

# A Multi-Strategy Approach for Location Mining in Tweets: AUT NLP Group Entry for ALTA-2014 Shared Task

Parma Nand, Rivindu Perera

Natural Language Processing Group

Auckland University of Technology

Auckland 1010, New Zealand

{pnand, rperera}@aut.ac.nz

## Abstract

This paper describes the strategy and the results of a location mining system used for the ALTA-2014 shared task competition. The task required the participants to identify the location mentions in 1003 Twitter test messages given a separate annotated training set of 2000 messages. We present an architecture that uses a basic named entity recognizer in conjunction with various rule-based modules and knowledge infusion to achieve an average F score of 0.747 which won the second place in the competition. We used the pre-trained Stanford NER which gives us an F score of 0.532 and used an ensemble of other techniques to reach the 0.747 value. The other major source of location resolver was the DBpedia location list which was used to identify a large percentage of locations with an individual F-score of 0.935.

## 1 Introduction

The objective of the ALTA competition was to identify a wide range of location mentions in Twitter messages. Among the various types of information that can be mined from Twitter messages, location mining has attracted a lot attention because of its application in identifying the geographical location of a topic in the Tweet, whether it be an accident, a natural disaster, a concert or an open invitation party gone out of hand. Location mining is a special case of a generic NLP task called Named Entity Recognition (NER), which involves identifying and classifying entities into person, organization and location type entities. This particular competition focussed on only identifying location type entities, however had a wider scope compared to a typical NER task. The objective of the task was to identify the string which

specifies the most specific possible locations in the message. A typical NER tool such as the Stanford<sup>1</sup> and OpenNLP<sup>2</sup> only identifies word based locations which has to be then composed into noun phrases consisting of multiple words identifying the full location. In addition to the locations identified by noun phrases, the task required the identification of the specific locations defined by propositional attachments such as “40km south of tenetfield2” and “100 mi from nations capital”.

Our approach to solving the location mining task at hand was to use a NER system to attain a benchmark performance and then to use various techniques to fine tune the system to account for the noise in Twitter messages. The architecture forms part of a bigger project called the *Twitter Miner*<sup>3</sup>. This project is a higher level research project currently underway at Auckland University of Technology which is meant to extract information for various purposes from microblogging type texts such as Twitter.

For the NER, we used the Conditional Random Field (CRF) based, Stanford NER system as this has been tested to give the highest accuracy on Twitter messages (Lingad et al., 2013). In this paper Lingad et al. (2013) reported the re-trained Stanford NER to achieve an F-value of 0.902 compared to 0.576 for the pre-trained NER. We however found that, re-training the Stanford NER with bare training data provided by the organizers gave us F scores around the 0.4 mark at token level compared to approximately 0.57 for the pre-trained model. Analysis of the errors for both the models showed that due to loose capitalization in twitter messages, a lot of the locations could not be identified simply because they did not exist in the training data. Since the pre-trained model was trained with much larger training set

<sup>1</sup><http://nlp.stanford.edu/software/CRF-NER.shtml>

<sup>2</sup><http://opennlp.apache.org>

<sup>3</sup><http://staff.eleena.aut.ac.nz/Parma-Nand/projects.html>

it could detect a larger number of locations giving us a higher precision value. We tried incremental training of the Stanford NER, however ran into technical difficulties with memory and computational time required. Hence, we adopted the approach of ensemble system consisting of a pre-trained Stanford NER, knowledge infusion, regular expression identifier and use of rules. Instead of a pipeline architecture with no re-processing we adopted a parallel architecture to cater for the copious amount of noise in Tweets. The parallel architecture enabled us to revisit previous decisions and correct them.

The rest of the paper organized as follows. Section 2 gives an overview of related works for location mining in Twitter messages as well as some generic NER works. Section 3 describes the task description followed by our methodology. The results are detailed in Section 5 followed by conclusion.

## 2 Related Works

Location Mining is a subtask of the more generic information extraction task of named entity recognition. There are numerous works on NER in the formal domains, however this section gives an overview of the recent work specifically in the informal domain of microblogging, mostly for “Twittersphere”.

Ritter et. al. (2011) presented an NER system named *T-NER*, which uses Freebase as an information source to enhance supervision for the Stanford NER to classify entities into 10 classes, one of which was Geo-location. Their system achieved an overall F score of 0.66 and a Geo-location F score of 0.77. Li et al. (2012) present a random walk model which exploits the gregarious properties associated with Tweets in addition to the textual content to determine the named entities. The gregarious property is based on named entities mentioned together in media other than Twitter. The authors used Microsoft N-Gram and Wikipedia as the corpus to compute the gregarious property value. This system attained an F score of 0.419 on their data compared to 0.466 for the previously mentioned system and 0.423 for the Stanford NER on the same data.

The task of location mining specifically from Twitter messages has attracted a lot of attention because Twitter is current up to the minute hence can be used for information about upto date events

around the globe. One of its immediate use is for almost real time disaster detection so that services can be deployed as soon as possible. There are multiple other uses for Twitter location mining such as location based advertising and geography based sentiment mining. Lingad et al. (2013) present test results for using 4 off-the-shelf NER’s to determine locations with varying degrees of granularity from Country, State, City, Area, Suburb to Point of Interest. The results from this paper showed that the retrained Stanford NER was the the best performer at 0.872 and the standard Stanford NER 4-class classifier attained a value of 0.691. We used the results of this paper to choose the NER used to be used for our location miner.

Twitter messages may also have meta data indicating the location from which a Tweet was sent. This is only present in Tweets sent from mobile devices equipped with GPS hardware, however note that this feature can also be turned off for privacy reasons by the user. Ikawa et al. (2012) present a model which exploits the GPS location as well as the textual content of the Tweet. This model uses associations between locations and relevant keywords from past messages during training, which is then used to estimate where the new message was issued from. The identified location is then allocated a geographical square and the errors were calculated based on the distance within 10 kilometres. The study reports a precision value of 0.45 with a dataset of 20,535 messages out of which 16,380 were used for training. A large part of the error in this was that the location mentions in the text of the Tweet might not necessary correlate with the location of the user.

Mahmud et al. (2012) also present a system for predicting the home locations of Twitter users. Unlike the previous system, this paper uses an ensemble of statistical and heuristic classifiers to predict Google’s geo-coding bounding box for Twitter corresponding to the Twitter users. The paper reports accuracies at various granularities which range from 0.54 to 0.78 recall values.

Sankaranarayanan et al. (2009) present a system which also does location mining for a different purpose. The object of this work is to cluster Twitter messages according to the news content based on geographical locations. This work again uses ensemble learning, similar to our approach. It uses references to geographic locations in the text, called toponyms, to determine the co-

ordinates which then used to resolve using various techniques such as NER, POS tagging and a look up database containing countries, city, river etc. In addition to this the textual content of the Tweet is extended using its metadata about the location of the Tweeter. The metadata information is added as textual content of the Tweet which is then treated similarly as the rest to the message. The paper does not report any location specific accuracy, however illustrates another use for location mining and the enforces the use of ensemble architecture for the purpose the purpose.

Apart these there are several other papers (e.g., Kinsella et al., 2011; Li et al., 2011) who have done work on location mining focussed on either the location from where the Tweet was sent or the geographical location of the Tweeter. Something that is common in these works is that they all use some kind of ensemble of techniques rather than any one particular technique. This paper reports the results of a system with a similar architecture which uses an ensemble of techniques, however the task being tackled there is slightly different. It is more akin to location extraction from the text rather than anything to do with the location of the Tweeter or the location from which the Tweet was sent. They may be the same for some Tweets, however the tasks are quite distinct in that a Tweeter can tweet about a location quite different from his or her registered location or current location.

### 3 Shared Task Specification

The organizers provided a set of 2000 Twitter Id's with identified locations for training and another set of 1003 Twitter Id's for testing. The locations in the test set were not released, hence the participants had to make their own test sets based on the training data. The participants were allowed two uploads per day to test the accuracy, which was calculated with a confidential subset of the testing data giving the participants an approximation of the level of accuracy achieved. The output was submitted as a comma separated (csv) file with tweet Id's and the locations separated by a space on each line. The order of the locations on a line was immaterial so "New Zealand" and "Zealand New" are both correct.

The overall goal of the task was to identify all mentions of locations in the text of a twitter message. This includes all single and contiguous word

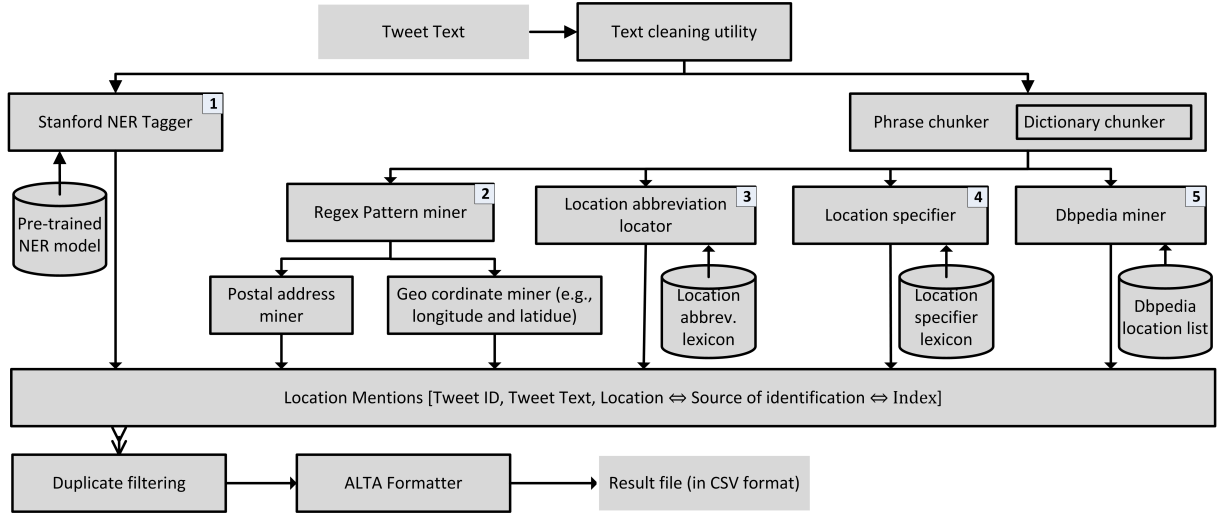
mentions such as "guatemala" and "new zealand" and similar mentions as abbreviations such as "nz" for New Zealand. The following examples give an overview of the wider scope of the task.

- "#eqnz" - location is "eqnz"  
In this case we need to check for location within strings, however extract the whole token if one is found.
- "http://www.abc.net.au/melbourne/" - location is "http://www.abc.net.au/melbourne/"
- "cork city...#Cork" - locations - "cork city cork2"  
In this case locations such as city had to be identified and words in locations appearing more than once had to be tagged with the count.
- "Morrison's Island - S Terrace" - location is "morrison's island s terrace"  
Punctuations need to be removed leaving the strings as they appear in the text.
- "U.S. EPA on Twitter" - location is "NONE"  
In this case "U.S." is not a location, but a user.
- "Our house" - location is "our house"  
Common nouns with possessive pronouns need to be identified since they are specific locations.
- "Southwest towns of San Marcos" - location is "Southwest towns of San Marcos"  
In this case we need to retain the preposition "of".
- "60 miles east starting from Stamford CT" - location is "60 miles east Stamford CT"  
This involves removing the verb and the preposition from the location.

### 4 Methodology

The schematic representation of the method is shown in Fig. 1. We utilized the parallel processing strategy with multiple modules to identify locations.

The proposed location mining architecture adheres to the parallel processing of five major modules. These modules take the cleaned twitter text generated by text cleaning utility. The Named Entity Recognition (NER) module attempts to identify location entities using pre-trained model file.



**Figure 1:** Schematic representation of the Twitter location mining framework

The Regular Expression (Regex) pattern miner is composed of two subcomponents; postal address miner and geo-coordinate miner. Locations can also be specified using abbreviations. This is more common in microblogging type texts. To address this need, the framework utilizes a module that can identify the location abbreviations (e.g., chch, au, nz) using a location abbreviation lexicon. The location specifier module contains the logic to identify locations based on list of location specifiers such as beach, coast and street. We have also utilized the exponentially growing Linked Data cloud in location mining through DBpedia miner.

All these modules are configured to write the results to a single data structure called “*Location Mention*” consisted of tweet id, tweet text, location phrases. Each location phrase has another three attributes; location text, source of identification and the index. The source of identification attribute is used to identify from which module the phrase is identified as a location and index specifies the index of the location phrase in the tweet text. The framework is also equipped with two other modules to merge locations and to format the final result according to the ALTA formatting guidelines. The following sections explore aforementioned modules in detail.

#### 4.1 Cleaning Twitter text

Compared to normal text, microblogging type text contains different features. The most common features that can be seen in Twitter are hashtags and mentions. Existence of such features can disturb the location identification process. To overcome

such noise we employed a text cleaning utility that can identify these features and convert them back to the normal textual form. The library is based on one of our previous work (Nand et al., 2014) in Twitter text mining that uses Hidden Markov Model (HMM) based tagging.

#### 4.2 Location mining using NER

Named Entity Recognition (NER) is a widely used technique to extract elements in text into predefined categories such as organizations, cities, person names. The proposed framework utilizes NER to identify cities and countries from the tweet text. We have integrated the Stanford NER toolkit (Finkel and Manning, 2009) with MUC-3 class pre-trained model to locate cities and countries. Fig. 2 shows an example scenario for NER based location mining. This Tweet example and the rest of the following examples are taken from ALTA test set unless specified otherwise.

```
Tweet ID: 255741958625062912
Cleaned Text: Aus Weather Warnings on
Twitter: Queensland : Fire Weather Warning:
http://t.co/knVkydSB
Locations: Queensland
Source: [NER]
```

**Figure 2:** Example scenario for Named Entity Recognition based location mining

#### 4.3 Phrase and dictionary chunking

Except for the NER module which takes the raw Tweet text, all other modules are designed to pro-

cess chunked text. In essence, the framework first chunks text based on a predefined rule set. This rule set is implemented with the use of LingPipe library<sup>4</sup> which makes it a scalable module. Table 1 shows the list of rules used for this task. The chunked text is analysed again to check whether there is a possibility to chunk it further. The framework utilizes a dictionary based chunking technique for further chunking if needed. Each token appearing in the phrase is checked whether it can be chunked into multiple words using a lexicon list. The lexicon list is a combination of words extracted from dictionary and a location abbreviation lexicon (see Section 4.5).

Rule	Description
(NN)(NNP)(NNS)	Adjacent noun, singular proper noun or plural noun phrases
(JJ)(NN/NNP/NNS)	Adjective with a noun phrase
(JJR)(NN/NNP/NNS)	Comparative adjective with a noun phrase
(JJS)(NN/NNP/NNS)	Superlative adjective with a noun phrase

**Table 1:** Phrase chunking rules

#### 4.4 Regex Pattern miner

Locations can also be represented in postal addresses and geo coordinates such as longitude and latitude. Therefore, the framework is equipped with a regular expression pattern miner which takes a predefined template and identifies whether a given text phrase is a location or not. An example scenario of Regex pattern miner is depicted in Fig. 3.

#### 4.5 Location abbreviation miner

Abbreviations are more abundant in microblogging type texts due to space limitations. Essentially, when mining locations it is important to consider abbreviations that can specify a locations (e.g., chch  $\Rightarrow$  Christchurch, nz  $\Rightarrow$  New Zealand). These abbreviations can appear in the text as a single word or in a combination with another word. The latter case makes it difficult to identify because the combination can form different textual

<sup>4</sup><http://alias-i.com/lingpipe/index.html>

```

Tweet ID: 255729069977636864
Cleaned Text: Weather Underground on Twitter:
Watching area near 9.3N 48.3W for tropical develop-
ment Invest 98L: winds 30 mph moving W at 24 mph
http://t.co/ml1PrSJW hurricane
Locations: 9.3N 48.3W
Source: [R-GEO]
.....
Tweet ID: 264171210596831232
Cleaned Text: SEQ incidents on Twitter: Reports
of Fire Services Incident ? near 333 Manly Road Manly
West http://t.co/c2c0GAcm
Locations: 333 Manly Road Manly West
Source: [R-POSTAL]

```

**Figure 3:** Example scenario of Regular Expression based pattern mining

```

Tweet ID: 269641336326615041
Cleaned Text: QLD Times on Twitter: Deebing
Heights house badly damaged by a fire this morning:
http://t.co/U7X06IOh. QldFire
Locations: QLD Qldfire
Source: [ABBREV]

```

**Figure 4:** Example scenario of location abbreviation based mining

representations. This is addressed in our framework by performing the abbreviation mining in two steps. First, phrase chunker attempts to identify phrases from the given text as discussed in Section 4.3. Then single words appearing in these phrases are chunked again using a location abbreviation lexicon list and a set of English words. If it is found that the phrase contains an abbreviation of a location (e.g., nzquake  $\Rightarrow$  nz, quake), the complete phrase is tagged as a location. Fig. 4 shows an example case taken from the test dataset.

#### 4.6 Location specifier based identification

Tokens like beach, street and coast are generally used to specify locations. As a preprocessing task we used the ALTA training data to create a location specifier list. This list is further enriched with another set of location specifiers extracted using a thesaurus. We iteratively searched *The-saurus.com*<sup>5</sup> with seed words selected from the ALTA training set based location specifier list. From this we were able to list 183 location speci-

<sup>5</sup><http://www.thesaurus.com/>

- X? [km,mi,ft,miles] [south][north][east][west] [of,from] Y?
- X? [km,mi,ft,miles] (away) from Y?
- X? near Y?
- X? [km,mi,ft,miles] [outside] [of,from] Y?

**Figure 5:** Sample set of templates created to identify locations which are prepositional attachments

Tweet ID: 260164937983340546  
 Cleaned Text: Daily Examiner on Twitter:  
 Clarence Valley bushfire update. <http://t.co/TxnorR7X>  
 Locations: Clarence Valley  
 Source: [SPEC]

**Figure 6:** Example scenario of location specifier based identification

fiers. This list is used with a template based matching to identify whether an extracted phrase is a location. An example scenario is shown in Fig. 6 which shows how the specifier “Valley” is used to identify a location.

In addition to these predetermined location specifiers, this module also selects locations which use prepositional attachments. This identification is accomplished using a set of templates that we created based on the ALTA training data and twelve Wikipedia pages related to locations. Sample set of templates are shown in Fig. 5. The complete template set is composed of 28 unique templates.

#### 4.7 DBpedia based location mining

DBpedia<sup>6</sup> based location mining module uses the most extensive knowledge about locations compared to the other four processes described previously. DBpedia is a Linked Data resource which is built based on the Wikipedia<sup>7</sup> text. In general a Linked Data resource is made up of triples which represent vast domain knowledge and categorized into predetermined classes.

In DBpedia, the ontology class “Place” is used to categorize all location specific entities. Under this main entity class “Place”, there are 149 sub classes that denote places such as theatres, lakes, pyramids, etc. Since DBpedia does not offer a database from which we can easily filter out these locations, we have created a database with all the

<sup>6</sup><http://dbpedia.org/About>

<sup>7</sup>[http://en.wikipedia.org/wiki/Main\\_Page](http://en.wikipedia.org/wiki/Main_Page)

Tweet ID: 264099474564075520  
 Cleaned Text: 702 ABC Sydney on Twitter: The bushfire at Lake Macquarie near Teralba is now under control. Homes are no longer under threat  
 Locations: Lake Macquarie  
 Source: [DBPEDIA]

**Figure 7:** Example scenario for DBpedia based location mining

information required to filter only locations from the DBpedia data files. Table 2 shows few records from this database. The phrases which match with the entity literal value (shown in Table 2) were tagged as locations. Fig. 7 shows an example scenario of DBpedia based location mining where the phrase “Lake Macquarie” is identified as a location.

#### 4.8 Merging module

Since the framework follows the parallel processing architecture utilizing five individual modules, there was a need for a merging module which can generate accurate representations of multiple location identifications which ultimately point to one location in Tweet text. The ultimate goal of this merging module is to present the most informative location as the final result. This was accomplished using a text merging utility that takes the index of each location mention and merges them according to the same order of tokens appearing in the tweet text. An example scenario is shown in Fig. 8 which depicts the process of merging two location identifications; one identified by DBpedia miner and other identified by location specifier list based identification module. In this example scenario, compared to the location identified by the DBpedia miner, the location specifier based module has identified a more informative location. In this case we analyse the index of the tokens from both process and since the DBpedia based identification is a subset of the specifier list based identification, we merge two and output the result as “40km south of Tenterfield”.

#### 4.9 ALTA formatter

ALTA shared task requires special formatting of the result as a comma separated file which was accomplished by this module. In essence, the formatter was based on the following four rules:

- remove all punctuations from the phrase

Ontology class	Entity literal value	Link	Data file
Theatre	Stephen Joseph Theatre	<a href="http://dbpedia.org/page/Stephen_Joseph_Theatre">http://dbpedia.org/page/Stephen_Joseph_Theatre</a>	Stephen_Joseph_Theatre.rdf
Lake	Lake Macquarie	<a href="http://dbpedia.org/resource/Lake_Macquarie_(New_South_Wales)">http://dbpedia.org/resource/Lake_Macquarie_(New_South_Wales)</a>	Lake_Macquarie.rdf
Museum	BritishMuseum	<a href="http://dbpedia.org/resource/British_Museum">http://dbpedia.org/resource/British_Museum</a>	British_Museum.rdf
Airport	Glasgow Airport	<a href="http://dbpedia.org/resource/Glasgow_Airport">http://dbpedia.org/resource/Glasgow_Airport</a>	Glasgow_Airport.rdf
Mountain	Mount Vesuvius	<a href="http://dbpedia.org/resource/Mount_Vesuvius">http://dbpedia.org/resource/Mount_Vesuvius</a>	Mount_Vesuvius.rdf

**Table 2:** Sample set of records from DBpedia entity database

<p>Tweet ID: 255914885928583168  Cleaned Text: Live Traffic NSW on Twitter:  TENTERFIELD: NewEnglandHwy closed in both  directions 40km south of Tenterfield due to a bushfire.</p> <p>Locations: Tenterfield  Source: [DBPEDIA]</p> <p>Locations: 40km south of Tenterfield  Source: [SPEC]</p> <p>Output: 40km south of Tenterfield</p>
---

**Figure 8:** Example scenario for merging two locations

- if locations are repeated in a tweet, number them from the second occurrence
- if there is no location for the tweet, then mark it as *NONE*
- lowercase all extracted location phrases

The resulting phrases were converted to a Comma separated file with two fields; tweet id and the location phrases.

## 5 Results

Our location miner achieved an average F-value of 0.747 which was in the second place compared to the winner which had an F score of 0.77807. The F-value was calculated based on the “bag of words criteria” for each tweet as described in the Section 3. If a tweet did not contain any locations, the participants were required to label them with “NONE”. The overall results for the test dataset is shown in Table 3.

The precision and recall values were computed for individual Tweets, which were then averaged to compute the overall F-value. Hence the precision and recall values for a Tweet with no location mentions was taken as 1.0 for “NONE” to indicate no location. Any other strings instead of “NONE”

Test Data Property	Value
No. of Tweets	1003
No. of Location tokens	3179
No. NONE Tweets	115
Av. Recall	0.7279
Av. Precision	0.7905

**Table 3:** Dataset details and Results

Module	Number	F Score
Stanford NER	1003	0.532
Postal Address Miner	6	0.167
Geo coordinate	8	0.242
Location Abbreviation Miner	607	0.710
Location Specifier	202	0.884
DBpedia Miner	1867	0.935
Total	3680	–

**Table 4:** The number of locations resolved by individual modules with respective F score.

were counted as false positives. Table 3 shows that there were 115 Tweets (11.5%) with no locations. It should be noted that this strategy for accuracy approximation will tend to boost the F-value if a large number of Tweets have no location mentions.

Table 4 gives the number of the locations specified by the individual modules of the ensemble system used for the task. The total is much higher than the total number of locations because of the false negatives and some of the locations were identified by more than one module. The majority of false negatives were identified by the *Location Specifier* module which was primarily based around rules based on the use of the prepositions

of place such as, “at”, “on” and “in”. The *DBpedia miner* module was able to identify a total of 1867 locations which is even higher than then the *Stanford NER* module. Furthermore, *DBpedia miner* has achieved the highest F score of 0.935 compared to other four modules. The results illustrate that as comprehensive information sources become available, their use in conjunction with machine learning algorithms can be effectively used for improved accuracy.

## 6 Conclusions and Future Work

This paper showed the use of an ensemble approach to solve the problem of location mention identification. We presented an ensemble architecture that uses a basic general purpose NER, with a combination of various rule based modules in conjunction with *DBpedia* knowledge base to achieve an F score of 0.747. A critical aspect of any ensemble architecture is how to combine the results at the end. This was also illustrated by the *Merger* module which takes outputs from the various ensemble modules and combines them into a noun phrase location phrase as was required by the shared task specification. The design of the architecture enables us to exclude a class of location mentions and also to include any new ones that a task at hand might dictate. The final results can also be easily modified to output single token locations or full noun phrase locations. In future we intend to further improve the accuracy and to classify the locations into types such as country, site and address for specific applications.

## References

- [Finkel and Manning2009] Jenny Rose Finkel and Christopher D Manning. 2009. Joint Parsing and Named Entity Recognition. In *Proceedings of the North American Association of Computational Linguistics (NAACL 2009)*.
- [Ikawa et al.2012] Yohei Ikawa, Miki Enoki, and Michiaki Tatsubori. 2012. Location inference using microblog messages. In *Proceedings of the 21st international conference companion on World Wide Web - WWW '12 Companion*, page 687, New York, New York, USA, April. ACM Press.
- [Kinsella et al.2011] Sheila Kinsella, Vanessa Murdock, and Neil O'Hare. 2011. "I'm eating a sandwich in Glasgow". In *Proceedings of the 3rd international workshop on Search and mining user-generated contents - SMUC '11*, page 61, New York, New York, USA, October. ACM Press.
- [Li et al.2011] Wen Li, Pavel Serdyukov, Arjen P. de Vries, Carsten Eickhoff, and Martha Larson. 2011. The where in the tweet. In *Proceedings of the 20th ACM international conference on Information and knowledge management - CIKM '11*, page 2473, New York, New York, USA, October. ACM Press.
- [Li et al.2012] Chenliang Li, Jianshu Weng, Qi He, Yuxia Yao, Anwitaman Datta, Aixin Sun, and Bu-Sung Lee. 2012. *TwNER*. In *Proceedings of the 35th international ACM SIGIR conference on Research and development in information retrieval - SIGIR '12*, page 721, New York, New York, USA, August. ACM Press.
- [Lingad et al.2013] John Lingad, Sarvnaz Karimi, and Jie Yin. 2013. Location extraction from disaster-related microblogs. *Proceedings of the 22Nd International Conference on World Wide Web Companion*, pages 1017–1020, May.
- [Mahmud et al.2012] Jalal Mahmud, Jeffrey Nichols, and Clemens Drews. 2012. Where Is This Tweet From? Inferring Home Locations of Twitter Users. In *ICWSM*.
- [Nand et al.2014] Parma Nand, Ramesh Lal, and Rivindu Perera. 2014. A HMM POS Tagger for Micro-Blogging Type Texts. In *Proceedings of the 13th Pacific Rim International Conference on Artificial Intelligence (PRICAI 2014)*.
- [Ritter et al.2011] Alan Ritter, Sam Clark, Mausam, and Oren Etzioni. 2011. Named entity recognition in tweets: an experimental study. pages 1524–1534, July.
- [Sankaranarayanan et al.2009] Jagan Sankaranarayanan, Hanan Samet, Benjamin E. Teitler, Michael D. Lieberman, and Jon Sperling. 2009. *TwitterStand*. In *Proceedings of the 17th ACM SIGSPATIAL International Conference on Advances in Geographic Information Systems - GIS '09*, page 42, New York, New York, USA, November. ACM Press.