2010.
- [9] James M. Torson. Interactive image cube visualization and analysis. In VVS '89: Proceedings of the 1989 Chapel Hill workshop on Volume visualization, pages 33–38, New York, NY, USA, 1989. ACM.
- [10] Fuan Tsai, Chun-Kai Chang, Jian-Yeo Rau, Tang-Huang Lin, and Gin-Ron Liu. 3d computation of gray level co-occurrence in hyperspectral image cubes. In EMM-CVPR'07: Proceedings of the 6th international conference on Energy minimization methods in computer vision and pattern recognition, pages 429–440, Berlin, Heidelberg, 2007. Springer-Verlag.
- [11] Thomas Villmann, Erzsébet Merényi, and Barbara Hammer. Neural maps in remote sensing image analysis. *Neural Netw.*, 16(3-4):389–403, 2003.
- [12] Jianting Zhang, Le Gruenwald, and Michael Gertz. Vdm-rs: A visual data mining system for exploring and classifying remotely sensed images. *Comput. Geosci.*, 35(9):1827–1836, 2009.

## DRAFT $(r_2)$ : Visualizing Relative Gene Expression Between Many Populations

#### S. Birch<sup>\*</sup>

October 6, 2010

#### Abstract

A new approach of comparing the gene expression profiles of many populations employing parallel coordinates is proposed. The objective in mind is twofold: one, that it aids in identifying anomalies in the data; two, that the visualization be suitable for publication and easy to interpret. These objectives are acheived through novel coordinate semantics and visual optimization techniques for clarity. The technique will be quantitatively evaluated by cross-validation with known results found from other techniques and qualitatively for ease of use and pattern-finding ability.

#### 1 **Background & Significance**

The space of a population's gene expression profiles (with n populations and mgenes) can be seen as a set of m points in  $\mathbb{R}^n$ , where each point is positioned on each population axis at the expression rate of its corresponding gene (so as to effectively represent a feature vector of expressivity in the space), such as shown in Figure 2. To explore the relationships between different populations each pair of coordinates is examined for a total of  $\binom{n}{2}^1$  comparisons. For the case of n = 2 these comparisons are trivially visualized on a scatterplot, as in Figure 1. These graphs are interpreted by looking at the distance of the points from y = x, indicating the degree that genes are expressed differently in the two populations (where a point perfectly on the line is expressed equally in both populations), and their absolute position on both axes, indicating the magnitude of expression. Further dimensions can be added by extending the plot to three dimensions or by encoding extra dimensions in color or size. For large numbers of populations, however, the problem becomes much more difficult to visualize as the pairwise comparisons scale  $\mathcal{O}(n^2)^2$ . Nevertheless, these comparisons are critical to understanding how gene expression varies across condition and person.

<sup>\*</sup>sbirch@cs.brown.edu

<sup>&</sup>lt;sup>1</sup>i.e. n choose 2, the number of unique pairs selectable from n axes. Given by  $\frac{n!}{2!(n-2)!}$ . <sup>2</sup> $\binom{n}{2} = \frac{n!}{2!(n-2)!} = \frac{1}{2}(n-1)(n-2) = \frac{1}{2}n^2 - \frac{3}{2}n + 1 = \mathcal{O}(n^2)$ 



Figure 1: Gene expression in two populations  $(Foxp^- \text{ versus } Foxp^+)$  on a loglog scale. Genes on y = x are expressed with the same magnitude. Figure 2d in [2].

Current approaches focus on either an analytic approach, such as clustering [CITE] or by visualizing standard scatterplots after applying dimensionality reduction techniques to the data [CITE]. Popular techniques for dimensionality reduction include principal component analysis, multidimensional scaling and singular value decomposition [1]; the first is shown in Figure 3. A problem is that dimensional reduction functions by choosing data to throw away. Slonim notes: "recall that data-reduction and visualization tools are projecting many thousands of dimensions into two or three may prevent frustration if the reduced data fail to capture the expected aspects of a data set" [1].

#### 2 Specific Aims

This proposal hypothesizes that visualizing the high-dimensional space with parallel coordinates will allow for effective analysis of larger numbers of populations. This approach will present all the data at once, eliminating the choice of data to omit and the potential dangers associated with that choice. The primary goal is therefore clarity and organization of the massive amount of data to be represented. This appears to be unprecedented in the literature. The proposed visualization contructs one axis for each pair of populations and the distance metric is plotted for each gene. The distance metric encodes how far the gene is from being equally expressed in both populations. Given  $M_A$  and  $M_B$ , magnitudes of expression in populations A and B respectively, it is given by:



Figure 2: Plot of (random) genes in three populations' expression magnitude space. Note that the scale is arbitrary.



Figure 3: Three dimensional PCA plot of the gene expression profiles of various cancer types. Figure 2a in [1].



Figure 4: Geometric interpretation of the distance metric on the projected plane. Note the basis formed by R and O.

$$\delta(M_A, M_B) = \left\| proj_E \left( \begin{array}{c} M_A \\ M_B \end{array} \right) - \left( \begin{array}{c} M_A \\ M_B \end{array} \right) \right\|_2$$

Where E is the line A = B. A geometric interpretation on the A - B plane is given in Figure 4.

These axes are arranged in parallel and ordered to maximize the clarity of the resulting graph. Dimension reordering has been studied in the context of reducing clutter in [3]. Whether the order is more important to preserve the locality of populations (in order to make population-wide trends apparent) or to minimize clutter is a novel question we hope to address. Finally, additional clarity may be lent by "bundling" genes together to reduce the space used by similar genes. [4] has implemented "visual clustering" using bundling techniques by curving the lines between axes (see an example in Figure 5).

#### 2.1 Benefits

The most significant benefit of parallel coordinates is that it scales to an arbitrary number of dimensions while still being interpretable in two dimensions (which also allows for publication). The cost of this is that interpretation is often more difficult as direct comparison is only explicitly shown between adjacent axes. This problem is resolved here by using  $\binom{n}{2}$  axes, allowing for the direct comparisons to be shown on the axis itself. The cost of this is that all the magnitude information in the R coordinate is lost. Because the goal is to look at relative expression however, the absolute magnitude isn't relevant (it is also possible to encode it in another dimension, such as color)<sup>3</sup>.

 $<sup>^{3}</sup>$ Note how this differs from the data omitted in clustering or dimensionality reduction, however: this is an a priori choice about the meaning of the visualization rather than an a posteriori selection of data to represent.



Figure 5: Before (a) and after (b) "visual clustering," as presented in [4].

#### 2.2 Evaluation

The result of the proposal will be evaluated quantitatively and qualitatively at the completion of the timespan. Experts in bioinformatics will qualitatively asses the software for ease of use with regard to finding anomalies in the data. Quantitatively, the method will be cross-validated using known anomalies from other techniques.

#### 3 Timeline & Methods

#### Week 1

Acquire data. Set up environment, choose technology platform.

#### Week 2:

Initial visualization prototypes: parallel coordinates alone and with relative axes (without optimizations). Evaluate feasibility and clarity of graphs.

#### Week 3:

Optimize visualization with regard to dimension reordering.

#### Week 4:

Optimize visualization with regard to bundling and clustering.

#### Week 5:

Interactivity, user tests, evaluation.

#### Week 6

Write up and disseminate results.

#### References

- D. Slonim, "From patterns to pathways: gene expression data analysis comes of age," *Nature Genetics*, vol. 32, pp. 502–508, 2002.
- [2] J. Hill, J. Hall, C. Sun, Q. Cai, N. Ghyselinck, P. Chambon, Y. Belkaid, D. Mathis, and C. Benoist, "Retinoic acid enhances Foxp3 induction indirectly by relieving inhibition from CD4+ CD44hi Cells," *Immunity*, vol. 29, no. 5, pp. 758–770, 2008.
- [3] W. Peng, M. Ward, and E. Rundensteiner, "Clutter reduction in multidimensional data visualization using dimension reordering," in *Information Visualization*, 2004. INFOVIS 2004. IEEE Symposium on. IEEE, 2005, pp. 89–96.
- [4] H. Zhou, X. Yuan, H. Qu, W. Cui, and B. Chen, "Visual clustering in parallel coordinates," in *Computer Graphics Forum*, vol. 27, no. 3. John Wiley & Sons, 2008, pp. 1047–1054.

#### Interactive Maps for Functional Brain Connectivity Queries

PI: Steven R. Gomez<sup>\*</sup> Co-PI: Ryan Cabeen<sup>†</sup> Collaborator: Jeff Chi-Tat Law<sup>‡</sup>

September 30, 2010

#### Abstract

We propose to build an exploratory map environment for visualizing neural connectivity queries in the human brain, and integrating these queries with analytics tools for tracking user hypotheses and evidence regarding functional connections. This work will help brain researchers more effectively target neuronal relationships for further connectivity experimentation. We will conduct a preliminary evaluation of our proposed tool with neuroscientists studying brain connectivity, and analyze design choices for dense circuit visualizations in this domain.

<sup>\*</sup>Department of Computer Science, Brown University, steveg@cs.brown.edu

<sup>&</sup>lt;sup>†</sup>Department of Computer Science, Brown University, rpc@cs.brown.edu

<sup>&</sup>lt;sup>‡</sup>Schnitzer Lab (Depts. of Biology and Applied Physics), Stanford University, lawjeffw@gmail.com

#### 1 Introduction

Researchers in neural circuitry are concerned with understanding connections between spatially distributed and functionally differentiated parts of the brain. Insight from these connections may give way to new medical treatments and diagnostic techniques for neurological disorders. While scientists have access to a growing collection of experimental connectivity data collated and disseminated on the web, the scale and complexity of these data make it difficult to gain insight from individual, textual queries – the standard interface to database access in current tools [3, 4]. Our work proposes a visualization tool for neural connections that communicates many connections as a 2D map with anatomical landmarks, allowing the user to filter connections of interest and quickly retrieve curations about these connections.

#### 1.1 Circuit Analysis

In [5], the authors give a broad overview of graph theory relevant to brain circuit analysis. The motivation for analysis like this is described in [2], which outlines a research agenda for studying brain connectivity and the tool-building to facilitate that research. An application of such tools for the neural network of C. elegans is demonstrated at [1], which shows a lightweight graph visualization of cell connections and allows the user to walk through this graph interactively. The motivation and data for this work comes from a study [6] investigating a hypothesis about evolution toward "wiring optimization" in the neuron circuit. For our purposes, this work demonstrates the kind of hypothesis-exploration-validation environment we want to support for domain scientists studying human brain circuitry.

There are many visualization tools that try to expose structure in biological systems for analysis by domain scientists. In [9], proteomic data and interaction networks are presented in an interactive graph framework. Tools for interacting with brain data – the domain of the proposed project – have used multiple linked views of DTI models, spanning more concrete/anatomical representations (e.g. streamtubes) and abstract ones (e.g. dendrograms) [7].

#### 1.2 Information Visualization

The project will consider current methods in information visualization, including embeddings of higher dimensional data into map representations, graph drawing techniques and interaction. Jianu et al. have created 2D visualizations of structural brain connectivity (DTI streamtubes) and disseminated them using the familiar interface of the Google Maps API [8]. We will explore embedding functional connectivity (i.e. projections) into maps like these in creating a tool for 'full network' connectivity, as described in [5], for brain researchers.

We also must take care in building interaction methods on network abstractions. With graphs at the scale of the human brain, representing individual axons will be a "data deluge" and may hide interesting or insightful patterns in the connection data. In [12], the authors examine filtering methods for user navigation of details in dense graphs. Users should also be able to interact with and rearrange graph layouts to investigate topological hypotheses, and recent work in this kind of subgraph interaction and manipulation is explored in [10].

#### 2 Contributions

We propose a map-style environment that allows the user to explore, filter, and query neuronal connections by anatomical coordinates and function. Our contributions are as follows:

- 1. A novel visualization tool for multiple neural connection display over a 2D anatomical map. A module for tracking hypotheses and evidence from queries will be integrated in the environment to facilitate scientific workflow using the tool.
- 2. A preliminary evaluation using a "think aloud" protocol with neuroscientists at Stanford to assess the effectiveness of our tool. We will develop a rubric for insight measurement similar to the evaluation in [11] to quantify user performance during a session with our tool.

The proposed project is *significant* because it will release a major bottleneck in access to collated brain circuit data by allowing visual, multi-connection analysis supported by sense-making analytics features. The *impact* of this work is high because our tool will be evaluated and integrated by the research team of Drs. Jeff Law and Mark Schnitzer at Stanford University. We will extend this evaluation to brain researchers at Brown if time permits. Furthermore, the proposal establishes an interdisciplinary collaboration between computer science and biology labs at these institutions, and may lead to further deployment of computational tools to scientific domains, which is healthy for each discipline.

#### 3 Plan

We will deliver the following:

- Week 1: BAMS circuit data cleaned and written into database (relational or other).
- Weeks 2–3: Graph drawing for connections over map with anatomical landmarks. Filtering and interaction in the graph tool.
- Week 4: Hypotheses-tracking and analysis support module.
- Week 5: User evaluation and plan for iterative improvements.
- Week 6: Final report and tool dissemination.

Our aim is to continue after the semester with a quantitative evaluation that demonstrates the usefulness of our tools, and prepare a submission for IEEE InfoVis 2011 (deadline: late March) describing our results.

#### 4 Facilities

The Linux workstations provided by the Computer Science Department at Brown will be used to develop and deploy our software; workstations in the Schnitzer Lab will be sufficient for evaluating our tool. No other facilities will be required.

#### References

- [1] Wormweb. http://wormweb.org/.
- [2] Jason W. Bohland, Caizhi Wu, et al. A proposal for a coordinated effort for the determination of brainwide neuroanatomical connectivity in model organisms at a mesoscopic scale. *PLoS Comput Biol*, 5(3):e1000334, 03 2009.
- [3] M. Bota, H.-W Dong, and L.W. Swanson. Brain Architecture Management System. Neuroinformatics, 3(1):15–48, 2005.
- M. Bota and L.W. Swanson. Online workbenches for neural network connections. Journal of Comparative Neurology, 500(5):807–814, 2007.
- [5] Ed Bullmore and Olaf Sporns. Complex brain networks: graph theoretical analysis of structural and functional systems. *Nature Reviews Neuroscience*, 10(3):186–198, 2009.
- [6] Beth L. Chen, David H. Hall, and Dmitri B. Chklovskii. Wiring optimization can relate neuronal structure and function. *Proceedings of the National Academy of Sciences*, 103(12):4723– 4728, March 2006.
- [7] R. Jianu, C. Demiralp, and D.H. Laidlaw. Exploring 3D DTI Fiber Tracts with Linked 2D Representations. *IEEE Transactions on Visualization and Computer Graphics*, 15(6), 2009.
- [8] Radu Jianu, Cagatay Demiralp, and David H. Laidlaw. Exploring brain connectivity with two-dimensional neural maps. *IEEE TVCG*, 2010. In Review.
- [9] Radu Jianu, Kebing Yu, Vinh Nguyen, Lulu Cao, Arthur Salomon, and David H. Laidlaw. Visual integration of quantitative proteomic data, pathways and protein interactions. *IEEE Trans. on Visualization and Computer Graphics*, September 2009.
- [10] Michael J. McGuffin and Igor Jurisica. Interaction techniques for selecting and manipulating subgraphs in network visualizations. *IEEE Transactions on Visualization and Computer Graphics*, 15:937–944, 2009.
- [11] Trevor M. O'Brien, Anna M. Ritz, Benjamin J. Raphael, and David H. Laidlaw. Gremlin: An interactive visualization model for analyzing genomic rearrangements. In *Proceedings of IEEE InfoVis*'10, Salt Lake City, UT, USA, 2010.
- [12] Frank van Ham and Adam Perer. Search, Show Context, Expand on Demand: Supporting Large Graph Exploration with Degree-of-Interest. *IEEE Transactions on Visualization and Computer Graphics*, 15:953–960, 2009.

## Steven R. Gomez

Department of Computer Science Brown University Box 1910 Providence, RI 02912 USA	<i>office:</i> CIT 423 <i>cell:</i> (802) 272-4941 <i>e-mail:</i> steveg@cs.brown.edu www.cs.brown.edu/~steveg	
Scientific visualization, visual computing and analytics, human-computer interaction, computer graphics, vision		
Brown University, Providence, Rhode Island US	A	
Second-year graduate student in doctoral prog	ram in Computer Science	
<ul> <li>Courses completed: Computer Graphics, Computer Vision, Computational Pho- tography, Distributed and Parallel Computing, Networking and Distributed Sys- tems</li> </ul>		
Dartmouth College, Hanover, New Hampshire USA		
B.A. magna cum laude, Computer Science, Ju	ne 2007	
<b>M2S, Inc.</b> , Lebanon, NH – <i>Software Engineer</i> – 2008–09 Development for several tools, including virtual stent grafts for Preview aortic aneurysm visualization software.		
Brown University Graduate Fellowship – 2009–10 Rufus Choate Scholar (top 5% of class) – 2006–07 John G. Kemeny Undergraduate Computing Prize – First Place, 2006, for <i>BlitzChat</i> Dartmouth Presidential Research Scholar – 2005–06 Class of 1928 Endowed Scholarship – 2003–07		
Steven R. Gomez, Radu Jianu, and David H. Laic Interface for White Matter Tractography. In Procee	Ilaw. A Fiducial-Based Tangible User edings of <i>ISVC</i> , 2010.	
Steven R. Gomez. Interacting with Live Preview Frames: In-Picture Cues for a Digital Camera Interface. In Proceedings of <i>ACM UIST</i> (Poster), 2010.		
Keller, R., Hunt, M., Jones, S., Morrison, D., Wolir Gary: Design Abstractions for a Jazz Improvisatio <i>Comput. Sci.</i> 193 (Nov. 2007), 47-60.	n, A., and Gomez, S. 2007. Blues for on Assistant. <i>Electron. Notes Theor.</i>	
Languages: Java, MATLAB, C, C++, Scheme, Perl, Tcl/tk, PHP, SQL, X/HTML, Languages: Java, MATLAB, C, C++, Scheme, Perl, Tcl/tk, PHP, SQL, X/HTML, Languages: Applications: Eclipse, Netbeans, Xcode, Apache, CVS, Subversion, MS Office, Adobe CS CS Operating Systems: Mac OS X, Linux, Windows 98/2000/XP		
	<ul> <li>Department of Computer Science Brown University Box 1910</li> <li>Providence, RI 02912 USA</li> <li>Scientific visualization, visual computing and an computer graphics, vision</li> <li><b>Brown University</b>, Providence, Rhode Island US Second-year graduate student in doctoral prog • Courses completed: Computer Graphics, C tography, Distributed and Parallel Computi- tems</li> <li><b>Dartmouth College</b>, Hanover, New Hampshire U B.A. <i>magna cum laude</i>, Computer Science, Ju</li> <li><b>M2S, Inc.</b>, Lebanon, NH – <i>Software Engineer</i> – 2 Development for several tools, including virtual ste- visualization software.</li> <li>Brown University Graduate Fellowship – 2009–100 Rufus Choate Scholar (top 5% of class) – 2006–00 John G. Kemeny Undergraduate Computing Prize Dartmouth Presidential Research Scholar – 2005 Class of 1928 Endowed Scholarship – 2003–07</li> <li>Steven R. Gomez, Radu Jianu, and David H. Laid Interface for White Matter Tractography. In Proceed Steven R. Gomez. Interacting with Live Preview I Camera Interface. In Proceedings of <i>ACM UIST</i> ( Keller, R., Hunt, M., Jones, S., Morrison, D., Wolin Gary: Design Abstractions for a Jazz Improvisatio <i>Comput. Sci.</i> 193 (Nov. 2007), 47-60.</li> <li><i>Languages:</i> Java, MATLAB, C, C++, Scheme, Pe <i>Applications:</i> Eclipse, Netbeans, Xcode, Apache, CS <i>Operating Systems:</i> Mac OS X, Linux, Windows S</li> </ul>	