# Safeguarding Abila through Multiple Data Perspectives
## VAST 2014 Grand Challenge Award: Effective Analysis and Presentation

Parang Saraf[*]        Patrick Butler[†]        Naren Ramakrishnan[‡]

Discovery Analytics Center
Department of Computer Science
Virginia Tech

## ABSTRACT

We introduce a system for visual analysis of news articles, emails, GPS tracking data, financial transactions and streaming micro-blog data. This system was developed in response to the 2014 VAST Grand Challenge and comprises of several interfaces for mining textual, network, spatio-temporal, financial, and streaming data.

**Index Terms:** H.5.2 [Information Interfaces and Presentation (e.g., HCI)]: User Interfaces—Interaction styles (e.g., commands, menus, forms, direct manipulation)

## 1 INTRODUCTION AND PROBLEM OVERVIEW

The VAST 2014 grand challenge describes a hypothetical scenario wherein some of the employees of an imaginary organization, GAStech have gone missing and it is speculated that an environmental activist group, Protectors of Kronos (POK) is responsible behind the disappearance. The challenge requires analysis of a variety of datasets in order to provide crucial leads to law enforcement agencies about suspicious persons, locations, organizations, events, etc. The provided dataset includes unstructured news articles, email headers from GAStech's company email, GPS tracking data of company cars assigned to employees, financial transaction data of employees' credit card & loyalty card, streaming data from a micro-blogging website and control room data of law enforcement agencies. The approach required to solve the challenge is akin to completing a jigsaw puzzle. Analysis of each of the datasets reveals partial information about the plot. An analyst needs to put all these pieces together in order to reveal the complete plot.

## 2 SYSTEM DESIGN

We developed a web-based visual analytics system that provides a set of tools and widgets specifically designed to analyze the given datasets. The system makes use of Google maps to display geospatial data and the Javascript-based graphical libraries `d3.js` [1] & `nvd3.js` for charts and graphs.

### 2.1 Unstructured News Articles

Figure 1 demonstrates the interface used for analyzing news articles. The interface makes use of a Python based search engine, Whoosh that allows logical queries and returns relevant news articles where similar articles are grouped together. The returned news articles are analyzed further using the Stanford's Named Entity Recognition tool [3], the results of which are displayed using word clouds. A time series plot shows keyword frequency over time, thereby giving insights into correlations between searched keywords.

---

[*]e-mail: parang@cs.vt.edu
[†]e-mail: pabutler@vt.edu
[‡]e-mail: naren@cs.vt.edu

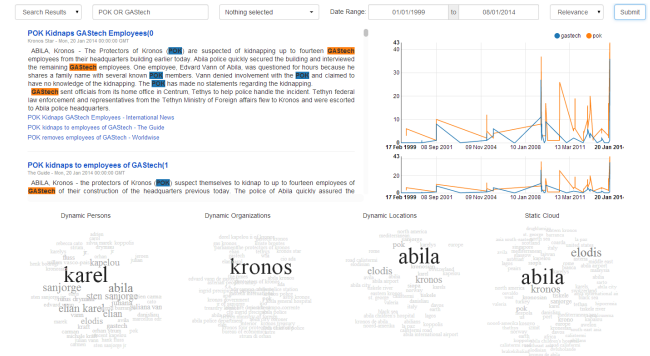Figure 1: News Interface

### 2.2 Email Headers

The email header information (see Figure 2) is analyzed using the following three widgets: radial graph for examining the underlying email network structure, co-occurrence matrix for identifying users who frequently exchange emails between them and a search interface for querying the emails. The co-occurrence matrix uses spectral co-clustering [2] to group users together and was implemented using the Python based scikit-learn package [4].
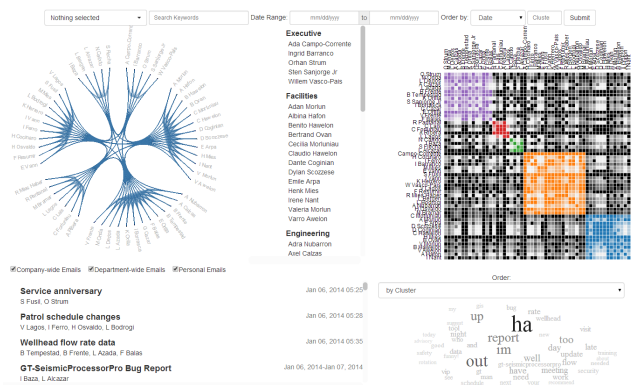


Figure 2: Email Header Interface

### 2.3 GPS Tracking Data

Three different interfaces were built to process GPS tracking data. The first interface (see Figure 3) allows for playback of GPS coordinates. Using this interface an analyst can visualize the movement of cars over time.

The second interface (see Figure 4) examines user's location information along with financial transaction data to geo-locate recreational establishments on a map.
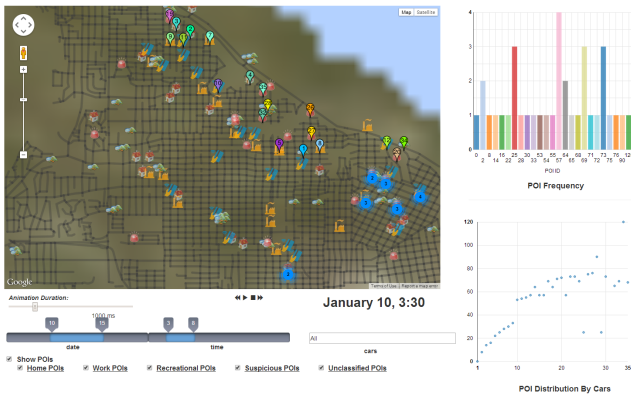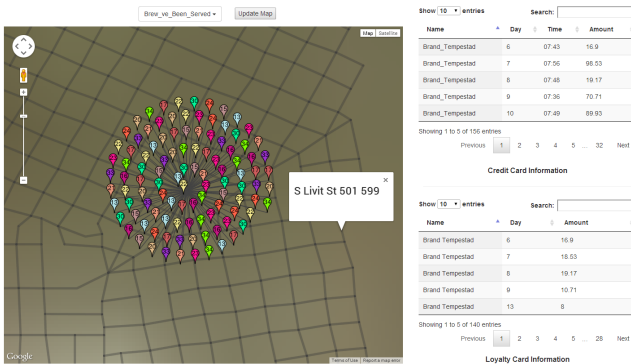
Figure 3: Location Playback Interface



Figure 4: Recreational POIs Interface

The third interface (see Figure 5) provides widgets to explore locations which are frequently visited by users. These locations, also known as Points of Interest (POI) can be classified as home, work, recreational and suspicious based on their distribution and frequency over time.
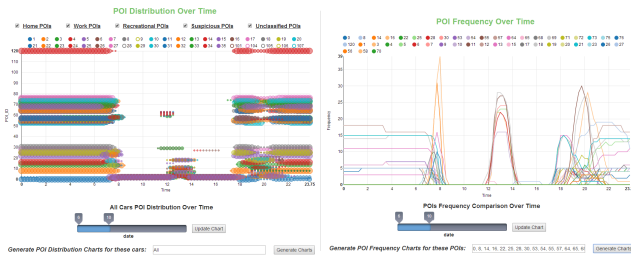


Figure 5: Points of Interest Analysis Interface

## 2.4 Financial Transaction Data

The interface for analyzing user spending data (see Figure 6) provides three widgets: employee vs employee spending comparison, employee spending distribution, and establishment sales distribution. As the name suggests, employee vs employee spending comparison compares employee spending across company as well as across departments that he/she works for. Employee spending distribution displays the spending of an employee across all establishments as well as across days. Similarly, establishment sales distribution compares the sales of a particular establishment across all employees as well as over days.
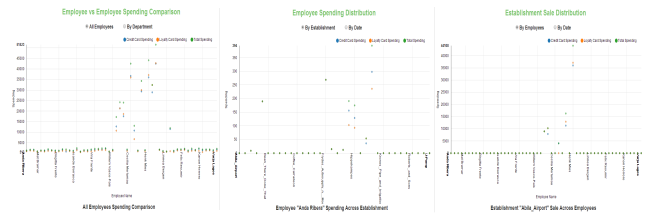


Figure 6: Financial Data Interface

## 2.5 Streaming Micro-blogging and Control Room Data

Figure 7 presents the interface for visualizing streaming data. In this interface, micro-blog records are displayed in different shades of red, where the gradient of the color defines the frequency with which a message is re-posted. Darker the color higher the frequency of re-posted message. An analyst can visualize geo-tagged posts on a map and can filter the micro-blog and control room data by searching for keywords. A line chart shows the frequency of searched keywords over time.
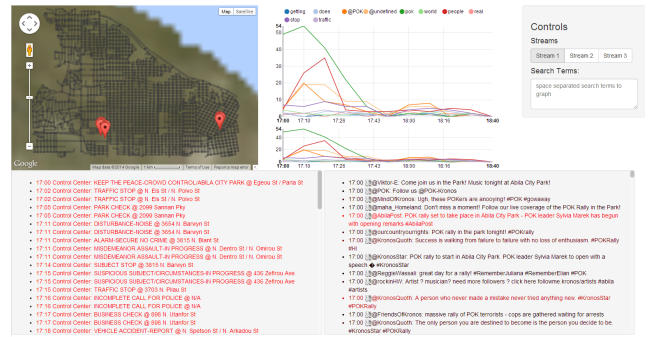


Figure 7: Streaming Data Interface

## 3 RESULTS

Using these different interfaces, we were able to unearth reliable evidence suggesting that Protectors of Kronos is indeed responsible for the disappearance of employees. Additionally, a few GAStech employees involved in suspicious activities were also identified.

### REFERENCES

[1] M. Bostock, V. Ogievetsky, and J. Heer. D3: Data-driven documents. *IEEE Trans. Visualization & Comp. Graphics (Proc. InfoVis)*, 2011.

[2] I. S. Dhillon. Co-clustering documents and words using bipartite spectral graph partitioning. In *Proceedings of the seventh ACM SIGKDD international conference on Knowledge discovery and data mining*, pages 269–274. ACM, 2001.

[3] T. G. Jenny Rose Finkel and C. Manning. Incorporating non-local information into information extraction systems by gibbs sampling. In *Proceedings of the 43nd Annual Meeting of the Association for Computational Linguistics (ACL 2005)*, pages 363–370, 2005.

[4] F. Pedregosa, G. Varoquaux, A. Gramfort, V. Michel, B. Thirion, O. Grisel, M. Blondel, P. Prettenhofer, R. Weiss, V. Dubourg, J. Vanderplas, A. Passos, D. Cournapeau, M. Brucher, M. Perrot, and E. Duchesnay. Scikit-learn: Machine learning in Python. *Journal of Machine Learning Research*, 12:2825–2830, 2011.