

# A Digital Library for Crowds on the Real-Time Social Web

Jeff A. McGee  
Texas A&M University  
College Station, TX 77843  
jeffamcgee@gmail.com

Krishna Y. Kamath  
Texas A&M University  
College Station, TX 77843  
kykamath@cs.tamu.edu

## ABSTRACT

In this project, we want to build a digital library for transient crowds in highly-dynamic social messaging systems like Twitter and Facebook. A crowd is a short-lived ad-hoc collection of users, representing a “hot-spot” on the real-time web. For example, an event like super-bowl might result in formation of crowds that are discussing various happenings in the game. Successful detection of these hot-spots can positively impact related research directions in online event detection, content personalization, social information discovery, etc. We will build a framework that allows an analyst or curious user to find interesting crowds and see how they evolve. This framework will have three main parts: a database to store crowds, users, and their messages; a set of crowd detection algorithms and filters; and a tool for searching for crowds by topic, geography, or user-name.

## 1. INTRODUCTION

## 2. RELATED WORK

Twitter<sup>1</sup> is a microblogging platform which is fast gaining popularity [21] among broad sections of society and has a global outreach spreading from developed, urban nations like the United States where it has a high adoption rate [13], to developing countries in parts of Asia and South America.

Along with working on understanding microblogging usage and communities [13], the main author - Akshay Java was one of the first few who dealt with the measurement of the usage and nature of communities in microblogging. In his latter study [12], he presented his observations of the microblogging phenomena and user intentions by studying the content, topological and geographical properties of such communities. He found that microblogging provides users with a more immediate form of communication to talk about their daily activities and to seek or share information.

---

<sup>1</sup><http://www.twitter.com>

When it comes to visualization of the social crowd, a tool which in some ways shows similar objectives as us and stands out among the rest is Vizster [8]. It is a visualization system for which can be used for exploration and navigation of the social networks of any scale. The system provides facilities to improve the functionalities like discovering people or connections which can help in promoting awareness of community structure and exposure to information. It maintains a fun element to it while providing these given features. It was aimed to provide specific end-users, functionalities like situated exploration and visual analysis of big-scale communities. It is based on a design of network layouts which have features like supporting visual search and analysis, exploring connectivity in large graphs and automatic identification and visualization of structures of communities.

Another related study was on how web-based retrieval from twitter can be used as knowledge base [5] which shows that the visualization system we try to attempt can be helpful. Marc Cheong et al. in this study tried to analyze demographics of the Twitter users and attempted to provide and analyze the trend in Twitter into what’s going on in the real world. They enabled us to understand the underlying characteristics behind the trend setting in Twitter and thereby providing a new perspective on the contributors of a trend. They used visualization techniques along with artificial intelligence-based data mining methods to classify the Twitter messages dealing with the trend topic.

Identifying highly dynamic ad-hoc collections of users what we refer to as crowds in massive social messaging systems like Twitter and Facebook is important [14]. Kamath et al. in the study suggest an efficient locality-based clustering approach for identifying crowds of users in near real-time compared to more heavyweight static clustering algorithm. The study consisted of 711,612 users and 61.3 million messages, and tells what approaches to follow to efficiently and effectively identify Twitter-based crowds.

In the other study based on the previous one, Kamath et al [15], they add another salient feature - a novel crowd tracking and evolution approach for linking crowds across time periods. Unlike the more static and long-lived group-based membership offered on many social networks their goal was to support the discovery of organic and highly-temporal group affiliation, which they refer to as “transient crowds”. A transient crowd according to them is a potentially short-lived ad-hoc collection of users bound together by some common thread - which can be communication-based, location-based or interest based. We are inspired by this idea of formation and detection of crowds. We try to

implement detection of crowd in Crowdy in a similar fashion.

Some other studies related come from areas like social network visualization, short text visualization, progressive disclosure, and user modeling. Considering social network visualization, studies like [4] by Boyd et al described features of on-line social networking sites and the history of the social networking sites, and the changes during their developments. Deeper understanding of the social networks can help boost the research in social network visualization. In Moreno's following work [19], he emphasized the use of pictorial representations to map social relations and the embedded patterns. He considered people as nodes in the network and relations between people as the edges in the network. Laumann et al [17] applied computational procedure, multidimensional scaling to locate nodes in visualizing network. The authors are also the first to plot the networks in three dimensions. Other examples of good visualizations can be Facebook Friend Wheel [1] which visualizes a Facebook user's network of friends in a wheel format and Trendsmap [2] which plots real time trending topics over two-dimensional (latitude, longitude) world map, and you can zoom in to see the trending topics for a specific location or an area.

In the area of short text visualization, studies like [3] tell that these social web-sites serve to fulfill their users' desires to interact with and help others. The authors also proposed a 3-level framework for understanding why members of online communities participate or not. According to the 3-level framework, users make comments or short-text post to have their voice be heard. Leggett and Shipman [18] proposed and described the characteristic of interactive communication in seven dimensions in a similar way to examine characteristics of communication for a commenting system like roles, voices, interaction, indirection, history, narrative and media. In [11], they tried to understand how an on-line community itself perceives the relative quality of its own user-contributed content and further used the information to promote valuable opinion from the flood of comments that attached to each social object. We compared three different visualization approaches including the traditional list view, the quantitative aggregation view and the personalized space view. We learned that, to best present massive meta-data such as user contributed comments to the end user, a personalized filtering model is needed and that a 2D space positional layout is desired such as in the case of Plurk and Opinion Space systems.

In the area of progressive disclosure, Phan et al built a prototype of a system that visualized packet traces from a computer network [22] which is similar to the problem investigated. They introduced the concept of progressive multiples which combines progressive disclosure with small multiples. In Scalability in System Management GUIs [6], the authors use semantic zooming to show information about network configuration at various levels of detail, which can tell about the way progressive disclosure can be used for system administrators. In An Interactive Ambient Visualization for Code Smells [20], the authors describe a method for alerting uses about code that may need refactoring. They created a plug-in for Eclipse that detected questionable code and alerts the programmer. This shows progressive disclosure for software developers.

Involving user models and visualization into the picture, studies like [23] tell that the individualization often depends

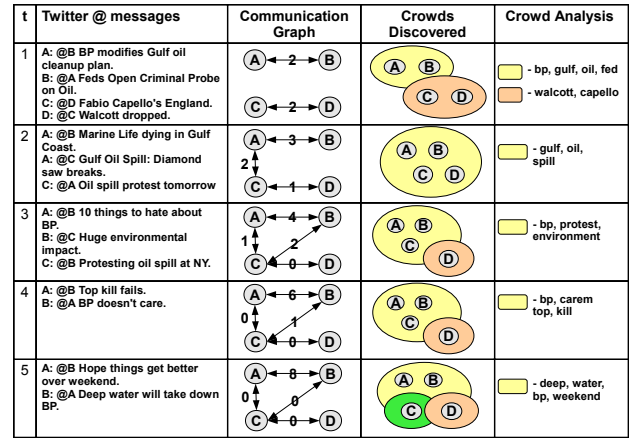


Figure 1: Example of crowd discovery and tracking in Twitter.

on the user's expectation of the user interface which might be adapted or adaptable or adaptive user interface. Papers like [10] explain that the user model selection is crucial for any tool and how the visualization can help in the discovering and building knowledge. Data and knowledge visualization module of such a system can be used to help the user interaction in other modules. This helps in better user interaction experience irrespective of which module of the system he/she is interacting with. The user's knowledge, preferences, goals etc. are generally included in the user model when visualization is concerned in such tools [9]. Also, critical factors like cognitive load on the user ought to be considered when such tools are made to ensure effectiveness. The identification of target users is also a important factor to be considered. Model-based tools like TADEUS can include explicit task and user models in their system, over potential to work in integrated user interface development environments and support the entire user interface life cycle. Although research papers like [23], talk about how different software tools can be used for helping the user designers create interfaces with consideration of user models, there have been efforts to combine user model with task and domain models to help developers in creating and designing user interface [24]. Systems like DB-USE along with TADEUS help us in understanding of how user model can be combined with other models to help in creating/designing better user interface for applications which are mainly concerned with visualization.

### 3. CROWD DISCOVERY ALGORITHM

To illustrate the problem of crowd discovery, consider the simple example in Figure 1. At time  $t=1$ , users A and B send messages to each other, as do users C and D.<sup>2</sup> The associated communication graph shows an edge between the two pairs, where for simplicity the edge is annotated with the number of messages between the users (2, in both cases). Further, suppose we identify crowds based purely on graph

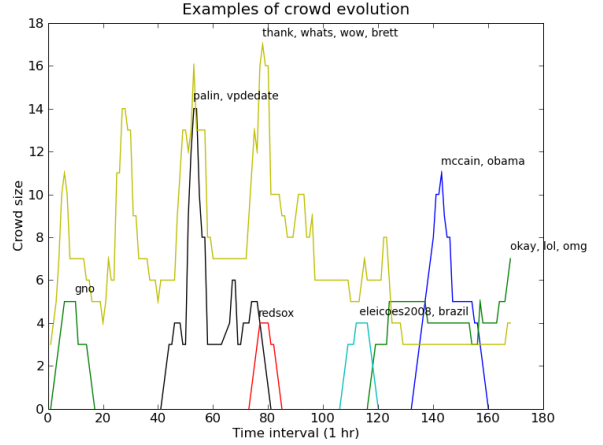
<sup>2</sup>For simplicity, the example discretizes time so that all messages between users occur in steps. In practice, the proposed algorithm relaxes this assumption and can handle arbitrary message sending times.

connectivity. So for time  $t=1$ , we see there are two crowds discovered  $\{A,B\}$  and  $\{C,D\}$ . For each crowd, we can characterize the semantics of their communication with simple keywords extracted from the content of the tweets: (“oil”, “gulf”) and (“walcott”, “capello”). At time  $t=2$ , the communication graph is updated with a new edge (connecting User A and User C), and the existing edges are decayed by one (again, a simplifying assumption for the purposes of this example). A single crowd is discovered since all users are connected via edges with non-zero edge weights. At time  $t=3$ , User D leaves the main crowd since no messages to or from User D have been observed since time  $t=1$ . This process continues until time  $t=5$  when User C also leaves the main crowd due to inactivity. Note that crowds are discovered from communication graph only and not from the content of the messages. As an example of crowd tracking, we can track the evolution of the yellow crowd across time periods, observing the changes it goes through as it grows in size from  $t=1$  to  $t=2$  and then reduces to two users by  $t=5$ .

The grouper deals with the transient crowd discovery and tracking in systems like Facebook and Twitter. For practical crowd discovery and tracking in a large time-evolving communication network, however, we face four key challenges:

- First, systems like Facebook and Twitter are extremely large (on the order of 100s of millions of unique users), placing huge demands on the computational cost of traditional community detection approaches (which can be  $O(n^3)$  in the number of users [7]).
- Second, these services support a high-rate of edge addition (new messages) so the discovered crowds may become stale quickly, resulting in the need to re-identify all crowds at regular intervals (again, incurring the high cost of community detection). The bursty nature of user communication demands a crowd discovery approach that can capture these highly-temporal based clusters.
- Third, the strength of association between two users may depend on many factors (e.g., recency of communication), meaning that a crowd discovery approach based on graph clustering should carefully consider edge weights. With no decay at all (meaning that edges are only inserted into the network but never removed), all users will tend towards a single trivial large crowd. Conversely, overly aggressive edge decay may inhibit any crowd formation at all (since edges between users may be removed nearly as soon as they are added).
- Fourth, crowds may evolve at different rates, with some evolving over several minutes, while others taking several days. Since crowds are inherently ad-hoc (without unique community identifiers – e.g., Fans of LA Lakers), the formation, growth and dispersal of crowds must be carefully managed for meaningful crowd analysis.

With these challenges in mind, we propose to discover and track transient crowds through a communication based clustering approach over time-evolving graphs that captures the natural conversational nature of social messaging systems. Two of the salient features of the proposed approach are (i) an efficient locality-based clustering approach for identifying crowds of users in near real-time compared to more heavy-weight static clustering algorithms; and (ii) a novel crowd



**Figure 2: Examples of the crowds discovered in the dataset.**

tracking and evolution approach for linking crowds across time periods.

To support transient crowd discovery in Twitter-like services with 100s of millions of participants, we propose to leverage the inherent locality in social messaging systems. Concretely, we identify two types of locality that are evident in Twitter-like messaging systems: (i) temporal locality and (ii) spatial locality.

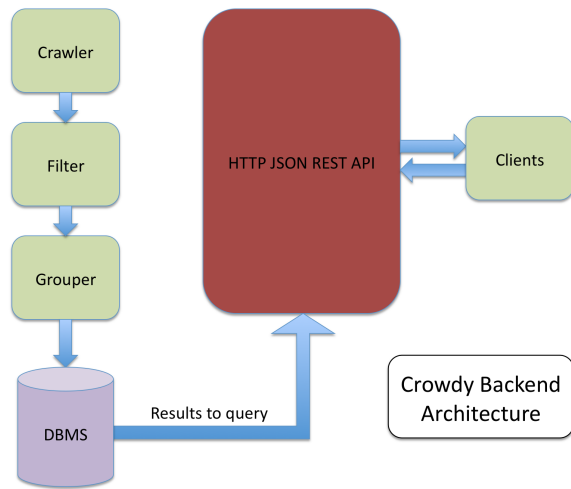
**Temporal Locality:** Transient crowds are intuitively short-lived, since they correspond to actively communicating groups of users. Hence, the composition of a crowd at a point-in-time should be impacted by recent messages as opposed to older messages. As more users interact with the crowd, the crowd should grow reflecting this *temporal locality* and then shrink as users in the crowd become inactive (that is, their last communication with the crowd becomes more distant in time).

**Spatial Locality:** Intuitively, transient crowds are made up of a very small percentage of users compared to the entire population of the social network. Hence, new messages (corresponding to the addition of edges to the communication network) should have only a local influence on the crowds that exist at any given time. That is, changes in a small region of a graph should not affect the entire graph. In a dataset of 61 million Twitter messages described in [16], we have confirmed the existence of this *spatial locality* by finding that only about 1% of users are within two hops, meaning that an edge insertion has only a local effect.

Hence, we can take advantage of both, local changes to the overall communication network (spatial locality) and recent changes to the network (temporal locality), for supporting efficient transient crowd discovery. The complete algorithm is described in [16]

### 3.1 Sample Crowds Discovered

We illustrate some of the discovered crowds and their lifespans in Figure 2, with an annotation next to the crowd peak showing the topic of discussion. We see a crowd (shown in black) discussing Sarah Palin and the Vice-Presidential debate from the 40<sup>th</sup> hour to 80<sup>th</sup> hour that peaks around



**Figure 3: Examples for trending phrases discovery problem.**

the time of the actual debate. We observe that crowds that talk about general everyday things have a greater lifespan than crowds discussing specific events. For example in Figure 2, a crowd (annotated with *thank, whats, wow*) discussing everyday things lives through the entire week, while, during the same period we observe several event-specific crowds, like crowds discussing the Red Sox, Sarah Palin, and Girl's Night Out (gno) forming and dispersing. These event-specific crowds start forming just before the event and die a few intervals after the completion of that event. This distinction between the crowds discovered clearly indicates two types of Twitter usage: first, it is used as a platform to discuss and debate specific events, and second, it is also used as a means of everyday communication.

## 4. ARCHITECTURE

An architecture for the proposed solution is shown in Figure 1. The various modules of this architecture are as follows:

**Crawler:** This module deals with interacting with social media websites to collect information. For this project we will be collecting data from Twitter using Twitter Streaming APIs. We collected tweets from [xxxx] users. To collect this set of users we started with a set of 5 users and used snowball sampling to collect more users. We then used the “follow” parameter of the filter method from Streaming API to generate a data stream for crowd detection.

**Grouper:** The grouper module analyses the data collected, to determine crowds and adds them to the DBMS. This is the central module that detects crowds of different types - communication, content, geo-graphical, etc. In this project we will be discovering crowds based on communication only. The algorithm to discover crowds is described in Section 3. Once crowds are discovered, this module formats crowds to a well defined crowds object and adds it to a DBMS.

**Filter:** This module is used to filter un-necessary data from the framework. It can appear both before and after the grouper module depending upon its functionality. In the

figure we have showed the module appear before grouper only. When the module appears before the grouper, it is used to weed out information unnecessary for the grouper. For example in case of crowd detection, it remove information like spam messages. When the module appears after the grouper, it cleans the data generated by the grouper. For, example in case of crowd detection, it analyzes crowds discovered by the grouper and removes un-necessary crowds, or crowds that don't meet specific quality guarantees. The module can also have additional functions that support features like search.

**HTTP JSON RESTful APIs:** This module provides client an interface to the digital library. The module exposes JSON APIs that the clients can use to browse and search the crowds archive.

### 4.1 Tools

The tool we used to develop this framework are as follows:

- **PyLucene:** To provide the search functionality we use Apache Lucene. It provides methods for text indexing and searching. To integrate Lucene with our framework we use PyLucene which is a Python wrapper for Lucene. This is used in the filter module of our architecture.
- **CherryPy:** This is an object-oriented web framework for Python. The JSON REST API that interface the digital library with clients is hosted using this web framework.
- **NetworkX:** This is a library for manipulating graphs in Python. We use this to determine quality of crowds generated by crowd detection framework and filter out crowds of poor quality. This sits in the filter module.
- **Beanstalk:** Beanstalk is a simple, fast work queue. We use it to connect the crawlers, filters, and groupers together.
- **MongoDB:** MongoDB is a scalable, high-performance, open source, document-oriented database. We store the crowds discovered by the framework in this DB. The API module interacts with this database to retrieve collections as required by the user query.

## 5. CONCLUSION

## 6. REFERENCES

- [1] Facebook wheel.  
<http://apps.facebook.com/friendwheel/>.
- [2] Real-time local twitter trends.  
<http://trendsmap.com/>.
- [3] J. Bishop. Increasing participation in online communities: A framework for human-computer interaction. *Comput. Hum. Behav.*, 23:2007, 2007.
- [4] D. M. Boyd and N. B. Ellison. Social Network Sites: Definition, History, and Scholarship. *Journal of Computer-Mediated Communication*, 13(1):210–230, Oct. 2008.
- [5] M. Cheong and V. Lee. Integrating web-based intelligence retrieval and decision-making from the twitter trends knowledge base. In *Proceedings of the*

- 2nd ACM workshop on Social web search and mining, SWSM '09*. ACM, 2009.
- [6] A. Dieberger, E. Kandogan, and C. Kieliszewski. Scalability in system management GUIs: a designer's nightmare. In *CHI'06 extended abstracts on Human factors in computing systems*, pages 99–104. ACM, 2006.
  - [7] G. W. Flake, R. E. Tarjan, and K. Tsioutsoulouklis. Graph clustering and minimum cut trees. *Internet Mathematics*, 1(4):385–408, 2004.
  - [8] J. Heer and D. Boyd. Vizster: Visualizing Online Social Networks. *Information Visualization, IEEE Symposium on*, 0, 2005.
  - [9] S. Henderson, B. Biller, M. Hsieh, J. Shortle, J. Tew, R. Barton, Z. Mikovec, I. Maly, and P. Slavik. Visualization of user's activities in a specific environment.
  - [10] T. Ho and T. Nguyen. Visualization support for user-centered model selection in knowledge discovery in databases. In *ictai*, page 228. Published by the IEEE Computer Society, 2001.
  - [11] C.-F. Hsu, E. Khabiri, and J. Caverlee. Ranking comments on the social web. In *CSE '09: Proceedings of the 2009 International Conference on Computational Science and Engineering*, pages 90–97, 2009.
  - [12] A. Java and T. Finin. Mining social media communities and content. In *Doctoral Dissertation in University of Maryland at Baltimore County Catonsville*. ACM, 2008.
  - [13] A. Java, X. Song, T. Finin, and B. Tseng. Why we twitter: understanding microblogging usage and communities. In *Proceedings of the 9th WebKDD and 1st SNA-KDD 2007 workshop on Web mining and social network analysis*, WebKDD/SNA-KDD '07. ACM, 2007.
  - [14] K. Kamath and J. Caverlee. Identifying hotspots on the real-time web. In *Proceedings of the 19th ACM international conference on Information and knowledge management*, CIKM '10. ACM, 2010.
  - [15] K. Kamath and J. Caverlee. Transient crowd discovery on the real-time social web. In *Proceedings of the fourth ACM international conference on Web search and data mining*, WSDM '11. ACM, 2010.
  - [16] K. Y. Kamath and J. Caverlee. Transient crowd discovery on the real-time social web. In *Proceedings of the fourth ACM international conference on Web search and data mining*, WSDM '11, pages 585–594, New York, NY, USA, 2011. ACM.
  - [17] E. O. Laumann and L. Guttman. The relative associational contiguity of occupations in an urban setting. *American Sociological Review*, 31:169–178, 1966.
  - [18] J. J. Leggett and F. M. Shipman, III. Directions for hypertext research: exploring the design space for interactive scholarly communication. In *HYPERTEXT '04: Proceedings of the fifteenth ACM conference on Hypertext and hypermedia*, pages 2–11, New York, NY, USA, 2004. ACM.
  - [19] J. L. Moreno. *Who Shall Survive?* Beacon House Inc, 1953.
  - [20] E. Murphy-Hill and A. Black. An interactive ambient visualization for code smells. In *Proceedings of the 5th international symposium on Software visualization*, pages 5–14. ACM, 2010.
  - [21] T. O'Reilly and S. Milstein. *The Twitter book*. O'Reilly, 2009.
  - [22] D. Phan, A. Paepcke, and T. Winograd. Progressive multiples for communication-minded visualization. In *Proceedings of Graphics Interface 2007*, pages 225–232. ACM, 2007.
  - [23] E. Schlungbaum. Individual user interfaces and model-based user interface software tools. In *Proceedings of the 2nd international conference on Intelligent user interfaces*, pages 229–232. ACM, 1997.
  - [24] V. Tran. UI generation from task, domain and user models: the DB-USE approach. In *Proceedings of the 2nd ACM SIGCHI symposium on Engineering interactive computing systems*, pages 353–356. ACM, 2010.