# Real-time Crisis Mapping of Natural Disasters using Social Media

**Stuart E. Middleton, Lee Middleton and Stefano Modafferi,** *University of Southampton IT Innovation Centre*

## ABSTRACT

We present a social media crisis mapping platform for natural disasters. We take locations from gazetteer, street map and volunteered geographic information (VGI) sources for areas at risk of disaster and match them to geo-parsed real-time tweet data streams. We use statistical analysis to generate real-time crisis maps. Geo-parsing results are benchmarked against existing published work and evaluated across multi-lingual datasets. We report two case studies comparing 5-day tweet crisis maps to official post-event impact assessment from the US National Geospatial Agency (NGA) compiled from verified satellite and aerial imagery sources.

## Author Keywords

Geoparsing; Real-time; Crisis Mapping; Natural Language Processing; NLP; Volunteered Geographic Information; VGI; Natural Disaster; Crisis Management

## ACM Classification Keywords

H.2.8 [Database Management]: Database applications–Text mining; I.2.7 Artificial Intelligence: Natural Language Processing - Text analysis

## General Terms

Algorithms; Design; Experimentation

## INTRODUCTION

In today's society the ubiquitous use of mobile communication devices has seen social media sites, such as Facebook, Twitter and YouTube publishing microblogs, images and videos in real-time from people experiencing natural disaster events often live and in-situ. In the humanitarian sector this has sparked great interest [1] in developing innovative approaches to utilize social media for events such as earthquakes, floods and tornados to both inform the public and assist civil protection authorities in focussing response efforts.

Recent natural disaster events have seen humanitarian organizations and networks of volunteer's setup live web-based manual crisis mapping sites [1] such as for the Haiti 2010 earthquake, Russian 2010 wild-fires, New York's 2012 hurricane Sandy and Oklahoma's 2013 tornado. These organizations check and filter crowd-sourced information from news reports, social media and civil protection agency alerts, and present it live on web-based crisis maps for the general public to see. Challenges [1] for these organizations include automating the huge task of real-time data fusion of large volumes of multi-source heterogeneous information and maintaining the trust & credibility in this data.

In this article we present a real-time crisis mapping platform capable of geoparsing Tweet content. Our novel approach exploits readily available location information from gazetteers, street maps and volunteered geographic information sources. Our goal is to improve geoparsing precision of street-level tweet incident reports and empirically quantify how accurate resulting social media crisis maps can be during natural disaster events. To our knowledge this is the first time an analysis has been published which directly compares street-level Twitter-based crisis maps to a verified ground truth based on post-event expert assessment. Such results are needed to help disaster management agencies assess the value of social media crisis mapping.

Currently real-time Geospatial Information Systems (GIS) [2] mostly map social media microblog reports using geotag metadata with long/lat coordinates. This approach turns social media into a crowd-sourcing virtual sensor network, allowing maps of twitter messages to be plotted. According to the US Geological Survey (USGS) [3] the main benefit of Twitter-based detection systems over sensor-based systems is their fast detection speed and low cost. Social media GIS systems can be combined with conventional GIS systems deploying hardware-based sensors, such as in-situ seismic sensors or remote sensing aerial photography & satellite imaging. Overall the aim is to build up a coherent situation assessment picture [3] [4] [5] and present it to emergency responders, civil protection authorities and the general public to help coordinate response efforts and improve overall awareness.

Unfortunately only about 1% of all tweets actually contain geotag metadata, see figure 1, and of this 1% the geotags are a mixture of genuine mobile devices (using GPS) and Twitter's default of the user's home location. In addition the tweeted text can contain references to one or more locations geospatially distant to the location of the device sending the tweet; this does not matter when mapping course-grained earthquake regions but does matter for finer grained maps such as flood inundation or tornado damage. Figure 1 shows the tweet breakdown during 48 hours of the 2012 hurricane Sandy which we recorded using our Twitter crawler. We have observed from our crawled tweet datasets that during events people tweet about flooding/damage to specific buildings, roads and geographic features such as local parks, rivers and beaches. Tweet reports are a mixture of a few first-hand reports and many re-tweets and comments on third party incident reports.

Geo-parsing systems [6] [7] [8] can parse text documents to extract likely geographic tokens or 'named entities' (e.g. places or regions such as 'New York'). When coupled with a geocoder, which can lookup location names on a map and return their geotag, this provides a way to associate geotags

for locations mentioned in microblog reports. Such systems often use a technique called named entity recognition. First the text is tokenized to extract sentences and words. Each token (i.e. word) is classified using a language-specific parts of speech (POS) classifier, identifying a lexical category (e.g. 'ADJ' adjective, 'N' noun, 'NP' proper noun). Lexical patterns can then be used to identify groups of tokens that are likely to refer to named entities. Challenges [6] for named entity recognition include acquiring enough labelled training data, handling poorly structured text from sources like Twitter and multi-language scalability.

### REAL-TIME CRISIS MAPPING PLATFORM

We are interested in mapping real-time tweet flood reports for 'at risk' coastal areas near known geological fault lines which have the potential to cause a Tsunami. Real-time monitoring is important as early wave impact assessments can be used to warn people on coastline further away allowing them to get to safety. Another key issue for decision makers in early warning control centres [9] is keeping crisis map false alarm rates to a minimum, since this undermines credibility in the data source.

To evaluate how accurate geo-parsing of locations from Twitter data can be we compare location matches from our platform against expert manual labels for tweet datasets covering disaster zones located in the US, New Zealand, Italy and Turkey. To evaluate how social media crisis maps compare to expert impact assessments we directly compare both our flood tweet crisis map, for hurricane Sandy 2012, and our tornado tweet crisis map, for the Oklahoma tornado 2013, to official US National Geospatial Agency (NGA) post-disaster impact assessment maps compiled from verified satellite and aerial imagery.

Our system differs from existing crisis mapping approaches in that we geo-parse tweet text in real-time rather than only using the tweet's geotag. This means we can access all crawled tweets rather than the 1% with a geotag. We avoid the need for language and location specific training sets by pre-loading available gazetteer, street map and volunteered geographic information (VGI) data for areas 'at risk' of disaster. This allows us to work at a building and street level resolution as opposed to only working with higher level administrative regions. Finally we make use of statistical analysis techniques to identify a 'baseline noise signal' and use this to reduce false positives in our crisis maps.
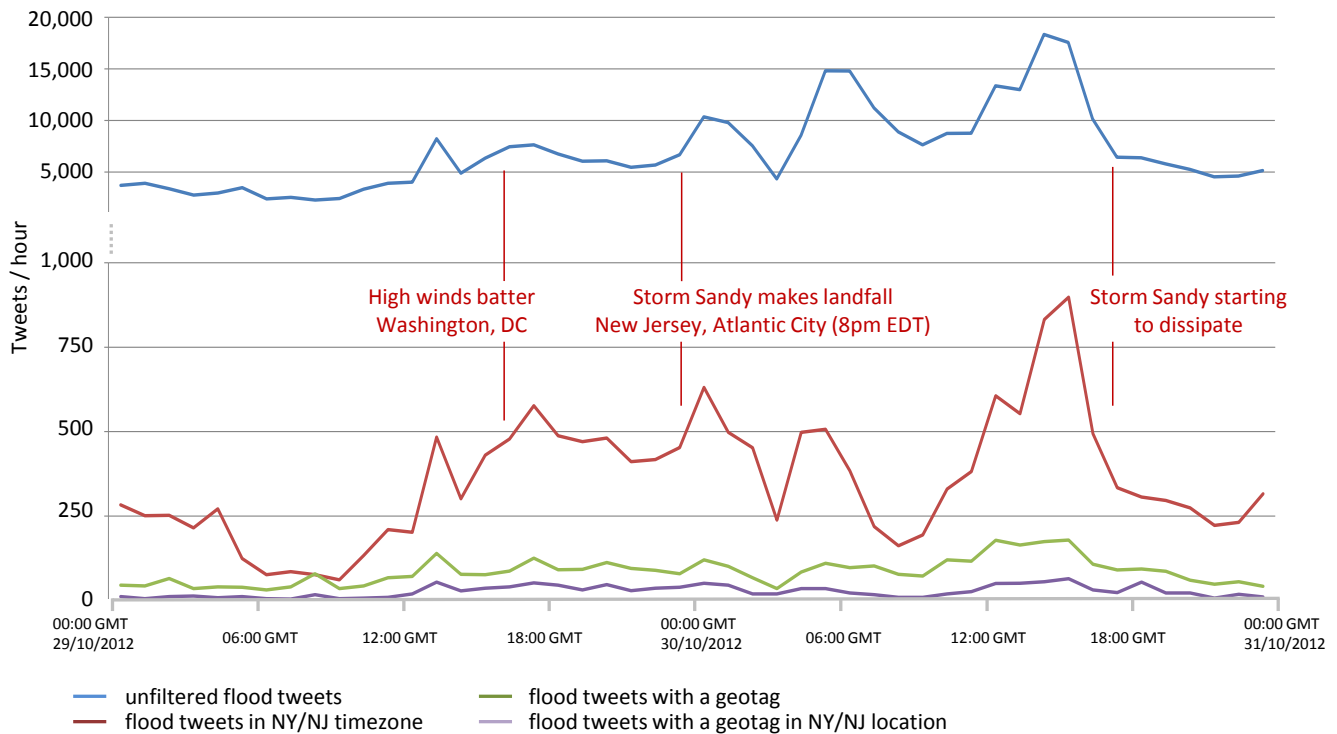


**Figure 1. Twitter Streaming API tweet traffic recorded using 'flood' keywords over 48 hours as hurricane Sandy made landfall between 29 & 30th Oct 2012. Peak tweet traffic was 18,000 tweets per hour, with 5% of tweets using the New York timezone, 1% of tweets containing a geotag and 0.3% containing a geotag located in New York / New Jersey**

## SYSTEM ARCHITECTURE

Our social media crisis mapping platform, see figure 2, is split into a set of offline and real-time services. The offline services prepare a geospatial database for the local region of interest and calculate baseline statistics during a historical period when no disasters occurred. The real-time services crawl tweets live from Twitter, identifying mentions of known locations and displaying them as a live street and/or region level crisis map.
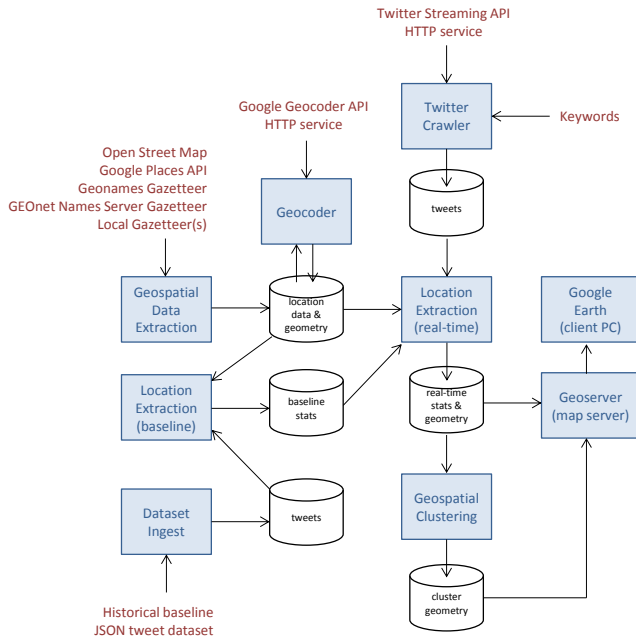


**Figure 2. System Architecture and Information Flow**

During the offline phase we use a set of geospatial data extraction tools to download geospatial data. We use OpenStreetMap to access street level information. We use GooglePlaces API to access volunteered geographic information (VGI) such as buildings and local features. We use the global gazetteers Geonames and GEOnet Names (GNS), as well as local gazetteers, to get region names. This geospatial information is stored in a MySQL database, along with any OpenGIS shape data for later visualization on a map. When downloading geospatial data we request the language native to the local area of interest.

A batch process geocodes each location's address, returning a well formed address string and a specific coordinate on a map. We use the GoogleGeocoding API for our geocoding. For place data (e.g. buildings, rivers) geocoding allows us to fill in blank address fields, or correct them where they contain an error or inconsistency. We have found building data uploaded by the general public varies greatly in its use of the name, street and address fields. For street data geocoding parses the address field into its sub-components, providing us with extra short name variants in addition to the official road name. This is important as people often use short names or abbreviations in Tweets.

The last offline step is to create baseline match statistics for each location in the database over a historical period when no disaster event occurred. Baseline match statistics are a useful tool to reduce false positives associated with location names that pop up often in Twitter conversations (e.g. 'Hollywood'). This baseline is used as a threshold above which location matches can be considered relevant. An ingest tool is used to import the historical dataset to a MySQL database.

The real-time system is driven by a twitter crawler tasked with a set of keywords. We use a set of European multi-lingual keywords for the event type we are looking to record (e.g. for flooding we use 'flood', 'tsunami', 'inondation', 'sel', 'alluvione' etc.). The TwitterStreaming API is used to receive tweets and we continually store them into our MySQL database, splitting SQL tables into 1 month blocks to ensure a fast table query response. We use regex expressions to check for retweets, looking for prefixes like 'RT', as the Twitter retweet metadata is unreliable.

We filter tweets outside of the local region's timezone to help restrict our analysis to people located in the affected area, as opposed to people located in another state / country commenting on news reports. We also filter retweets, which usually do not report new information and thus tend to artificially inflate a locations frequency count.

Our real-time location extraction service runs in parallel to the crawler, processing tweets as they arrive in the database. This service pre-loads locations for the spatial area of interest, tokenizes each of them and creates an in-memory hashtable of tokens ready for efficient real-time matching. Baseline statistics are also pre-loaded into memory. As new tweets are read from the database they are cleaned, tokenized and named entity matching performed, matching location tokens to tweet text tokens. Location matches are logged to a rolling in-memory buffer of configurable size, usually between 6 and 24 hours long, which forms the basis of a rolling sample window. The sample period is usually between 1 and 5 minutes, ensuring we have up-to-date statistics in the database for map display. All match statistics are saved to the database as soon as they are ready along with the OpenGIS geometry to plot on the crisis map.

We run a parallel geospatial clustering service to continuously cluster spatial areas of high activity and produce an easily visible polygon map overlay. This service applies a standard hierarchical clustering algorithm to compute clusters from location geometry.

The mapping visualization is performed using Geoserver, an open source map server. Map layers are driven from the geometry columns in MySQL database tables, plotting buildings (points), streets (lines), regions (points) and clusters of activity (polygons). We render our maps using Google Earth, although Geoserver supports a variety of

mapping formats such as Web Mapping Service (WMS) and OpenLayers. Examples of our rendered crisis maps can be seen in figures 4 and 5.

## GEO-PARSING TWEETS TO GET LOCATIONS

Both our real-time analysis and offline baseline location extraction services use the same geo-parsing algorithm. We support English, Italian, Portuguese and Turkish, languages native to the coastal regions in the North-Eastern Atlantic and Mediterranean region of our Tsunami early warning use case. Since we know a-priori the spatial region of interest we pre-load all possible location entities into an efficient in-memory lookup table. This avoids the need to use named entity recognition approaches, such as a parts of speech classifier coupled with a context grammar, to extract free text location phrases and then geocode them at run-time. Most online geocoding services, including the GoogleGeocoder API, have strict usage rate limitations making Geocoding on the fly not practical for the throughput of tweets we receive.

During system start-up we take each location in our spatial region of interest and tokenize it into 1-gram tokens using the Natural Language ToolKit's (NLTK's) [10] Treebank word tokenizer, then compute a final n-gram token from a sequential combination of the 1-gram tokens. A n-gram token is a phrase made up of 1 to N words, in our case up to a maximum of 5 words. For example, the address 'London Street' generates a two 2-gram token 'london street'. For buildings and street addresses we use our own multi-lingual corpus of building and street types, along with common variants and abbreviations. This allows us to expand token sets to include common usage variants of certain phrases. For example 'London Street' becomes 'london street' + 'london st'.

We remove any tokens that match the NLTK toolkit's multi-lingual stopword list, holding words with low information value such as 'the'. We also remove tokens that match NLTK toolkit's name corpus of common male and female names, avoiding false matches like 'Chelsea' which is both a location and a girl's name. We use weak token stemming to remove plurals as locations are proper names and stronger stemming would cause false positives. We filter any place and address tokens that are identical to region names, since a region match is most likely in this case. We reject any 1-gram token phrases for place and street names as these tend to be ill-defined (e.g. 'station') and prone to over matching. Lastly we compute a 1-gram 'hashtag' token by removing all spaces since hashtags are often used in Tweets (e.g. '#newyork').

During live real-time tweet processing we remove URL's and email addresses from tweets that might generate false tokens. We then use the NLTK toolkit's Punkt sentence tokenizer before executing the Treebank word tokenization as before. We compute all sequential combinations of n-gram tokens from each tweet's text and use this as the basis for location token matching. Our location match algorithm first checks for places tokens, then streets and finally regions. At each stage we remove previously matched tokens from the tweet token pool to avoid text with street names like 'london street' being used to also match a region like 'london'.

In performance testing our location extraction algorithm performs three times quicker than the peak levels of tweet throughput found in our recorded datasets. The processing speed, end to end including all of the database I/O, was about 270,000 tweets/hour for a 20,000 location dataset on a 8Gbyte RAM 2.5GHz CPU laptop. Our performance scales much better than linearly as more locations are added to the database.

We first evaluated multi-cultural geo-parsing accuracy on some tweet datasets recorded by our crawler over the last 2 years. These tweet datasets were manually annotated with places, streets and regions. These datasets vary in native language and size of area affected, with localized blackouts in Milan and widespread floods in New York. The Istanbul earthquake caused no building or street damage, hence we only matched region labels. We counted true/false positives and negatives, where a true positive occurred if the matched location was the same as the expert label, and computed precision, recall and F1 measures.

It can be seen from the results in figure 3 (bottom) that all locations reported a high match precision, but the Turkish dataset had an unusually low recall for region data. This was largely due to the way location names are written in Turkish grammar. For example 'izmir' is a Turkish location, but may appear as 'izmirda', 'izmirdan' or 'izmira' depending on if someone is going to, from or into a location. This result highlights potential limitations of our language-independent matching process.

We benchmarked the accuracy of our geo-parsing on the well-studied Christchurch 2011 New Zealand earthquake. Our 'gold standard' for comparison is a recently published [7] tweet geo-parsing system based on state-of-the-art language specific named entity recognition, lexio-semantic heuristics and a spelling checker. We tested using the same dataset of 2,000 manually labelled tweets, provided by Carnegie Melon University, with annotations showing places, streets and regions. We used the same experimental conditions as [7], including the GNS Gazetteer and a local gazetteer, as well as additional map data from OpenStreetMap. Figure 3 (top) shows our benchmarked results broken down into places, streets and regions. Whilst our approaches F1 measure is similar to the gold standard the precision at street level is much higher. This is an attractive result as a low false positive rate is an important requirement [9] for control room staff to avoid wasted time during a crisis situation.
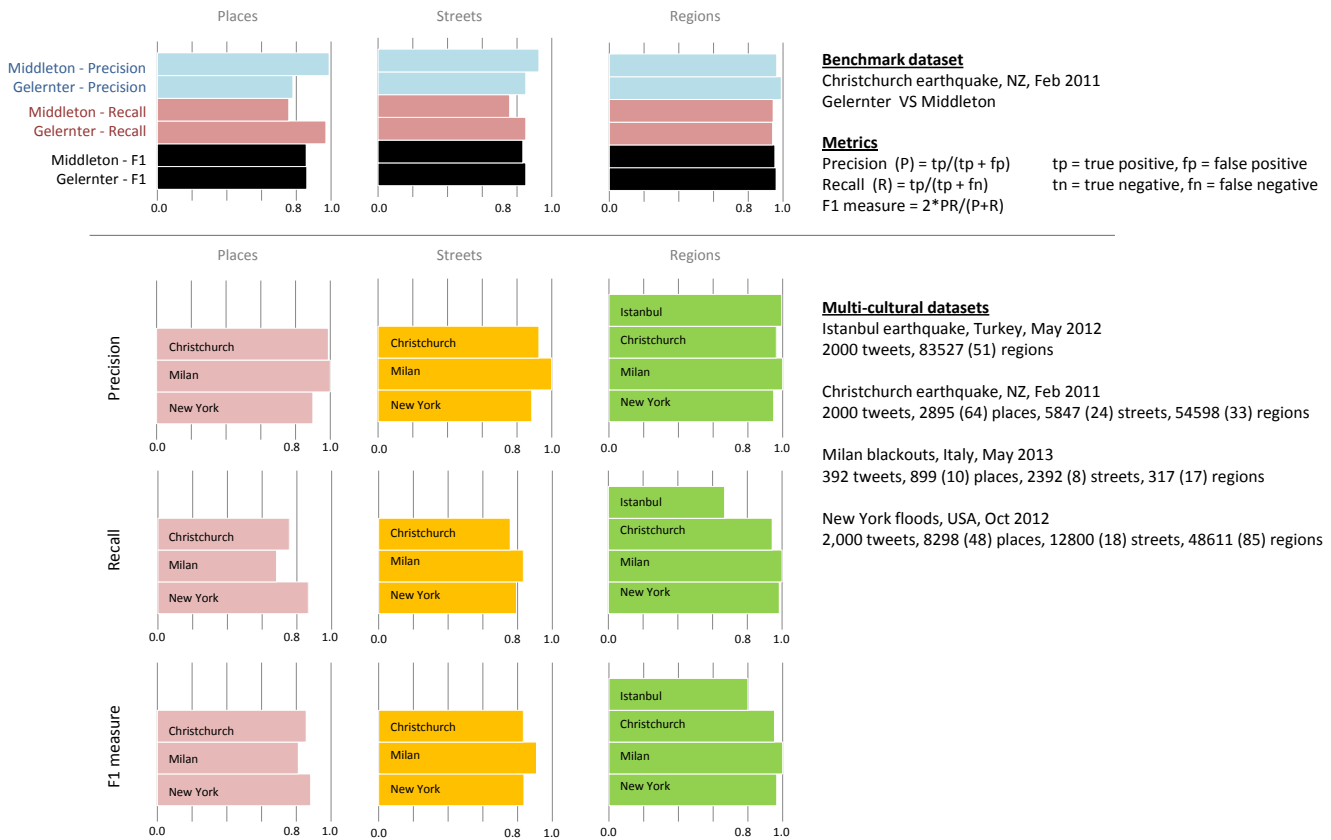
**Figure 3. Geo-parsing evaluation results. Benchmark results (top) for Christchurch Feb 2011 earthquake comparing published Gelernter [7] results to ours [Middleton] using the same 2,000 tweet labelled dataset. Multi-cultural datasets (bottom) show our results for a variety of European locations. Each dataset has a set of labelled tweets and a set of locations, of which a small subset of locations appeared in tweets e.g. Istanbul database has 83,527 regions across Turkey, of which 51 were mentioned in tweets**

## STATISTICAL ANALYSIS OF LOCATION MATCHES

We calculate a statistical baseline for each location to allow us to compute later a threshold level for tweet mentions before which each location is displayed on the crisis map. We use a configurable sample period (e.g. 5 minutes) and sample window (e.g. 6 hours) over which to calculate our statistics. For each location we count the number of tweets per sample period where the location is mentioned, using a historical dataset for the baseline in which we know no disasters happened. We then calculate a simple moving average and triangular weighted moving average across the dataset as a whole for the moving sample window. The case studies reported later both use a 1 month baseline tweet dataset with just under 1 million tweets each.

The same per-location match statistics are calculated for a moving sample window of real-time tweet data. The deviation of real-time metric values from baseline metric values is calculated every sample period, and compared to a configurable threshold before displaying each location on the crisis map.

Our central hypothesis is that locations mentioned many times in a sample window are more likely to be coherent and credible disaster related location reports than those with only 1 or 2 mentions. In the case studies reported next we use the simple moving average metric, and show how raising the threshold level for acceptance increases precision. Ultimately this threshold value will be tailored to suite each crisis management control room, reflecting the error-tolerance of the final end user decision makers.

## CRISIS MAPPING CASE STUDIES

We conducted two case studies to evaluate the quality of our tweet maps. The first event studied was hurricane Sandy (Oct 2012), which caused major flooding in New York and New Jersey. The second event was the Oklahoma tornado (May 2013), which devastated the town of Moore south of Oklahoma.

For the New York flooding event we ran our crisis mapping with a sample window of 6 hours and a sampling rate of 5 minutes. Three maps were computed using a high/medium/low threshold setting for the allowed deviation of simple moving average from baseline (dev_sma). The tweet dataset covered 5 days, contained 597,022 tweets (15,175 after timezone & retweet filtering) of which 4,302 had a location mention. Our New York location database has 8,298 places, 12,800 streets and 48,661 regions available for matching. We have all regions (cities, suburbs, neighbourhoods etc.) for New Jersey from

our gazetteers and coastal street data from OpenStreetMap and GooglePlaces covering all of Manhattan.

Each map was compared to a ground truth storm surge map from the official post-event impact assessment produced by the US National Geospatial Agency (NGA). Figure 4 shows the post-event storm surge map alongside our 5 day tweet map. To empirically evaluate our map we segmented it into a 8x8 grid and compared each grid cell to the ground truth map. True positives were reported for any cell that has both a tweeted location reported and some storm surge activity on the NGA impact assessment map. We counted the number of true/false positives and negatives and calculated precision, recall and F1 measures. As expected, when we increase the mapping threshold (dev_sma) the map precision increases at the expense of recall.



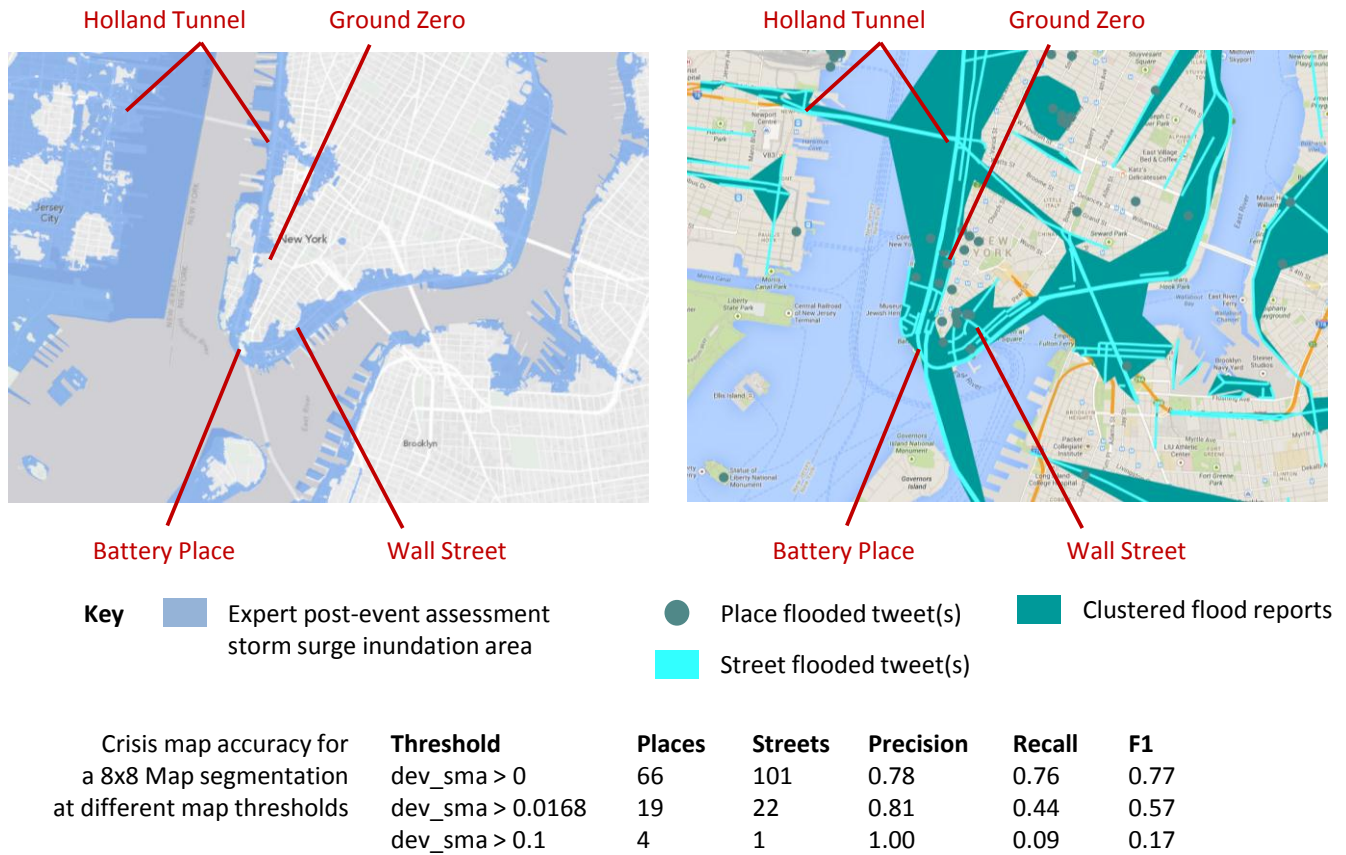| Crisis map accuracy for a 8x8 Map segmentation at different map thresholds | Threshold | Places | Streets | Precision | Recall | F1 |
|---|---|---|---|---|---|---|
| | dev_sma > 0 | 66 | 101 | 0.78 | 0.76 | 0.77 |
| | dev_sma > 0.0168 | 19 | 22 | 0.81 | 0.44 | 0.57 |
| | dev_sma > 0.1 | 4 | 1 | 1.00 | 0.09 | 0.17 |

**Figure 4. Crisis map comparison for New York's 2012 flooding. The left image is the ground truth post-event NGA impact assessment showing storm surge inundation. The right image is a 5 day tweet flood map (dev_sma > 0) for tweets between 29-10-2012 to 02-11-2012. Red annotations show major incidents. Source: FEMA Modelling Task Force (MOTF) storm Sandy impact analysis field-verified interim high resolution report, Nov 2012. Mapping courtesy of ArcGIS ESRI portal and Google Maps.**

Plaza Towers Elementary School

Plaza Towers Elementary School

Briarwood Elementary School     Medical Centre

Briarwood Elementary School     Medical Centre

**Key** | Expert post-event assessment tornado damage area | ● Place damaged tweet(s) | ■ Clustered damage reports

■ Street damaged tweet(s)

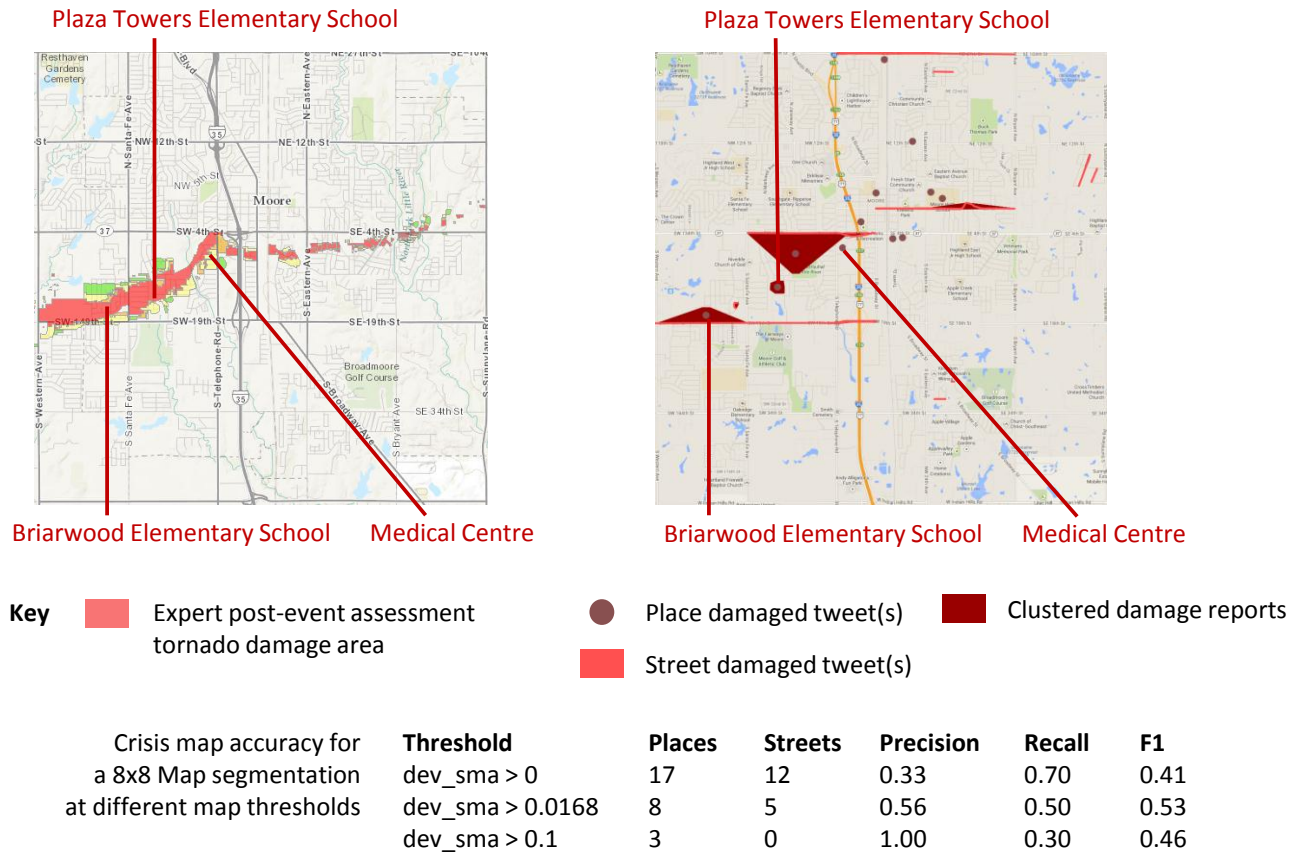| Crisis map accuracy for a 8x8 Map segmentation at different map thresholds | Threshold | Places | Streets | Precision | Recall | F1 |
|---|---|---|---|---|---|---|
| | dev_sma > 0 | 17 | 12 | 0.33 | 0.70 | 0.41 |
| | dev_sma > 0.0168 | 8 | 5 | 0.56 | 0.50 | 0.53 |
| | dev_sma > 0.1 | 3 | 0 | 1.00 | 0.30 | 0.46 |

**Figure 5. Crisis map comparison for Oklahoma's 2013 tornado. The left image is the ground truth official post-event NGA impact assessment showing building damage.  The right image is a 5 day tweet damage map (dev_sma > 0) for tweets between 20-05-2013 to 24-05-2013. Red annotations show major incidents. Source: National Geospatial Agency (NGA) damage assessment using aerial (FEMA, BAE) & satellite images (World View 1), May 2013. Mapping courtesy of ArcGIS ESRI portal and Google Maps.**

For the Oklahoma tornado event we also ran our crisis mapping using a 6 hour sample window and sampling period of 5 minutes, generating three maps with the same threshold values as the New York case study. The tweet dataset covered 5 days, contained 877,527 tweets (92,300 after timezone & retweet filtering) of which 42,434 had a location mention. Our Oklahoma location database has 625 places, 3,930 streets and 18,599 regions available for matching.

Our ground truth map was a US National Geospatial Agency post-event impact assessment showing structural damage across the town of Moore. Figure 5 shows this post-event damage assessment alongside our 5 day tweet map. We segmented each map into an 8x8 grid as before and compared each cell, counting the true/false positives and negatives. The results again show that we can raise the overall map precision at the expense of recall by raising the mapping threshold level.

**CONCLUSION AND FUTURE WORK**
Both our case studies demonstrate that high precision (i.e. 90% or higher) geo-parsing from real-time Twitter data is possible by exploiting large databases of pre-loaded location information for 'at risk' areas. Such data is readily

available online from mapping services, volunteered geographic information sources and gazetteers. These case studies also show that crisis maps generated from social media data can compare well to gold-standard post-event impact assessments from national civil protection authorities. This matches well with the requirements of use cases such as Tsunami early warning centres, where real-time crisis mapping with minimal rates of false positives are needed.

When applying our approach in the future it is important to consider the spatial size and significance of the natural disaster, as the quality of the crisis map is directly related to the number of people tweeting information about the disaster zone. Large scale news worthy events will usually receive more tweets than events in small localized areas, or areas in remote locations with limited mobile communication. However, as the uptake of social media around the world increases with time we feel the role this type of social intelligence has to play in assisting civil protection authorities will also increase.

For next steps we are experimenting with approaches for language-specific context filtering, to be applied as a type of secondary filter on the sub-set of tweets that match the

initial geo-parsing stage. This context filter would look at the natural language context in which locations are mentioned and try to classify patterns associated with specific classes of response, such as flooded transport systems, positive/negative reports, cries for help and reports with high levels of urgency. We will also look at using retweet's for adding a credibility value to original reports.

Currently each instance of our location extraction process looks for location matches in a single region of interest. We will in the future scale our approach across a computing cluster to handle many spatial regions of interest simultaneously. This offers the possibility of country-wide area map coverage and/or collections of processes able to be adaptively tasked to monitor new spatial areas on-demand.

We plan in the next few months to deploy a prototype as part of the award winning TRIDEC project, allowing potential end users to assess the social media crisis management platform first-hand and start the process of user-evaluation and progress towards adopting this early stage technology.

## ACKNOWLEDGMENTS

## BIOGRAPHY
Stuart E. Middleton is a senior research engineer at the University of Southampton IT Innovation Centre. His main research interests are social media, sensor systems, data fusion and ontologies. Stuart has a PhD in Computer Science from the University of Southampton. Contact him at sem@it-innovation.soton.ac.uk

Lee Middleton is a senior research engineer at the University of Southampton IT Innovation Centre. His main research interests are computer vision, machine learning, and pattern analysis. Lee has a PhD in Electrical and Electronic Engineering from the University of Auckland, New Zealand. Contact him at ljm@it-innovation.soton.ac.uk

Stefano Modafferi is a senior research engineer at the University of Southampton IT Innovation Centre. His main research interests are information modelling and software architectures. Stefano has a PhD in Information Engineering from the Politecnico di Milano. Contact him at sm@it-innovation.soton.ac.uk

## REFERENCES

1. P. Meier, "New Information Technologies and their Impact on the Humanitarian Sector", International Review of the Red Cross, Dec 2011, Vol 93, No. 844

2. T. Sakaki, M. Okazaki, Y. Matsuo, "Tweet Analysis for Real-time Event Detection and Earthquake Reporting System Development", IEEE Transactions on Knowledge and Data Engineering, April 2013, Vol. 25, No. 4

3. P.S. Earle, D.C. Bowden, M. Guy, "Twitter earthquake detection: earthquake monitoring in a social world", Annals of Geophysics, 2011, Vol. 54, No. 6

4. N.R. Adam, B. Shafiq, R. Staffin, "Spatial Computing and Social Media in the Context of Disaster Management", IEEE Intelligent Systems, Nov 2012, Vol. 27, Issue 6, pp. 90-96

5. J. Yin, A. Lampert, M. Cameron, B. Robinson, R. Power, "Using Social Media to Enhance Emergency Situation Awareness", IEEE Intelligent Systems, Nov 2012, Vol. 27, Issue 6, pp. 52-59

6. X. Liu, F. Wei, S. Zhang, M. Zhou, "Named Entity Recognition for Tweets", ACM Transactions on Intelligent Systems and Technology, Jan 2013, Vol. 4, No. 1, Article 3

7. J. Gelernter, S. Balaji, "An Algorithm for Local Geoparsing of Microtext", GeoInformatica, Springer, Jan 2013

8. G. Shi, K. Barker, "Extraction of geospatial information on the web for GIS applications", 10th IEEE International Conference on Cognitive Informatics & Cognitive Computing (ICCI*CC ), Aug 2011, pp. 41-48

9. A. Zielinski, S.E. Middleton, O. Necmioğlu, M. Hammitzsch, "Spatio-Temporal Decision Support System for Natural Crisis Management with TweetComP1", Proc. EWG-DSS 2013 Workshop "Exploring New Directions For Decisions In The Internet Age "F.Dargam, B.Delibasic, J.E.Hernández, S.Liu, J. Papatanasiou, R.Ribeiro, P.Zaraté (editors) May 2013, pp. 33

10. A. Bird, E. Klein, E. Loper, "Natural Language Processing with Python --- Analyzing Text with the Natural Language Toolkit", 2009, O'Reilly Media, ISBN: 978-0-596-51649-9