

# Top Health Trends: An information visualization tool for awareness of local health trends

**Sung Pil Moon**

Indiana University, IN, USA  
sungmoon@iupui.edu

**Yikun Liu**

Indiana University, IN, USA  
yikliu@iupui.edu

**Steven Entezari**

Indiana University, IN, USA  
sentezar@iupui.edu

**Afarin Pirzadeh**

Indiana University, IN, USA  
apirzade@iupui.edu

**Andrew Pappas**

MESH Coalition, IN, USA  
apappas.pu@gmail.com

**Mark S. Pfaff**

Indiana University, IN, USA  
mpfaff@iupui.edu

## ABSTRACT

We developed an intelligent information visualization tool to enable public health officials to detect health-related trends in any geographic area of interest, based on Twitter data. Monitoring emergent events such as natural disasters, disease outbreaks, and terrorism is vital for protecting public health. Our goal is to support situation awareness (SA) for personnel responsible for early detection and response to public health threats.

To achieve this goal, our application identifies the most frequently tweeted illnesses in a ranked chart and map for a selected geographic area. Automated processes mine and filter health-related tweets, visualize changes in rankings over time, and present other keywords frequently associated with each illness. User-centered visualization techniques of monitoring, inspecting, exploring, comparing and forecasting supports all the three stages of SA. An evaluation conducted with experts in health-related domains provided significant insights about awareness of localized health trends and their practical use in their daily work.

## KEYWORDS

Information Visualization; Health Trends Detection; Situation Awareness; Social Networking; Decision Aids.

## INTRODUCTION

Detecting and monitoring emergent crisis events such as natural disasters, disease outbreaks, and terrorism is vital to prevent negative impacts on public health. Researchers have applied many techniques to mine real-time data sources for early indications of such events. For example, to detect influenza outbreaks, previous attempts include analyses with Internet search query data (Polgreen, Chen, Pennock, & Nelson, 2008), emergency department visits (Yuan, Love, & Wilson, 2004), and telephone triage calls (Espino, Hogan, & Wagner, 2003).

Recently, researchers have started detecting and predicting emerging events using data from social network services such as Twitter due to its explosive growth rate, frequency, message volume, and public availability. Another stream of research using social media data focuses on geo-information visualization for foraging and sensemaking. Various information visualization tools in a variety of domains have been developed applying visual analytic techniques to effectively and efficiently present Twitter data. However, many of them have primarily focused on novel methods of capturing and rendering data on a map overlay, while there has been less attention on usability and overall ability to improve situation awareness (MacEachren et al, 2011; Hodgson, 2012a). Situation awareness is crucial to make effective decisions for appropriately controlling a broad range of complex, dynamic, and high-risk situations (Endsley, 1988).

This paper describes the Top Health Trends application, an information visualization tool for awareness of local health-related trends based on Twitter data. Our application displays the most frequently health-related tweets in both a ranking chart and a map. It also provides trend graphs of tweets and tag clouds of related keywords

mentioned with the tweets so that users of the tool can better aware of local situations. Results of evaluations with expert users in health-related domains are also presented. With this work, we make two main contributions:

- Providing insights about information visualization to be aware of local health-related trends combined using Twitter data and uses of visualization techniques to enhance each stage of situation awareness.
- Providing significant insights about awareness of localized health trends and their practical applicability to the work of public health professional.

The rest of the paper is organized as follows. We first introduce background information about trend detection with twitter data, information visualization tools with Twitter data, and situation awareness. Then descriptions of the tool we have developed are presented along with evaluations with expert users in health-related domains. The conclusion provides a summary, limitations, and future work.

## BACKGROUND

### TREND DETECTION WITH TWITTER DATA

Recently, there has been explosive growth of social media activity with hundreds of millions of active users. According to the Nielsen report (Nielsen Inc., 2011), social media and blogs reach 80% of all active Internet users in the US (245 million). In 2012, Twitter had 108 million accounts in the US, 465 million worldwide, and 175 million tweets daily (Hodgson, 2012b). Among those services, Twitter has received much attention as a new research data source due to message volume, frequency, and public availability (Culotta, 2010).

Achrekar, Gandhe, Lazarus, Yu, & Liu (2011) developed the Social Network Enabled Flu Trends (SNEFT) framework. This monitored Twitter messages to track and to predict the emergence and spread of an influenza epidemic in a population. They showed that Twitter messages could be used for almost real-time prediction of influenza activity trends two weeks earlier than current Influenza-Like Illness activity (ILI) reports by the Center for Disease Control and Prevention (CDC). Culotta (2010) identified that there was a high correlation (95%) between their influenza-related Twitter data analyzed and national health statistics (CDC weekly reports) to forecast future influenza rates. In the same vein, Signorini, Segre & Polgreen (2011) demonstrated that Twitter data could be used to estimate disease activity in real time, one to two weeks faster than the ILI reports. Gomide et al. (2011) also showed that their Dengue surveillance system using Twitter data could predict the dengue disease in Brazil one week earlier than the official data from Brazilian Health Ministry.

Early detection and prediction using Twitter data are not limited only to the disease outbreak domain. In the terrorism domain, Cheong & Lee (2011) proposed a framework that provides meaningful information visualization coupled with data mining and filtering techniques to identify terror threats based on Twitter data. Sakaki, Okazaki, & Matsuo (2010) proposed an earthquake detection system that used Twitter users as social sensors and monitored Twitter data to detect earthquakes, detecting earthquakes faster than the Japan Meteorological Agency (JMA) with high probability (96% of earthquakes). For marketing, Asur and Huberman (2010) analyzed tweets to outperform traditional box-office revenue prediction for movies. These examples support the validity of Twitter as a reliable data source for detecting ill-formed emerging events.

### INFORMATION VISUALIZATION TOOLS

Information visualization can be defined as the process of presenting information in ways that are naturally compatible with human visual capabilities (Gershon, Eick, & Card, 1998). Through information visualization tools, users can get benefits of improved efficiency, more interactivity, increased user satisfaction, reduced cognitive loads, and broadened insights about data and information. This has been demonstrated in areas including medicine (Chitarro, 2001), decision making (Lurie & Mason, 2007; Pfaff et al, 2010; Pfaff et al, 2012), and interaction design (Yi, Kang, Stasko, & Jacko, 2007).

Various information visualization tools using Twitter as a data source have been developed for a variety of purposes. Twitvis (D. Ferreira, Freitas, Rodrigues, and V. Ferreira, 2009) is an information visualization tool allowing Twitter users to identify and to discover social networks of both known and unknown users with related interests. To achieve this, their application provides two main visualizations: an egocentric network view showing a network of friends in a graph form, and the keyword view showing the relationship between the keywords of interests and users tweeting about them. TwitInfo (Marcus et al, 2011) is a microblog-based event tracking interface which can aggregate, summarize, and visualize tweets about user-specified events based on their social streaming algorithm for automatic identification of event tweet peaks. When a user clicks on a event

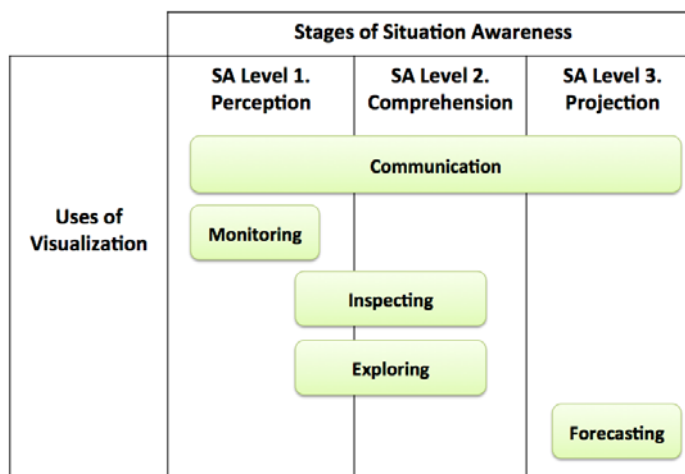
peak in the interface, it automatically filters and refreshes its map, links, tweet list, and sentiment graph in the time period of that peak. SensePlace2 (MacEachren et al, 2011) is an information visualization tool that is designed to support crisis management by providing overview and detail maps of tweets, filtering its place-time-attribute, and supporting an understanding of temporal and spatial patterns of activities, events, and attitudes.

## SITUATION AWARENESS

The term *situation awareness* (SA) is defined as “the perception of the elements in the environment within volume of time and space, the comprehension of their meaning, and the projection of their state in the near future” (Endsley, 1988, p. 97). Having higher SA is important to make effective decisions in urgent, dynamic, and high-risk situations. However, this may be hard to achieve in such situations. A number of researchers have investigated use of technologies to achieve or to improve SA such as the TSAS (Tactile Situation Awareness System; Rupert, 2000), an enhanced visualization tool for battlefield commanders, and the WIPER (Wireless Phone Based Emergency Response) system for managers in emergency operations centers (Madley, Szabo, and Barabasi, 2006).

Riveriro, Falkman, & Ziemke (2008) proposed an interactive methodology combining information visualization, data mining techniques, and human expert knowledge to detect abnormal anomalies in maritime traffic. However, in spite of these efforts, maintaining SA is still a challenging activity because of information overload, imperfect data from multiple data sources, complexities of systems and data mining techniques, and external factors such as high stress, time pressure, and inability to respond to unusual and unexpected incidents (Cameron, Power, Robinson, & Yin, 2012).

D’Amico & Kocka (2005) classified the use of visualizations to enhance the stages of situation awareness. There are five major uses of visualizations: monitoring, inspecting, exploring, forecasting, and communications (Figure 1). These uses are not isolated, but overlap, and should be implemented with proper user-interface techniques based on which of the three stages of situational awareness is being targeted.



**Figure 1: Relationship between the stages of situation awareness and the use of visualization (D’Amico & Kocka, 2005).**

## TOP HEALTH TRENDS APPLICATION

### OVERVIEW

Our application set out to improve the ability of public health officials to recognize public health crisis through Twitter feeds. The system, in its ideal state, would be able to take twitter feeds based on hundreds of key words and aggregate their data into actionable intelligence. We limit the definition of ‘intelligence’ to the notion that it enables users to acquire new knowledge, to enhance human cognitive performance, and to support users’ routine work tasks. This is accomplished by incorporating intelligently automated capabilities in our application with an effective combination of existing information visualization techniques. Feeds for this first version of the tool were set up using keywords extracted from the illness taxonomy of the Now Trending Challenge (ASPR, 2012) supplied by the Office of Assistant Secretary for Preparedness and Response at the United States Department of Health and Human Services. These keywords captured diseases such as pertussis, influenza and dengue fever; keywords were also included for potential public health crisis such as anthrax attacks and biological terrorism.

The application we aimed to build would identify the top trending health threats within geographic regions defined by the users, incorporating geolocated tweets. Output would give a level of how many tweets referenced a specific threat, the historical trends of that specific keyword, and what additional keywords were frequently included in tweets about a specific threat. To the end user in the public health profession, these data points allow users to interpret why a disease is trending and identify other contextual information from the twitter output.

Visualizing where trends are headed and how threats are impacting the health care system provides knowledge

not easily obtained by manually reading the content of tweets. Our tool produces a map identifying tweet clusters and allows users to drill down by opening individual tweets and user profiles. Allowing for regional and large area views allows users to compare cities during large events such as sports championships or major conferences. Although not part of the end product, we also planned to have alerting built into the product that could tie into schedules of mass public events. When tweets of a certain nature reached a threshold, the system could tweet at followers, email alerts to users, or text messages to public health personnel.

## SYSTEM DESIGN

Top Health Trends is a web-based application that mines Twitter data and produces a graphic analytical display of health-related tweets. It helps public health practitioners quickly identify trending health-related terms in a geographic area of interest by visualizing changes in rankings over time, and identifying associated keywords outside of the designated illness taxonomy.

Top Health Trends harnesses the power of the open-source Twitter API and MapQuest API to produce actionable intelligence and to support situation awareness among members of the health care, public health and emergency management communities. Top Health Trends stores Twitter data and provides analysis for the last seven days on a constantly rolling basis, updating every five minutes. Tweets are first selected according to the ASPR Now Trending Challenge illness term taxonomy. They are then tagged and processed for detailed term extraction and filtering. This is accomplished using the Natural Language Tool Kit (NLTK; <http://nltk.org>), which provides ready-to-use API entries for performing both cleansing and analyzing tasks.

The batch of tweets goes through several stages before final extraction. First, the tweets are grouped by keyword such that tweets regarding the same illness can be processed separately. All stop words are sifted so that only meaningful words are kept for further analysis. Typical stop words include pronouns (I, me, she), articles (the, a, an), conjunctions (and, or), and prepositions (by, on, in). Additional non-illness-specific words can also be added to the stop word list for stronger filtering, such as “feel,” “think,” “go,” or “get.” All words are converted to lower-case form and punctuation is removed, as well as words that do not comply with ASCII encoding. All words are then stripped of morphological affixes, leaving only the word stem (saw -> see, people -> person). Finally, distribution analysis is performed on all words to generate the ranked list of disease-related tweets.

Top Health Trends collects all tweets whether geolocation is enabled or not. Throughout our testing phase, the ratio of tweets with geotagged data to non-geotagged data was 1:100. The low incidence of geotagged tweets reduces the amount of data to analyze, but still provides more than enough tweets to identify current trends without any apparent systematic bias (i.e. we have no reason to believe that people with geolocation turned off tweet about health topics any differently than people who have that feature turned on).

## TWEET MAP

Analyzing trends using Top Health Trends is as easy as double clicking on the term in the Now Trending Ranking Chart (described further below). Upon the user's request, the website polls the Top Health Trends database for trending information related to the selected geographic area and updates the trend chart order based on the ranking of tweeted terms. The user can define a geographic area of interest by specifying a specific location in the “Current Location” bar or by using map controls to zoom and pan to a specific area (Figure 2). The location can be specified very flexibly, using any variety of common geographical expressions, such as zip code or city and state. The interactive map was based on the notion that interactive visualization tools that filter irrelevant information lead to decision quality improvements and reduced cognitive effort (Eick and Wills, 1995; Lurie and Mason, 2007). The interactivity provided in the

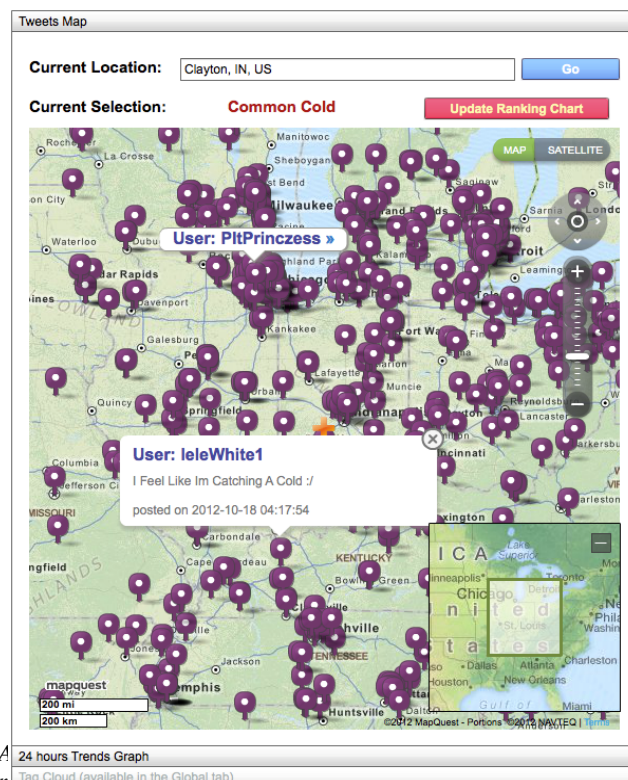


Figure 2: Tweet map of Indianapolis area showing all tweets related to common cold taxonomy.

tweet map follows the Keim’s mantra “Analyze first, Show the important, Zoom, Filter and analyze further, Details on demand” for knowledge discovery with combinations of analytical approaches with advanced visualization techniques (Keim et al., 2006, p. 16).

Tweets without geo-tags still provide valuable supplemental information to help interpret the Tweets known to come from a geographic region of interest. We created a unique screen to display the top illness taxonomy trends based on the global Twittersphere. This Top 5 Global List is provided to allow users to view all of the 25 ASPR taxonomy terms. This allows the end-user to compare their selected geographic area with more comprehensive trends. Of all the tweets available to the system, it currently only includes English-language tweets without hyperlinks. These criteria were determined early in our development due to the high number of false-positives associated with non-English tweets and hyperlinks. This tweet map enables visual exploring, monitoring, and inspecting to support the projection and comprehension stages of situation awareness.

**TOP TREND RANKING CHART**

To help users read and understand the data, Top Health Trends has adapted the familiar “Billboard Top 100” chart format (Figure 3). The use of this metaphor is intended not only to help the user form a structural understanding of the data and system but also to concisely represent the necessary situational awareness elements of the chart. Barr, Biddle, and Noble (2002) note that structural metaphors are crucial in that “the basic task of a user-interface is to help the user in coming to grips with the correct system image” (p. 26). This chart concisely displays the top five trending illnesses (numbers in red boxes), and for each illness, from top to bottom in the grey boxes, the previous day’s ranking, the number of days on the top five chart, and the peak rank during the current week. Our use of the Billboard Top 100 chart metaphor allows the user to assess the current situation, understand underlying trends, and support the users predictions in a familiar and accessible setting. All of this information is specific to the geographic location specified by the user on the map. A legend at the bottom of the screen shows users how to read the chart.

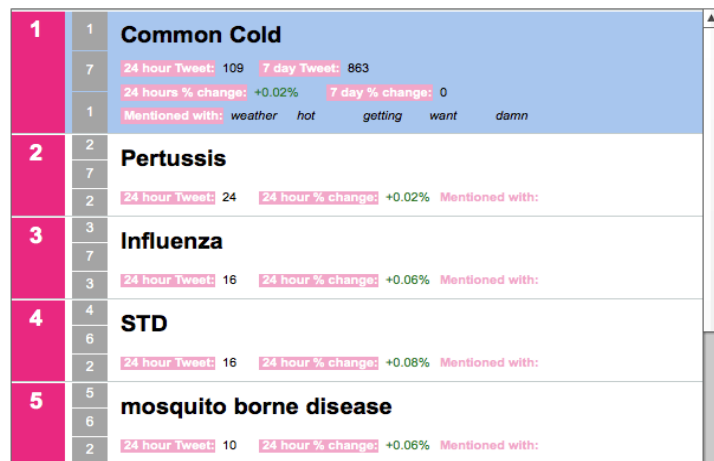
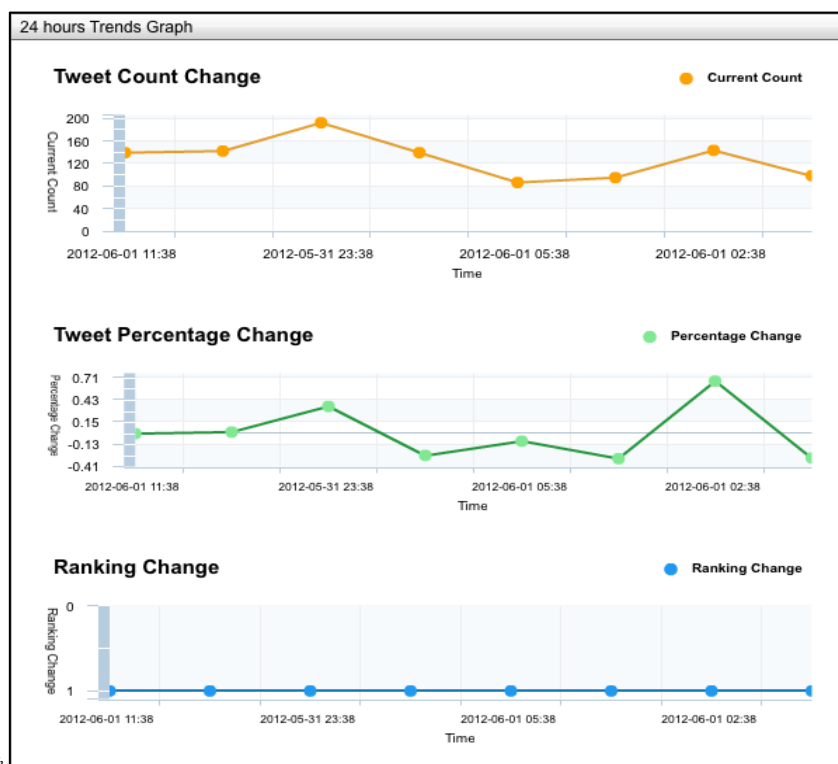


Figure 3: Local tweets trend in a “Billboard Top 100” format

**TREND GRAPH**

The user is also given the opportunity to drill down for specific content and trend information for a given disease term. By clicking on the name of any given disease, the map shown to the user will be augmented by points indicating the geolocation of those tweets. From here the user can gain a better understanding of why these diseases are trending. By clicking on the “Trends Graph” (Figure 4) option, the user will be presented with line charts displaying the tweet count, percentage change, and ranking change history over the past 24 hours in 3 hour increments (this increment was recommended



Proceedings of the 10<sup>th</sup> International ISCRAM Conference – Baden-Baden, Germany, May 2013  
 T. Comes, Figure 4: Graphic analysis of the seven-day trend for Common Cold, including daily count, percent change, and overall ranking change.

by an epidemiologist at the Indiana State Department of Health as most relevant interval to assess change). Line graphs were selected for these change graphs as they are among the most common ways to show trends (Robertson et al, 2008). The pattern of data is shown clearly in this type of graph, and they show how trends change over time. It enables users to make predictions of future values and it is also easy to read. This trend graph affords monitoring, inspecting, and forecasting supporting all three stages of situation awareness.

## RELATED KEYWORDS

Recognizing that tweets containing terms from the ASPR taxonomy are frequently unrelated to actual health issues, Top Health Trends also presents the user with the most frequently appearing terms that are tweeted along with illness taxonomy terms (Figure 5). This information provides useful context to help users determine the

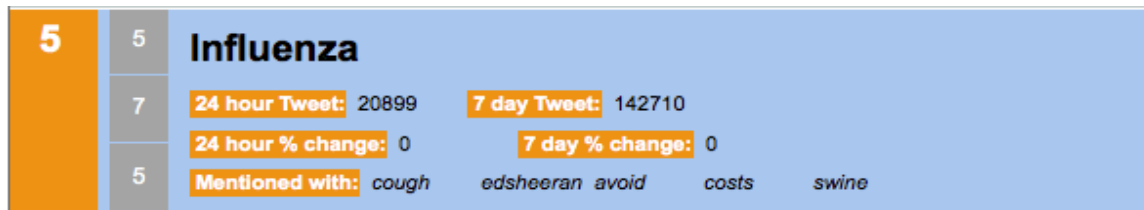


Figure 5: The related keywords mentioned with a selected illness provides context for the illness ranking.

accuracy of ranked terms. For example, there may be a group of Twitter users mentioning a celebrity who has HIV while another group of Twitter users could be discussing symptoms of HIV. Most of the existing tools that use Twitter data cannot distinguish these contexts. However, our tool provides users contextual hints so that they are able to better understand the context of the tweets and drill down for a closer inspection. An example that occurred during development was during the 2012 National Basketball Association Playoffs. Mentions of “flu” were often accompanied by basketball-related words like “Kobe” and “Michael Jordan” because of Kobe Bryant’s flu-ridden performances, not because any of those tweeting had flu symptoms. By revealing these related keywords, the user can perform additional research that may provide alternative explanations for seemingly high numbers of tweets related to a given health condition. After thorough testing of many possible thresholds, we determined that a minimum of 50 tweets is necessary to return a fruitful list of related keywords.

## TAG CLOUD OF THE DISEASE

We included tag clouds of the related keywords into our application to support better sensemaking of the tweets (Figure 6). A tag cloud is a visual representation showing a set of related words closely placed together in different sizes and colors. It is popularly used how it represents data in a limited space and uses tag size to draw attention to the most important items (Hearst and Rosner, 2008). We implemented this function based on the open-source tag cloud component provided in the Adobe Marketplace and Exchange (Adobe Inc., 2012). In our application, the tag cloud section was included to provide a broad keyword overview in addition to the five top related keywords listed for each illness. This tag cloud section shows up to 20 related keywords and their frequency of appearance. For example, in Figure 6, when the user places the mouse over any dot or keyword, it automatically shows its related information on the top-left side in the tag cloud section including the disease name the user selected in the ranking chart, the chosen keyword, and its count. This tag cloud facilitates visual exploring and inspecting that support the comprehension and projection stages of situation awareness.



Figure 6: Tag cloud for ASPR term Gastroenteritis

## EVALUATION WITH DOMAIN EXPERTS

We conducted a small-scale preliminary evaluation with experts involved in various aspects of public health. We considered experts as people with at least one year of working experience in a health-related domain. Our main goals were to understand possible uses of the application by domain experts, identifying how our application can aid their daily work, and finding unexpected problems and opportunities for improvement for the next generation of the tool. Rather than a rigorous usability evaluation typical of a more mature software product, this evaluation was designed as a way to check-in early with domain experts and verify whether the

tool addressed a genuine need in their professional work. Further, it was an opportunity to refine and prioritize the wish-list of new functionality in future versions of the tool, as well as bug-test the software with real users.

## PARTICIPANTS

Four participants were recruited from the MESH Coalition (a non-public-private partnership for healthcare providers to effectively respond to emergency events), Indiana University School of Medicine, and Indiana University Health. Participants were two EMS experts, one doctor, and one pharmacist. Three participants were female and one participant was male. Two participants said they used computers for their work more than 20 hours in a week while one used computers between 5 and 10 hours and the other between 15 and 20 hours.

## PROCEDURES

The procedures of the evaluation were to interact with the application through 19 simple tasks (Table 1), which gave users opportunities to experience features of the application, followed by a post-task questionnaire asking their comprehension of the information provided in the application, degree of helpfulness for being aware of local health trends, and degree of helpfulness to their daily work.

## RESULTS

The small sample size of this preliminary trial of the tool does not allow for any statistical analysis, but the responses on post-test questionnaires showed that users had positive feedback about the application on the whole. All participants were able to perform the 19 tasks correctly, and all found the application easy to understand, informative, visually pleasing, and easy to use. One participant mentioned it was “quite easy to use because it is self-explanatory.” They were also positive about the applicability of the application to their daily work in some degree. One participant specifically stated “This application has potential to be applicable in many areas. As a pharmacist, the tool will help me knowing when I need to order drugs.”

However, the two main concerns raised by multiple participants regarded the invisibility of data filtering techniques behind the visualization, and the credibility of intelligence generated by Twitter data. Participants were curious about how the Twitter data was filtered and how related keywords were associated each other. Because most frequently appeared keywords in tweets do not necessarily imply public health issues, some participants could not identify direct associations between a disease and its related keywords. To address this issue, some of the users suggested filtering the related keywords not only by the frequency, but also considering colloquial names and symptoms of the disease with an appropriate weighting strategy. Regarding the credibility issue, two participants specifically mentioned that it would be risky for expert users to make a decision if the information were solely based on unverified tweets. They suggested integration with official health data such as the U.S. Centers for Disease Control Influenza-like Illness (CDC ILI) report. Customizable components were suggested as a valuable feature aiding their daily work. Suggested features included customizable ranking and trend graphs, providing notifications as the number of tweets for certain disease exceeded a given threshold, and cross-referencing the application with external sources such as weather or news reports.

Component	Questions	Location
Ranking chart	<ul style="list-style-type: none"> <li>What is the name of the disease currently ranked number 3 in Indianapolis?</li> <li>How many days has the disease you selected been on the top 5 chart in this week?</li> <li>How many tweets the number 6 ranked disease have been counted for the past 7 days in this area (Indianapolis)?</li> <li>What percent change for 7 days did the disease you selected record?</li> <li>What are two keywords mentioned with the top disease?</li> </ul>	Indianapolis, IN
	<ul style="list-style-type: none"> <li>What is the name of the disease currently ranked 19th in the country?</li> <li>What is its peak rank?</li> <li>What are its 24 hour tweet count and 7 day tweet count?</li> </ul>	Global
Tweet map	<ul style="list-style-type: none"> <li>Go to New York, NY using the overview map component, and Zoom in to zoom level 10. Show all tweet icons mentioned about ‘Common Cold’ on the map</li> <li>How many 24 hour tweets are there of the disease ranked number 3 in New York?</li> <li>Go to San Francisco, CA using location textInput on the top right side of the map. Write any content of the ‘tweet’ mentioned about Influenza near San Francisco, CA.</li> </ul>	New York, NY San Francisco, CA
Trend graph	<ul style="list-style-type: none"> <li>Open the 24 hour trends Graph, Choose the disease currently ranked number 3.</li> <li>What was the rank of the disease 9 hours ago?</li> <li>When the time was the most significant change of tweet percent since yesterday?</li> <li>What was the lowest rank of the disease since yesterday? When was it?</li> </ul>	San Francisco, CA
Tag cloud	<p><i>Proceedings of the 10<sup>th</sup> International SIGRAPH Conference at Baden Baden, Germany, March 2010, 100,000. Comes, F. Fiedrich, S. Fortier, J. Geldermann and T. Müller, eds.</i></p> <ul style="list-style-type: none"> <li>What are five keywords mentioned with the disease you selected?</li> <li>What is the count of the top most mentioned keyword?</li> <li>What is the keyword least frequently mentioned?</li> </ul>	Global

## LIMITATIONS AND FUTURE WORK

A more rigorous evaluation would certainly have produced more detailed quantitative and qualitative results, but the time and expense of doing so was considered in the knowledge that we would most likely find results confirming limitations and future directions for the tool that had already come to light during the development process. Natural language learning was the most difficult challenge to overcome. Users often use slang, jargon, and abbreviations on Twitter, which impedes the tool's accuracy. Problems differentiating from common cold and cold weather, or gastrointestinal disease and a critique of poor workmanship were often evident. Our team wanted the application to begin learning the nuances of the Twitter style of natural language in order to select tweets that made more sense to public health than those meant solely for social interaction. Overcoming the language barrier would allow for a more precise selection of disease ranking tweets as well as a more accurate understanding of spatial spread of potentially related disease. This problem goes hand-in-hand with long-term aggregation of baseline term frequencies (longer than the development cycle for this first generation of the Top Health Trends tool). For example, terms related to the common cold will always occur far more frequently than the mention of anthrax, making the raw counts a poor indicator of which health issues need immediate attention. However, a single instance of anthrax is a far greater public health issue than 500 runny noses. Additionally, these baselines will vary significantly from region to region, causing geographically unique false-positives or false-negatives. For example, the Women's National Basketball Association team in Indianapolis, Indiana, is called the "Indiana Fever." Differentiating between health-related and basketball-related uses of the term "fever" in Indiana will require unique filtering methods not applicable in other geographic regions.

The value of this tool extends beyond the knowledge of what health terms are trending on Twitter. Often, public health entities work backwards to stop an outbreak weeks after it has begun. This is based on the delay from disease onset, clinical presentation, and epidemiological reporting to local or state authorities. Using Twitter as a

### **Table 1. Task questions during the experiment with domain expert participants**

surveillance tool allows clinicians to monitor in real time for unexpected spikes in symptoms or other disease-related events. With geolocated data, public health personnel can go to the location of a suspected outbreak and begin an investigation of what may be a rapidly evolving public health crisis, days before traditional methods suggest an outbreak. This allows for shortened time to response, saving people from secondary and tertiary infections and stopping the spread of a threat before excessive casualties result.

Currently, the Top Health Trends application relies on the tweets having geolocation information. Fewer than 1% of tweets have geolocation information (specific longitude and latitude), and only one-quarter of Twitter users provide a specific city location in their profile (Cheng, Caverlee, & Lee, 2010). While techniques are under development to infer the location of a tweet from its content, such as keywords and time stamps, this problem actually presents a gift in the form of data reduction. We have no evidence yet to suggest a systematic bias for eliminating tweets without specific geolocation information (that is, whether or not someone is tweeting about being ill is not expected to have any relationship with whether they have the geolocation function turned on). This smaller sample size dramatically increases the speed at which we can acquire, process, refine, and present tweets to the user. While the absolute number of tweets is smaller, the ranked frequencies of the keywords should remain the same, which in this application is most important. For this first version of the Top Health Trends tool, we determined that the significant effort for developing and implementing intelligent data mining techniques to extract additional geolocation information would not sufficiently improve the output. However, such strategies would be of great value when searching for terms of relatively low frequency, where any evidence of a particular illness could be useful, regardless of location.

The formal list of disease and symptom terms required for the first generation of this tool was based on terms for illnesses derived from a National Library of Medicine study and provided to Now Trending Challenge participants by ASPR. However, since tweeting is primarily an informal form of communication, we found minimal use of formal medical terms like "H3N2" or "pneumocystis." Terms like these would more likely be tweeted by a medical professional, or a person just diagnosed by a doctor. This has much less value for epidemiological surveillance than more realistic language (e.g. "I feel like crap and have this weird rash"), which are available from other sources and can be merged with the supplied list of formal terms. Sources of these additional terms include patient complaint data from the intake desk of emergency departments, as well as direct interviews of clinicians who directly interact with patients every day, and who are tasked with diagnosing a medical condition, in part, through a verbal report from the patient. We are currently planning a survey of area health professionals to develop a weighting system for the search terms, such that those which are the most likely indicators of a true health condition are weighted more heavily than those which have a higher tendency to be false positives. We would also like to supplement the taxonomy with additional terms recommended by health professionals who are familiar with the language used by patients regarding their health.

We would also like to improve accuracy of filtering algorithm for the tag clouds of related keywords. As



mentioned in the evaluations, our tag clouds are based on the frequency of additional terms co-occurring with terms in the supplied illness taxonomy. Using the tag cloud alone, users may not be able to recognize direct relationships between the disease and those co-occurring terms. Additional interactive tools could integrate additional sources of illness-related terms used in a variety of conversational settings. When users examine one of these keywords, the application could also present a pop-up or overlay of additional known terms and symptoms so that they can better determine the context of the tweets and determine the nature and validity of a possible threat to public health. The list of symptoms corresponding to a specific disease will be added based on previously validated taxonomies (Paul & Dredze, 2011). Integration with reliable official data sources for credibility and more customizable components will also be implemented following further consultation with public health professionals.

Lastly, the methods of data retrieval and visualization in this application can be readily repurposed for other uses. When the user requests information from the interface, an API call is made to the database. This call accepts the latitude/longitude boundaries, disease keyword, and time period of reference as parameters. Future applications, as well as future aspects of this application, can replace the disease keyword with other unique parameters to retrieve data within the boundaries of *any* specified topic of interest, not just public health.

## CONCLUSION

We have presented the Top Health Trends application designed to enable users to be aware of local health trends based on Twitter data. The system integrates a tweet map, top trend ranking charts, the trend graphs, and the tag cloud of related keywords. Although many limitations remain in this first version of the tool, results of evaluations with expert users in health-related domains showed that such an application has potential to be used in their daily work. We plan to continue to improve the accuracy and credibility of the application, as well as grow its feature set and support for interoperability with existing monitoring and surveillance systems used in the domain of public health.

## REFERENCES

1. Achrekar, H., Gandhe, A., Lazarus, R., Yu, S., & Liu, B. (2011) Predicting flu trends using Twitter data, *First International Workshop on Cyber-Physical Networking Systems*, 713-718.
2. Adobe Inc. (2012) Adobe Marketplace and Exchange. <http://www.adobe.com/cfusion/exchange/index.cfm>.
3. Asur, S. and Huberman, B. A. (2010) Predicting the future with social media. HP Labs. <http://www.hpl.hp.com/research/scl/papers/socialmedia/socialmedia.pdf>.
4. Barr, P., Biddle, R. and Noble, J. (2002) A taxonomy of user-interface metaphors, *Proceedings of the SIGCHI-NZ Symposium on Computer-Human Interaction (CHINZ 2002)*.
5. Cameron, M. A., Power, R., Robinson, B., and Yin, J. (2012) Emergency situation awareness from Twitter for crisis management, *Proceedings of the 21<sup>st</sup> Conference on WWW 2012 Companion*, 695-698.
6. Cheng, Z., Caverlee, J., and Lee, K. (2010) You are where you tweet: A content-based approach to geolocating twitter users, *Proceedings of the 19th ACM international conference on Information and knowledge management (CIKM '10)*, ACM, 759-768.
7. Cheong, M. and Lee, V. (2011) A microblogging-based approach to terrorism informatics: Exploration and chronicling civilian sentiment and response to terrorism events via Twitter, *Information Systems Frontiers*, 13, 1, 45-59.
8. Chitarro, L. (2001) Information visualization and its application to medicine, *Artificial Intelligence in Medicine*, 22, 81-88.
9. Gershon N., Eick S.G., and Card S. (1998) Information visualization, *ACM Interactions*, 5, 2, 9-15.
10. Culotta, A. (2010) Towards detecting influenza epidemics by analyzing Twitter message, *1st Workshop on Social Media Analytics (SOMA '10)*, Washington, DC, 115-122.
11. D'Amico, A. and Kocka, M. (2005) Information assurance visualizations for specific stages of situational awareness and intended uses: Lessons learned, *Visualization for Computer Security (VisSec 05)*, 107-112.
12. Eick, S.G. and Wills, G.J. (1995) High interaction graphics, *European Journal of Operational Research*, 81, 3, 445-459.

13. Endsley, M.R. (1988) Designing and evaluation for situation awareness enhancement, *Proceedings of the Human Factors Society 32nd Annual Meeting*, 97-101.
14. Espino, J. U., Hogan, W. R., and Wagner, M. M. (2003) Telephone triage: A timely data source for surveillance influenza-like diseases, *Proceedings of the AMIA Annual Symposium*, 215-219.
15. Hodgson, L. (2012a) Facebook 2012: The latest on everybody's favorite social network. <http://www.blogherald.com/2012/02/15/facebook-2012-infographic/>
16. Facebook, Inc. (2012) The power of Facebook advertising. <https://www.facebook.com/business/power-of-advertising>
17. Ferreira, D., Freitas, M., Rodrigues, J., and Ferreira, V. (2009) Twitvis – exploring twitter network for your interests, *UMA 2009*, Funcha, Madeira, Portugal, 1-8.
18. Field, K. and O'Brien, J. (2010) Cartoblography: Experiments in using and organizing the spatial context of micro-blogging, *Transactions in GIS*, 14 5-23.
19. Gomide, J., Veloso, A., Meira, W., Almeida, V., Benevenuto, F., Ferraz, F., and Teixeira, M. (2011) Dengue surveillance based on a computational model of spatio-temporal locality of Twitter, *Proceedings of the ACM WebSci 2011*, 1-8.
20. Hearst, M. and Rosner, D. (2008) Tag clouds: Data analysis tool or social signaller? *Proceedings of the 41st Hawaii International Conference on System Sciences (HICSS 2008)*
21. Keim, D., Mansmann, F., Schneidewind, J., & Ziegler, H. (2006). Challenges in visual data analysis, *Proceedings of the 10th International Conference on Information Visualization*, 9-16.
22. Kim, Y.J. and Hoffmann, C.M. (2003) Enhanced battlefield visualization for situation awareness, *Computer & Graphics*, 27, 6, 873-885.
23. Lurie, N. and Mason, C. (2007) Visual representation: Implication for decision making, *Journal of Marketing*, Vol.71, pp.160-177.
24. MacEachren, A.M., Robinson, A. C., Jaiswal, A., Pezanowski, S., Savelyev, A., Blanford, J., and Mitra, P. (2011) Geo-twitter analytics: Applications in crisis management, *Proceedings 25th International Cartographic Conference*, Paris, France.
25. Madley, G., Szabó, G., and Barabási, A. (2006) WIPER: The integrated wireless phone based emergency response system, *Proceedings of the International Conference on Computational Science*, 3, 417-424.
26. Marcus, A., Bernstein, M., Badar, O., Karger, D., Midden, S., and Miller, R. (2011) Twitinfo: aggregating and visualizing microblogs for event exploration, *Proceedings of the SIGCHI on Human Factors in Computing Systems*, 227-236.
27. Office of the Assistant Secretary for Preparedness Response (ASPR) (2012) Now Trending Challenge. <http://www.nowtrendingchallenge.com>
28. Paul, M. and Dredze, M. (2011) You are what you tweet: Analyzing Twitter for public health, *Proceedings of the 5th International AAAI Conference on Weblogs and Social Media (ICWSM 2011)*, Barcelona, Spain.
29. Pfaff, M. S., Klein, G. L., Drury, J. L., Moon, S. P., Liu, Y., and Entezari, S. O. (2012) Supporting complex decision making through option awareness, *Journal of Cognitive Engineering and Decision Making*, Advance online publication. doi: 10.1177/1555343412455799.
30. Pfaff, M., Drury, J., Klein, G., More, L., Moon, S., and Liu, Y. (2010) Weighing decisions: Aiding emergency response decision making via option awareness, *Proceedings of the IEEE International Conference on Technologies for Homeland Security*, 251-257.
31. Polgreen, M. P., Chen, Y., Pennock, M. D., and Nelson, D. F. (2008) Using Internet searches for influenza surveillance, *Clinical Infectious Diseases*, 47, 1443-1448.
32. Quincey, E.D. and Kostkova, P. (2009) Early warning and outbreak detection using social networking websites: The potential of twitter, *Proceedings of the eHealth*, 21-24.
33. Riveiro, M., Falkman, G., and Ziemke, T. (2008) Improving maritime anomaly detection and situation awareness through interactive visualization, *Proceedings of the 11<sup>th</sup> International Conference on Information Fusion*, Cologne, Germany, 1-8.

34. Robertson, G., Fernandez, R., Fisher, D., Lee, B., and Stasko, J. (2008) Effectiveness of animation in trend visualization, *IEEE Transaction on Visualization and Computer Graphics*, 14, 6, 1325-1332.
35. Rupert, A.H. (2000) Tactile situation awareness system: Proprioceptive prostheses for sensory deficiencies, *Aviation, Space, and Environmental Medicine*, 71, 9, 92-99.
36. Sakaki, T., Okazaki, M., and Matsuo, Y. (2010) Earthquake shakes Twitter users: Real-time event detection by social sensors, *Proc. of the 19<sup>th</sup> International World Wide Web Conference 2010*, Raleigh, NC, 851-860.
37. Signorini, A., Segre, M. A., and Polgreen, M. P. (2011) The use of twitter to track levels of disease activity and public concern in the U.S. during the influenza A H1N1 pandemic, *Plos ONE*, 6, 5.
38. Nielsen, Inc. Social media report: Spending time, money and going mobile. <http://blog.nielsen.com/nielsenwire/social/>
39. Hodgson, L. (2012b) Twitter 2012. <http://www.blogherald.com/2012/02/22/twitter-2012-infographic/>
40. Yi, J.S., Kang, Y, Stasko, J.T., Jacko, J.A. (2007) Toward a deeper understanding of the role of interaction in information visualization, *IEEE Transactions on Visualization and Computer Graphics*, 13, 6, 1224-1231.
41. Yuan, C.M., Love, S., and Wilson, M. (2004) Syndromic surveillance at hospital emergency departments – southeastern Virginia, *MMWR Morb Mortal Wkly Rep*, 53-56.