

Near Real Time Assessment of Social Media Using Geo-Temporal Network Analytics

Kathleen M. Carley^{ab}
and Jürgen Pfeffer^a

^aISR, SCS, Carnegie Mellon University

^bNetanomics

Pittsburgh, PA, 15213

Email: {kathleen.carley, jpfeffer}@cs.cmu.edu

Huan Liu

and Fred Morstatter

CIDSE

Arizona State University

Tempe, AZ, 85281

Email: {huan.liu, fred.morstatter}@asu.edu

Rebecca Goolsby

Office of Naval Research

Arlington, VA, USA

Email: rebecca.goolsby@navy.mil

Abstract—When a crisis occurs, there is often little time to evaluate the situation and determine how best to respond. We use rapid ethnographic methods centered on the construction of geo-temporally contextualized social and knowledge networks. By utilizing a combination of Twitter and news media, the consulate attack in Libya were examined in near real time. In this work we outline a procedure to extract key insights from the event as an event unfolds using a suite of tools developed by a team of researchers from two universities.

I. INTRODUCTION

As a crisis occurs, there is often little time to evaluate the situation and determine how best to respond. An example of such a crisis is the 2012 Benghazi consulate attack in Libya. How can the analyst or policy maker get early insight into a crisis as it unfolds? What information is available? How can that information be tracked? Finally, are there any early indicators or warning signs of these crises?

We ask, can these questions be addressed using a combination of traditional and social media? This paper addresses these questions by describing a near real time assessment activity that was occurring as the attack began and continued for 72 hours after the event. The data was collected in a few hours and the analysis done immediately. This process was repeated multiple times during this roughly 96 hour period. Herein we describe this process and illustrate the type of analyses done and visualizations constructed using the final images from roughly 72 hours after the event. The setting was at EUCOM, where the ASU-CMU research team was running a training session on social media exploitation under the auspices of the ONR. During training, the Libyan consulate was attacked. As a class exercise the team demonstrated how that event could be analyzed with the tools being taught. The analysts had received approximately 3 hours of training on TweetTracker and 6 hours on ORA (aka ORA-NetScenes), before they began producing results. This paper describes the process and results of this exercise. All images and data herein are based on the data collected and analyzed by the ASU-CMU team during those few days, most during the first 36 hours. A similar activity was conducted vis Hurricane Sandy and the 2013 Kenyan elections. Some of those results are reported herein.¹

¹Additional results can be seen at www.pfeffer.at/sandy and www.casos.cs.cmu.edu/projects/kenya

An ability to monitor social media and news data and use such data to rapidly characterize the socio-cultural landscape, i.e., the cultural geography, is critical in crises [1], and for the provision of humanitarian assistance and disaster response [2]. Carnegie Mellon University (CMU), Netanomics, and Arizona State University (ASU) have created a set of interoperable technologies that support the collection, analysis and visualization of on-line data – both social media and traditional media. A key feature of these tools is that they admit rapid ethnographic analysis of situations through the extraction of geo-temporal multi-dimensional networks often referred to as meta-networks [3]. The resulting process admits rapid assessment in near real time and preserves the processed data for more detailed exploration that can be conducted at leisure by the analyst.

II. TOOLS

There are four basic tools that are used in an interoperable fashion. See Figure 1 for a high level overview. These tools are TweetTracker, Tweet-to-ORA, REA, and ORA. TweetTracker [4] pulls tweets from the Twitter API in response to the filters provided by the analyst. Tweet-to-ORA converts the tweets extracted into a format that is importable by ORA. REA pulls news articles and associated tags from LexisNexis in response to the filters provided by the analyst and also converts them into a format that is importable by ORA. ORA [5], [6] is a dynamic social network analysis tool that allows the analyst to analyze and visualize semantic networks, social networks and other geo-temporal high dimensionality networks. ORA supports the analysis and visualization of tweets, e.g., by processing the hashtag and retweet network, and news article, e.g., by processing the social, knowledge, and task networks described therein.

A. TweetTracker

TweetTracker is a tool developed at ASU that allows analysts to collect and analyze tweets in real-time [4]. Analysts of TweetTracker specify the data they wish to collect in the form of parameters specific to the event they are interested in studying. The analyst specifies three different kinds of parameters: keywords, geographical boundary boxes, and tweeters. This is consistent with the way tweeters publish data on Twitter [7]. When tweeters publish tweets, they write a message of 140 characters or less. They also have the option to “geo-tag”

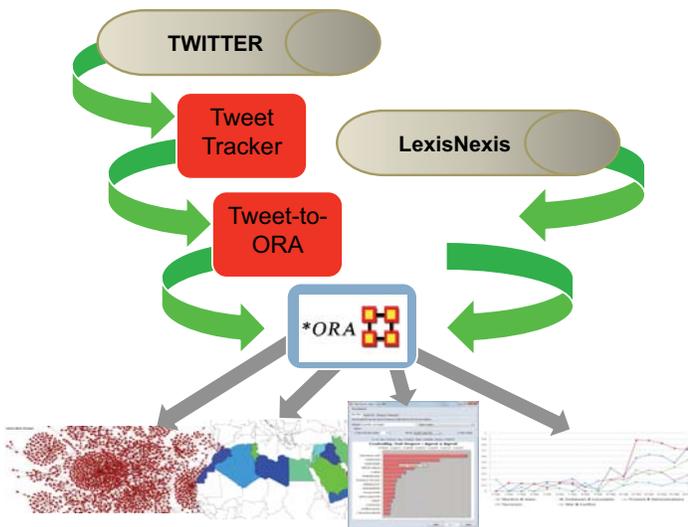


Fig. 1. High-Level View of System Interoperability, where color differentiates tool source – Red ASU, Blue CMU-Netanomics.



Fig. 2. Main window for TweetTracker.

their tweets. By geo-tagging, the tweeter shares with the world where the tweet was published. This is accomplished through the location sensor on the device (e.g., GPS on a mobile phone, IP on a web browser, etc.). Finally, TweetTracker allows the analyst to collect full timelines of Twitter users. TweetTracker has been used to collect Twitter data from Arab Spring protests, Occupy Wall Street, and many recent natural disasters. Figure 2 shows the main window for TweetTracker.

TweetTracker does not pull down the entire Twitter data set. Adhering to the limits set forth by Twitter, TweetTracker only extracts at most 1% of the tweets available and only those tweets that match the filters provided by the analyst. Imagine that there are 400 million tweets per day. At most 4 million will be extracted. If the filters provided generate more than 4 million tweets in a day, the set of tweets delivered will be arbitrarily capped by Twitter to 4 million tweets. Additionally, sometimes Twitter simply blocks data collection. Approximately 1% of all the tweets in Twitter have geo-tags and the same is true of the tweets collected. The tweets collected are a representative sample and sometimes a full collection of the tweets for those filters depending on the specificity of the filters [8]. TweetTracker tracks the tweets and retweets; however, it does not track the follower network.

Filters are words or phrases, geographic bounding boxes,

and tweeters. Words can be, but need not be, hashtags or user IDs. All filter parameters provided by the analyst are combined using an “or” function which casts a wide net and tries to extract as much data as possible from Twitter. The more specific the set of filters, the more likely the entire corpus of tweets related to those filters will be extracted.

B. Tweet-to-ORA

Tweet-to-ORA is a tool developed by collaboration between ASU and CMU which enables the analyst to export the information from TweetTracker into ORA. It extracts the timestamp, user ids, hashtags, and geo-location data from each tweet and puts them into a format that ORA can ingest. ORA imports this data, forming a dynamic meta-network in which there are a set of meta-networks by time period. In each meta-network, there are sub-networks: tweeter-to-tweeter retweet network, tweeter-to-location geographic network, tweeter-to-hashtag network, hashtag-to-hashtag co-occurrence network, and hashtag-to-location geographic network.

C. REA

Rapid Ethnographic Analyzer (REA) [9] is a process model developed at CMU that allows the analyst to extract news data from LexisNexis for use in ORA. REA does for LexisNexis news articles what the combination of TweetTracker and Tweet-to-ORA does for tweets. This system operates as a script in background. Using filters provided by the analyst, it downloads all articles and their tags in the specified time range that are available in LexisNexis. It then takes that data and creates a file for import into ORA that contains the following classes of nodes: agents (the people discussed), organizations (which are sub-categorized into specific organizations, industries, and other institutions), locations, and knowledge (these are the topics discussed). All networks connecting any of these classes with another class or itself are then constructed. The tie values are the counts of the number of times the tags co-occurred in the same article. The articles are also extracted and can be processed for more detailed information by text mining tools that produce networks such as AutoMap [10].

D. ORA

ORA, see Figure 3, is a tool developed by CMU and Netanomics, that allows analysts to fuse, analyze, visualize, and forecast the behavior of network data [5], [6]. Using ORA the analyst can identify key actors, key topics, key locations, characterize and visualize networks, assess changes in the networks and key locations in terms of where they are by using the geo-spatial mapping functions, and multiple other tasks. The system is organized to help create products about who, what and where are important when. The algorithms in ORA are from the fields of social network analysis [11], dynamic network analysis, link analysis and network science [6]. ORA employs both graph analytic and statistical network algorithms to assess, visualize and forecast behavior for geo-temporal networks. In addition, ORA supports 2D and 3D network visualization, geo-spatial network visualization, and traces of network activity across time and location.

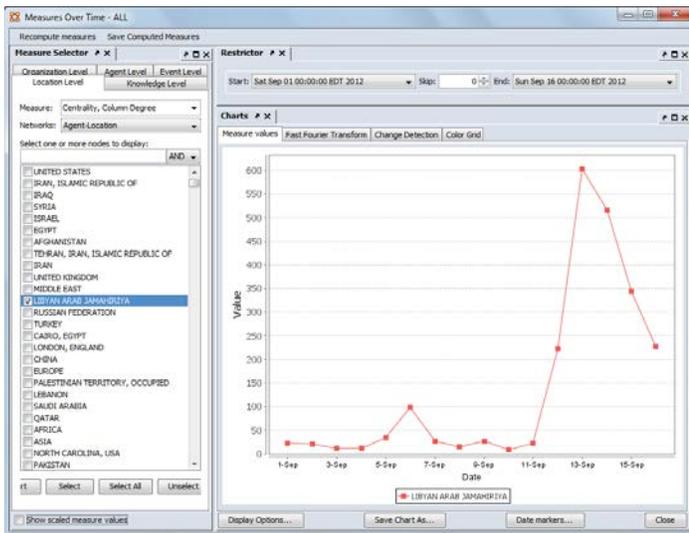


Fig. 3. Interface for temporal analysis in ORA.

TABLE I. DISTINGUISHING FEATURES OF DATA

Feature	Tweets	Tagged News
Size	140 Characters	1000+ Characters
Timing	As Generated	Lock-stepped Daily
Producer	Individuals & Corporations	Corporations
Edited	No	Yes
Tags	Hashtags by Author	Keytags Auto-Created
Source Language	Multiple Languages	English
Available for Collection	A Few Weeks	In Perpetuity
Items Collected	Millions	Hundreds of thousands

III. DATA

To assess a potential or actual crisis situation, two types of data are collected. First tweets. Illustrative tweets related to the embassy attack are shown in Figure 4. Second, news articles and the auto-tags for them created by LexisNexis were collected. Key differences between these types of data are shown in Table I.

A. Social Media: How Tweet Data was Collected

Beginning February 2nd, 2011 ASU began collecting data on Arab Spring activity in Libya using TweepTracker. We selected parameters that were expected to yield data relevant to the massive protest activity in the region. These keywords are: #libya, #gaddafi, #benghazi, #brega, #misrata, #nalut, #nafusa, #rhaibat, and ليبيا (“Libya”, in Arabic). ASU drew a geographic boundary box with the Southwest latitude/longitude point at (10.0, 23.4) and the Northeast point at (25.0, 33.0). Since the beginning of the collection through the time of this writing TweepTracker has collected over 5 million tweets pertaining to the activity in Libya. This data serves as a baseline. During the exercise at EUCOM students collected additional tweet data focused on the embassy attack.

B. Online News: How LexisNexis Data was Collected

CMU used REA to collect data on all 18 countries associated with the Arab Spring - see Figure 5. Starting from July 2010, approximately 600,000 news articles have been collected. This data serves as a baseline. During the exercise at EUCOM new data was collected using REA. The time period

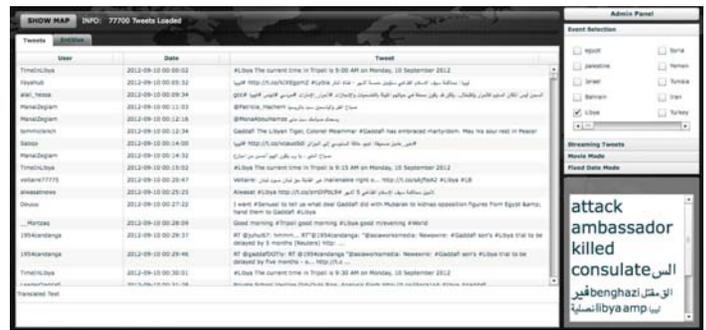


Fig. 4. Illustrative tweets as seen in TweetTracker. Tweets are in the left panel, sorting function in top right, and key concepts as a word cloud on lower right.



Fig. 5. Countries of interest for REA.

of interest is September 1-16, 2012. We collected 11,279 articles from 700+ major world publications in the LexisNexis database that discuss 18 Northern African and Middle East countries. All these newspaper and magazines are written in English. From these articles we extracted 192,913 index items that are grouped into the following categories: people, topics, organizations (including companies), and locations. LexisNexis is a professional provider of online information and offers access to articles of thousands of newspapers and news agencies worldwide. The LexisNexis SmartIndexing “applies controlled vocabulary terms for several different taxonomies”². For every article a couple of items are automatically indexed describing the content of the article (e.g. one article might be tagged with “Muammar Gaddafi”, “military operations”, and “human rights”). The items are standardized to avoid different items with identical meaning, e.g. Libya is named by its official name Libyan Arab Jamahiriya. We extract these index items and create networks based on co-occurrence of people, organizations, and locations in the same articles.[9]

IV. PROCEDURE

To study events like the Libyan embassy attack the analyst needs two types of data: a) baseline information and b) specific event information [12]. Baseline data can be continuously collected in background on general topics of interest. This data provides a background against which the event specific information can be calibrated. TweepTracker and REA are used to collect the background data and the specific event data. In both cases a “filter” needs to be created; i.e., a list of keywords that will be used to select the tweets and news items of interest. In general, this list should include the name of key political actors, or country of interest as well as general types

²<http://wiki.lexisnexis.com/academic/index.php?title=SmartIndexing>

of events of interest such as protest. Specific hashtags can be used as well. Keywords should be relatively specific phrases, rather than general words. For example, if interested in human trafficking, terms such as “human trafficking” and “sexual exploitation” will provide better results and less noise than “sex”. The second type of data is specific event information, crisis data. TweetTracker and REA are used to collect a second set of data during the crisis but using a more refined and crisis specific set of filters. The resulting set of data is within the realm of the baseline but narrower in scope. Data is collected continuously and can be analyzed by porting to ORA on demand. The ORA analyses and visualization take a few seconds to a few minutes depending on the size of the data.

Then the collected data is visualized to see general trends and to gauge the pattern and level of activity. Summary statistics may be generated such as the volume of tweets and articles relative to a specific search term. Simple visualizations and initial exploration of the tweets can be done in TweetTracker. The visualization and these summary statistics provide the analyst with a simple characterization of the data their filters have retrieved. After the filter has amassed a volume of tweets, the analyst runs Tweet-to-ORA. Next, the analyst imports the file produced by Tweet-to-ORA and that news data from REA into ORA. Then the analyst should visually inspect the data to identify any odd anomalies. Sometimes the keywords used in collecting the data need to be adjusted. As the analyst gets to know the data, obvious issues such as removal of irrelevant information can be dealt with. For example, ORA makes it easy to remove all data associated with actors or locations not of interest, anonymize the data, or merge data points together. This latter feature is important as many keywords and hashtags refer to the same thing.

ORA is then used for a more detailed evaluation; e.g., identifying key actors in the Twitter network with more than normal influence and identifying topics that are gaining in importance. The analyst can choose to use a narrow temporal window, e.g., an hour or a day, or a larger window, such as several days or a month. ORA forms a network within this window and supports dynamic analysis of changes across time and space. The analyst uses this network analytic capability to explore items of interest. If a specific tweeter or hashtag appears critical, the analyst can then go back to TweetTracker and explore the specific tweets associated with that tweeter or hashtag. Or, similarly, for news articles one can return to the URL for the news item and examine it.

V. RESULTS

On September 11th, 2012, the United States ambassador to Libya was killed in an attack on the U.S. consulate [13]. On September 12th, discussions of this attack exploded on social media. Using TweetTracker’s already-running Libya stream, we were able to capture tweets pertinent to this event. Since September 11th, the analyst has collected 114,515 tweets, with September 12th containing the largest spike in months of data. The 70,630 tweets collected on September 12th alone account for over 23% of all the tweets collected since May 1st. Figures 6 and 7, show the difference between all Libya tweets and just those involving Libyan Embassy. In Figure 6 we see that there are few tweets about Libya until the attack on the embassy. Figure 7 shows a definite temporal pattern to

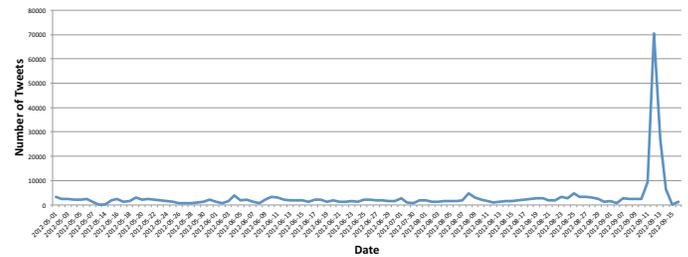


Fig. 6. Tweets per day mentioning Libya as displayed in TweetTracker.

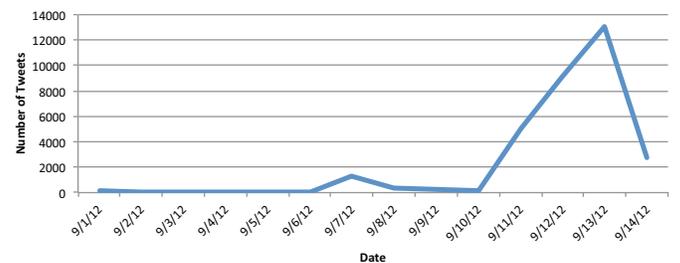


Fig. 7. Tweets per hour mentioning embassy as displayed in TweetTracker.

the tweets. Such patterns can be analyzed in ORA with Fourier analysis and over time trending algorithms. [14] These spikes are an alert that “something” is happening.

Next the analyst examines the news articles. Figure 8 shows the news articles associated with Libya. The sheer volume, i.e. the peaks, indicates activity in the region. There is little discussion of Libya until the embassy is attacked. In general, tweet data will lead news data just in volume by about a day, partially due to publishing deadlines [15].

The analyst next explores whether there was a geographical spread with respect to embassy attacks. Figure 9 shows related tweets segregated by country by hour in Arizona time. Notice that Libya is basically dormant until September 11, 2012 and

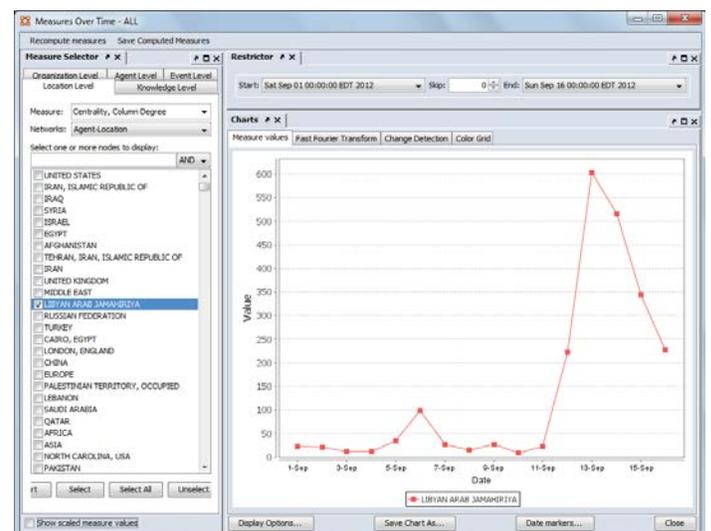


Fig. 8. News articles per day mentioning Libya as displayed in ORA.

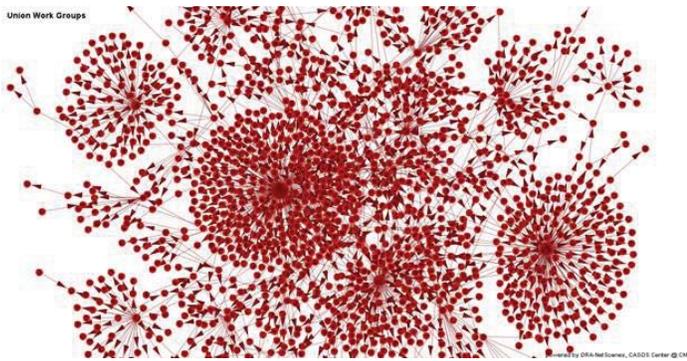


Fig. 13. Retweet network for Libya data as displayed in ORA.

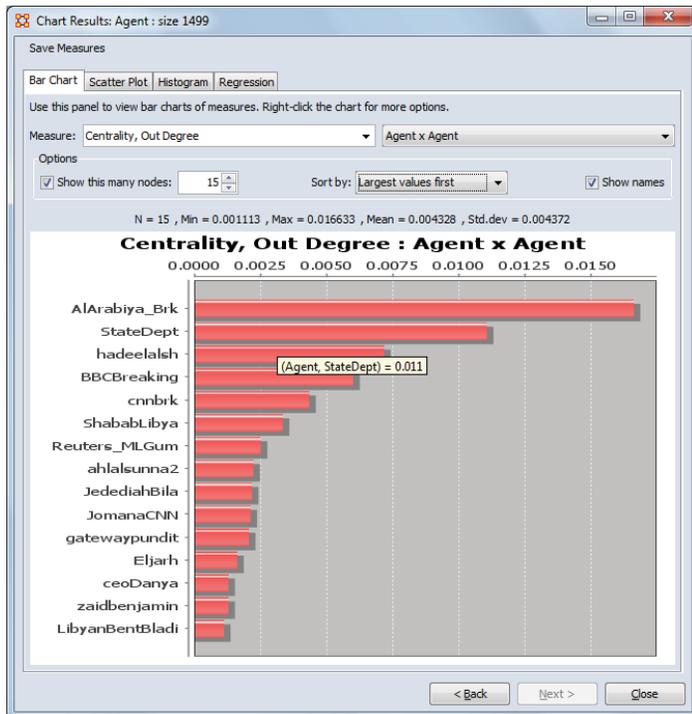


Fig. 14. Most respected tweeters as identifies using network analytics in the ORA Twitter report.

see that of the top six tweeters, those who are most frequently being retweeted, four are news agencies. Thus the tweets being most frequently spread are those by organizations not “the person on the street”. The other two most frequent are Hadeel Al-Shalchi, @hadeelalsh. A Middle East Correspondent for Reuters and the LibyanYouth movement, ShababLibya. On Sept 12, 2012 the most retweeted tweeter concerning Libya was AlArabiya_Brk with 636 retweets and then BorowitzReport with 632 retweets.

Now the analyst switches and examines what is being talked about. Figure 15 shows the core of the hashtag network. In this case there is a link just in case two hashtags appeared together in more than 20 tweets. This network breaks into two components - an Arabic and an English component. The Arabic hashtags only co-occur with Arabic hashtags and same for the English hashtags. This means in this data, the tweets are mono-language. Those hashtags that are connected

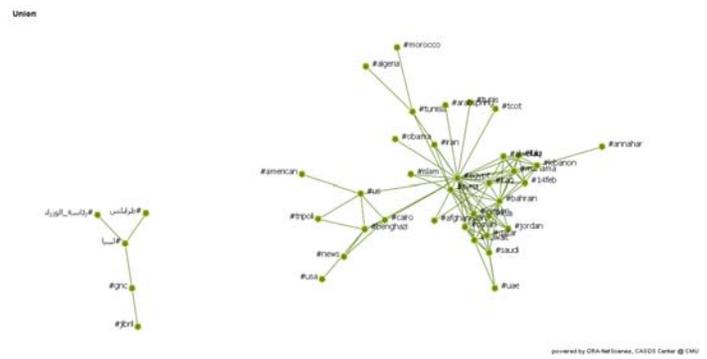


Fig. 15. Hashtag network for Libya data as displayed in ORA.

TABLE II. TOP HASHTAGS IN FIRST 24 HOURS IN ALL LIBYA TWEETS.

Hashtag	Number of Occurrences
#benghazi	644
#egypt	512
#secclinton	222
#gnc	193
#usa	190
#us	188
#ليبيا	168
#syria	159
#cairo	94
#tripoli	67

to large numbers of other hashtags (most central hashtags) are important in that they signal a central focus of concern. Notice that some of the most important hashtags are the Arabic word for Libya #ليبيا, #egypt and #syria. It is worth noting that the hashtag #benghazi is often linked to #cairo, #egypt, #us, #usa, #news, #tripoli. Suggesting that parallels are being drawn between this event and other revolutionary activities.

The top 5 hashtags during these first 24 hours are shown in Table II. The Arabic hashtag is the Arabic word for Libya. Note that the system currently does not clean the data so there are multiple hashtags identified that actually refer to the same topic - such as USA and US. The analyst can use ORA to merge these into a single node if desired. These hashtags, like the top hashtags in the Arab Spring are predominantly the names of cities or actors of import. The forgoing analysis takes about 1 hour to accomplish. In crisis events it is then repeated on demand, e.g., every six hours. The data is saved and automatically kept with the next increment of data. This supports more detailed followup analyses.

Analysts given this wealth of information can then follow up by addressing other questions such as:

- Is different information coming from the LibyanYouth movement than the news agencies?
- Which tweeter among these key actors are the “canaries” providing earliest information?
- When did discussion about the movie “Innocence of the Muslims” start and why?

An example of a follow up question is “What role did the movie Innocence of the Muslims play?” In Figure 16 the number of tweets mentioning the movie are shown In the tweets associated with Libya, while a few mentions did occur on the day of the attack, temporally most of those occurred

ters and other incidents, e.g. hurricane Sandy, the flooding in Thailand, the Kenya elections. On March 4, 2013 a presidential election was held in Kenya. There were numerous incidents of tribal violence since the last elections in 2007 and the analyst's questions for the weeks before the 2013 elections were: Is violence increasing as election time approaches and what is triggering it? What are the topics and events? Who is discussing them and what are they saying? Are the events similar than in the past? What can we expect for the weeks before and after the elections? News articles and Twitter data were analyzed similar as described in the previous sections. Figure 17 shows all geo-tagged tweets in the time period February 1–5, 2013 that discuss “Kenya”. As one can see, tweeters are located all over the globe. To get a better impression about what is discussed inside and outside of the country, the locations of the tweeters serve as filter for further analyses. These reveal that violence was not a topic discussed in Kenya four week before the elections (Figure 18).

VIII. A LOOK TOWARDS THE FUTURE

The data presented in this paper should not be interpreted as providing guidance on what happened during or in the immediate aftermath of the consulate attack or other incidents. Rather, it should be viewed as showing the strengths and limitations of this type of data. We present it more as guidance for what is possible and what can be done; and not as an assessment of the events. It is important to note that this tool-suite as is can support the analyst. Critical limitations were identified. For each of these, work is underway at various levels to meet the unfilled need.

The key limitations identified, in terms of immediate needs, are as follows. First, many of the tweets are from news broadcasting corporations; thus, it is difficult to disambiguate public sentiment from news-reporting bias. Future work needs to separate the two sources of tweets. Second, geo-spatial identification is poor. Most tweets are not geo-tagged. Basic research is needed to develop algorithms for inferring location when possible from non-geo-tagged tweets. Further, the current technologies need to be extended to differentiate tweets originating within and without the region of interest for analysis purposes. Third, either automated translation or language independent clustering of results and generation of filters is needed. Fourth, automated or semi-automated approaches for mapping the filters used for news and Twitter to common terms and for mapping the results to common terms is needed to support comparative analytics and data fusing. Basic research on cross-data source analytics is needed. Fifth, semi-automated support for creating filters is needed. We found that one of the most difficult tasks for analysts was identifying terms of interests and creating good filter lists. Finally, the entire systems needs to be increased in scale particularly the map generation functions. We note that the overall system is relatively fast, however the slowest part is generation of maps which is currently a little too slow for hourly updates.

TweetTracker's forthcoming ability to track company-produced hashtags, particularly news broadcaster links, and links to objects will support more in-depth analysis. Tweet-to-ORA will be integrated into TweetTracker rather a separate tool. REA will be incorporated into ORA. ORA will have one-step importing in the wizard for Tweet data from TweetTracker.

ORA's forthcoming alert function will allow users of this tool suite to identify which parts of the tweet or news stream to look at in greater depth. Finally, ORA will have a new reporting function overview specialized to tagged data from Tweets and LexisNexis. This higher level of functionality and the easier 1 step interoperability will make it easier for analysts to engage in these types of time critical assessments.

These and other features will enhance the analyst's ability to engage in these types of time critical analytics. The key is that these analyses began within 24 hours and supported continual updated assessments during the next 72 hours using existing tools and supported critical information assessment needs. As we move to the future, such tool suites will be critical in exploiting open source information so as to respond rapidly and effectively to crisis situations and disasters.

REFERENCES

- [1] D. G. Campbell, *Egypt Unshackled: Using Social Media to @#:) the System*. Amherst, NY: Cambria Books, 2011.
- [2] R. Goolsby, “Lifting elephants: Twitter and blogging in global perspective,” in *Social computing and behavioral modeling*. Springer, 2009, pp. 1–6.
- [3] K. M. Carley, M. W. Bigrigg, and B. Diallo, “Data-to-model: a mixed initiative approach for rapid ethnographic assessment,” *Computational and Mathematical Organization Theory*, vol. 18, no. 3, pp. 300–327, 2012. [Online]. Available: <http://dx.doi.org/10.1007/s10588-012-9125-y>
- [4] S. Kumar, G. Barbier, M. A. Abbasi, and H. Liu, “Tweetracker: An analysis tool for humanitarian and disaster relief,” in *Fifth International AAAI Conference on Weblogs and Social Media, ICWSM*, 2011.
- [5] K. M. Carley, J. Reminga, J. Storrick, and D. Columbus, “ORA User's Guide 2013,” Carnegie Mellon University, School of Computer Science, Institute for Software Research, Pittsburgh, PA, Technical Report CMU-ISR-13-108, 2013.
- [6] K. M. Carley and J. Pfeffer, “Dynamic network analysis (dna) and ora.”
- [7] L.S. (2011, September) What's in a tweet. Online. The Economist. [Online]. Available: <http://www.economist.com/node/21531066>
- [8] F. Morstatter, J. Pfeffer, H. Liu, and K. M. Carley, “Is the sample good enough? comparing data from twitter's streaming api with twitter's firehose,” in *International AAAI Conference on Weblogs and Social Media (ICWSM)*, 2013.
- [9] J. Pfeffer and K. M. Carley, “Rapid modeling and analyzing networks extracted from pre-structured news articles,” *Computational and Mathematical Organization Theory*, vol. 18, no. 3, pp. 280–299, 2012.
- [10] K. M. Carley, D. Columbus, and P. Landwehr, “AutoMap User's Guide 2013,” Carnegie Mellon University, School of Computer Science, Institute for Software Research, Technical Report CMU-ISR-13-105, 2013.
- [11] S. Wasserman and K. Faust, *Social Network Analysis: Methods and Applications*. Cambridge, MA: Cambridge University Press, 1994.
- [12] R. Goolsby, “Social media as crisis platform: The future of community maps/crisis maps,” *ACM Transactions on Intelligent Systems and Technology (TIST)*, vol. 1, no. 1, p. 7, 2010.
- [13] D. D. Kirkpatrick. (2012, September) Libya attack brings challenges for u.s. Online. [Online]. Available: <http://www.nytimes.com/2012/09/13/world/middleeast/us-envoy-to-libya-is-reported-killed.html?pagewanted=all>
- [14] I. McCulloh and K. M. Carley, “Detecting change in longitudinal social networks,” *Journal of Social Structure*, vol. 12, no. 3, pp. 1–37, 2011.
- [15] J. Pfeffer and K. M. Carley, “Social networks, social media, social change,” *Advances in Design for Cross-Cultural Activities Part II*, vol. 13, pp. 273–282, 2012.
- [16] L. C. Freeman, “Centrality in Social Networks: Conceptual clarification,” *Social Networks*, vol. 1, no. 3, pp. 215–239, 1979.