

# Where has this Tweet Come From? Fast and Fine-Grained Geolocalisation of Non-Geotagged Tweets

Pavlos Paraskevopoulos ·  
Themis Palpanas

the date of receipt and acceptance should be inserted later

**Abstract** The rise in the use of social networks in the recent years has resulted in an abundance of information on different aspects of everyday social activities that is available online, with the most prominent and timely source of such information being Twitter. This has resulted in a proliferation of tools and applications that can help end-users and large-scale event organizers to better plan and manage their activities. In this process of analysis of the information originating from social networks, an important aspect is that of the geographic coordinates, i.e., geolocalisation, of the relevant information, which is necessary for several applications (e.g., on trending venues, traffic jams, etc.). Unfortunately, only a very small percentage of the twitter posts are geotagged, which significantly restricts the applicability and utility of such applications. In this work, we address this problem by proposing a framework for geolocating tweets that are not geotagged. Our solution is general, and estimates the location from which a post was generated by exploiting the similarities in the content between this post and a set of geotagged tweets, as well as their time-evolution characteristics. Contrary to previous approaches, our framework aims at providing accurate geolocation estimates at fine grain (i.e., within a city). The experimental evaluation with real data demonstrates the efficiency and effectiveness of our approach.

**Keywords:** geotag, geolocation, Twitter, social networks

---

P. Paraskevopoulos  
University of Trento  
Telecom Italia - SKIL  
E-mail: p.paraskevopoulos@unitn.it

T. Palpanas  
Paris Descartes University  
E-mail: themis@mi.parisdescartes.fr

## 1 Introduction

Several social networks have emerged during the last decade. Social networks, such as Twitter [3], Facebook [1] and Google+ [2], give users the opportunity to express themselves and report details about their everyday social activities. The combination of this behavior with the widespread use of mobile smartphones and tablets, led to a very interesting phenomenon, where the activities reported within social networks are happening in real time, with individual users adding reports from several different locations (not just from their homes, or workplaces).

The above observation means that we now have access to datasets containing important information for the better and more detailed understanding of social activities. To that effect, several studies [32], including applications [28, 22, 35, 14, 7, 11, 6, 33, 39] and techniques [36, 31, 25, 34] have been developed that analyze datasets created through the use of social networks, in order to provide benefits to end users, businesses, civil authorities and scientists alike [27].

Note that several of these applications depend on the knowledge of the user location at the time of the posting. For example, this knowledge is necessary for applications that target to characterize an urban landscape, or to optimize urban planning [14], to identify and report natural disasters, such as earthquakes [28, 11], and to monitor and track mobility and traffic [7]. Such applications, which represent an increasingly wide range of domains, are restricted to the use of geotagged data<sup>1</sup>, that is, posts in social networks containing the geographic coordinates of the user at the time of posting.

Evidently, the availability of geotagged data, determines not only the possibility to use such applications, but also their quality-performance characteristics: the more geotagged data posts are available, the better the quality of the results will be (more accurately: the higher the probability for being able to produce better quality results). Nevertheless, the availability of geotagged data is rather limited. In Twitter, which is the focus of our study, the number of geotagged tweets is a mere 1.5-3% of the total number of tweets [19, 23, 15]. As a result, the amount of useful data for these applications to analyze is small, which in turn limits the utility of the applications.

In this study, we address this problem by describing a method for geolocalising tweets that are non-geotagged. Even though previous works have recognized the importance and have studied this problem [9, 18] (for a comprehensive discussion of this problem refer to [15]), their goal was to produce a coarse-grained estimate of the location of a set of non-geotagged tweets (e.g., those originating from a single user). The algorithms they propose operate at the level of postal zipcodes, cities, and geographical areas larger than cities. In contrast, we study this problem at a much finer granularity, providing location estimates for *individual* tweets first at the level of cities, and then at the level

---

<sup>1</sup> For the rest of this paper, we will use the terms *geotagged* and *geolocalised* interchangeably.

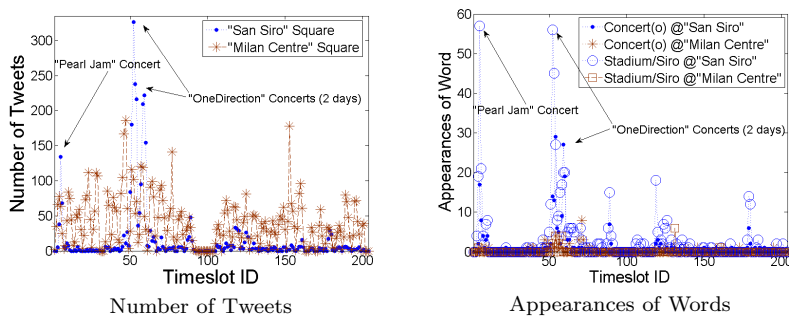


Fig. 1: Data generated from different neighborhoods (i.e., squares with side 1000 meters) in Milan (Italy), for time intervals of 4 hours, between June 20 and July 23, 2014.

of *city neighborhoods*, thus enabling a new range of applications that require detailed geolocalised data.

We illustrate and motivate some of these ideas in Figure 1. Figure 1(a) depicts the number of tweets posted from the neighborhood in which the “*SanSiroStadium*” is located, and from a neighborhood located in the center of the Milan (Italy), while Figure 1(b) shows the number of appearances of the keywords *concert* (in English and Italian) and *stadium/siro* in these neighborhoods. As these graphs show, the “San Siro” geolocation exhibits an unusually high activity during the time intervals that coincide with the concerts that took place in this stadium. Furthermore, during these concerts, the words *concert(o)* and *stadium/siro* originate from the “San Siro” geolocation much more frequently than a random geolocation in the city.

There are two main challenges that emerge when the granularity level becomes fine: first, to maintain high accuracy despite the wider range of possible locations available to the prediction algorithm; and second, to achieve high time performance despite the increased size of the search space of the algorithm. The framework we describe for the fine-grained geolocalisation of non-geotagged tweets is based on the careful evaluation of the similarities in the content between a new, non-geotagged tweet and a training set of geotagged tweets. The solutions we propose for this similarity evaluation make use of efficient-to-compute information retrieval and statistical measures, namely, Tf-Idf among the tweet contents, and correlation among the time series representing the volume of tweets in different candidate locations. The advantages of these measures are that they can effectively capture the most significant pieces of information needed to solve the problem, and that they have low time complexity.

The contributions we make in this paper can be summarized as follows.

- We describe and define the problem of fine-grained geolocalisation of non-geotagged tweets, which aims to operate on individual tweets, at the level of city-neighborhoods. We argue that the efficient solution of this prob-

lem will enable a multitude of applications that require detailed location information.

- We propose a framework for the solution of the above problem, which is based on the content similarities of tweets, as well as their time-evolution characteristics. The solution we describe is general, and essentially parameter free.
- Finally, we perform a detailed experimental evaluation of our approach, using real data from Twitter. The results demonstrate the efficiency and effectiveness of the proposed approach when compared to various alternatives.

The rest of the paper<sup>2</sup> is organized as follows. In Section 2 we present the related work. Section 3 formalizes the problem, and Section 4 describes our solution. We present our experimental evaluation in Section 5, and conclude in Section 6.

## 2 State Of The Art

Several works have studied the problem of geotagged tweet analysis. Balduini et al. [7] studied the movement of people by analyzing geotagged tweets. The authors analyzed tweets originating from London, and more precisely close to the Olympic stadium during the Olympic games. The results show that they could identify and track the movement of the crowd, especially during the opening ceremony. Some studies focus on the extraction of local events by analyzing the text in the tweets [12]. A recent study describes an approach of how to use social media data (including Twitter), in order to better understand and manage city-scale events, part of which involves the extraction of location information for tweets [6]. However, this is only done for tweets that are already geotagged, or tweets that mention the venue and/or event of interest, for a predefined set of venues and events. On the other hand, our framework is able to estimate the location of tweets regardless of the use of hashtags, or any other entity reference in the content, leading to a general solution that is effective even in cases where tweets referring to an unforeseen event (e.g., an accident), or tweets that do not explicitly mention the venue.

Abdelhaq et al. [4] use both geotagged and non-geotagged tweets for identifying keywords that best describe events. Then they keep only the geotagged tweets in order to extract the local events. Twitter posts have also been studied in order to identify the location of earthquakes [28], or fine-grained details on user activities (such as drinking alcohol) [16]. We note that in all the above studies, the tweets that are analyzed are already geotagged. In contrast, our focus is on non-geotagged tweets.

The identification of Points of Interest (POIs) with temporal awareness is the focus of a recent study [20]. The authors are analyzing tweets posted by Singaporean users, while using Foursquare check-ins referred in the tweets.

---

<sup>2</sup> This paper extends, and improves on our earlier results [26].

Another study proposed a framework that can automatically recognize POIs by correlating geotagged tweets with geotagged data deriving from Flickr [37]. The goal is to identify places, such as restaurants and hotels, that are not already part of databases such as LinkedGeoData, Geonames, Google Places, or Foursquare. The combination of data from Foursquare and the logs of executed applications on a smart-phone, has been used in order to predict the next location of the users [21].

The problem of using tweets in order to identify the location of a user, or the place that an event took place has been studied in the past. The “who, where, what, when” attributes extracted from a user’s profile can be used to create spatio-temporal profiles of users, and ultimately lead to identification of mobility patterns [38]. Cheng et al. [10] create location profiles based on idiomatic keywords and unique phrases mentioned in the tweets of users who have declared those locations as their origins.

The similarity between user profiles and location profiles has also been used in [9]. In this approach, they create user profiles for the active users, and extract the keywords that are characteristic of specific locations (i.e., they usually appear in some location, and not in the rest of locations). For the extraction of these keywords they initially assign weights using the Geometric-Localness (GL) method, and then prune them using a predefined keyword-weight threshold. This leads to a set of representative keywords for each location, which allows the algorithm to compute the probability that a given user comes from that location. A recent study evaluates the GL method, and compares it to other methods that solve the same problem. The experimental evaluation shows that the GL method achieves the best results [15].

Two studies that target to geotag tweets are presented in [13] and [24]. These two methods create chains of words that represent a location by using Latent Dirichlet Allocation (LDA) [8]. The latter study takes in addition into consideration the location a user has recorded as their home location. A study that predicts both a user’s location and the place a tweet was generated from is presented in [18]. In this study, the authors construct language models by using Bayesian inversion, achieving good results for the country and state level identification tasks. Finally, [30] presented a method for identifying the geolocation of photos by using the textual annotations of these photos.

Even though some of these studies are closely related to our work (e.g., [9, 18], which we further discuss in the experimental section evaluation section), we observe that they operate at a very different time and space scale. The profiles they create involve the tweets generated over a long period of time (up to several months), and the location that has to be estimated is the location of origin of the user, rather than the location from where a particular tweet was posted. Moreover, the space granularity used in these studies ranges from postal zipcodes to areas larger than a city. On the contrary, in our work we predict the location of individual tweets, at the level of city neighborhoods.

Two studies that target to geotag unique tweets are presented in [17] and [29]. The first method trains a model using past messages associated to locations, by extracting keywords that are connected to this location. In the later

**Algorithm 1** Tweet Geotagging Algorithm

---

**INPUT:** A training set of timestamped and geotagged tweets, a timestamped query-tweet ( $Q_t$ ) that is not geotagged.

**OUTPUT:** The most eligible candidate location.

---

```

1: for all  $i \in \{\text{candidate geolocations: Geolocs}\}$  do           ▷ process training dataset, for all
   locations
2:   for all  $t \in \{\text{time intervals}\}$  do                       ▷ and for all time intervals
3:      $Doc_{i_t} \leftarrow$  all tweets in location  $i$  at time interval  $t$ 
4:      $kwVector_{i_t} \leftarrow$  create vector of  $Doc_{i_t}$  keywords and their weights
5:    $kwVector_{Q_t} \leftarrow$  create vector of  $Q_t$  keywords and their weights           ▷ process
   non-geotagged tweet  $Q_t$ 
6:    $location \leftarrow \operatorname{argmax}_{i \in Geolocs} \{\text{similarity between } kwVector_{i_t} \text{ and } kwVector_{Q_t}\}$    ▷
   identify location of tweet  $Q_t$ 
7: return  $location$ 

```

---

study the authors develop a multi-indicator approach that combines information from the user’s profile and the tweets’s message for estimating both the location of a unique tweet and a user’s residence location. The main difference to our approach is that these methods rely on users that post many tweets in a time interval  $t$ , or on data from the user’s profile. In contrast, we target to geotag tweets even from users that have never posted before, or do not provide any profile data (such as their home location).

A recent survey presents methods relevant to location inference [5].

### 3 Problem Formulation

The problem we want to solve in this work is the estimation of the geographic location of individual, non-geotagged posts in social networks.

**Problem 1:** Given a set of geotagged posts  $P_{t_j}^{l_1}, \dots, P_{t_j}^{l_i}, t_1 \leq t_j \leq t_2$ , where  $l_i$  is the location the post was generated from and  $t_j$  is the time interval during which the post was generated at, and a non-geotagged post  $Q_{t_q}, t_1 \leq t_q \leq t_2$ , we wish to identify the location  $l$  from which  $Q$  was generated.

The timestamps  $t_1$  and  $t_2$  represent the start and end times, respectively, of the time interval we are interested in.

In the context of this work, we concentrate on fine-grained location predictions: we wish to estimate the location of a post at the level of a city neighborhood (which is usually much smaller than a postal zipcode). Furthermore, we focus on twitter posts, whose particular characteristics are the very small size (i.e., up to 140 characters long), and the heavy use of abbreviations and jargon language.

### 4 Proposed Approach

In this section, we describe our solution to the problem of fine-grained geolocalisation of non-geotagged tweets.

We provide a high level description of our approach in Algorithm 1. Our method is based on the creation of vectors describing the Twitter activity in terms of important keywords for each geolocation we have data from, and for the period of time we are interested in. The geolocations correspond to fine-grained spatial regions (in our study, they are squares with side length of 1000 meters). The time intervals correspond to brief time segments, during which posts on the same, or related topics may be observed (in our study, they are 4 hour intervals). The vectors represent the weights of each keyword, and are stored in *kwVector* for each geolocation and time interval. There are several ways to compute these weights: we consider the number of appearances of a keyword in a given geolocation, and the significance of a keyword, measured using Tf-Idf, for a given geolocation and the entire dataset.

In order to identify the geolocation for a non-geotagged tweet,  $Q$ , we compute the similarity between the vector of  $Q$  and the vector of each candidate geolocation. When calculating this similarity, we can additionally take into account the correlation between the local and the global activity time series, i.e., the evolution over time of the number of tweets in a given geolocation and all the geolocations, respectively. Finally, the algorithm returns the geolocation with the highest similarity value.

In the following sections, we elaborate on the methods discussed above.

#### 4.1 Grouping the Posts and Extracting Important Keywords

We start by processing the training set of geotagged posts. We group these posts according to the geolocation that they were generated from, and the time interval they belong to. After this grouping step, we calculate the concordance of the keywords in each group: the dictionary containing the number of appearances of each keyword in a geolocation. At the end, we have for each geolocation and time interval a vector of the important keywords, along with the corresponding weights. We call the algorithm that uses this method for generating the keyword vectors *TG (Tweet Geotagging)*.

We observe that concordance is a simple measure that only accounts for the frequencies of keywords, but fails to take into account their relative significance. Therefore, we also employ the Tf-Idf model:  $idf_{keyword} = \log(\frac{n}{k})$ , where  $n$  is the number of documents,  $k$  is the number of documents that keyword appears in, and  $tfidf_{i,keyword} = \frac{count}{l} * idf_{keyword}$ , where  $l$  is the total number of keywords in document  $i$ . Using Tf-Idf, we can calculate the significance of each keyword in our training dataset (according to the former equation above), and set the weight for a keyword in some geolocation, depending on the number of its appearances at this geolocation (according to the latter equation). This method leads to high weights for the keywords that appear at a small number of geolocations. As a final step, we sort the keywords according to their weight and prune the keywords with low weights, and therefore, only keep the significant keywords for each geolocation, which correspond to the keywords that best characterize the activity of the given geolocation at a particular time in-

**Algorithm 2** Similarity Calculation and Probability Extraction

---

```

1: procedure VECTORSIM(vectors for  $Q_t$  and candidate geolocations)
2:    $mag_{Q_t} \leftarrow \sqrt{\sum_{\forall j \in kwVector_{Q_t}} kwVector_{Q_t}[j]^2}$     $\triangleright$  Extract the magnitude of the
   Q-tweet kwVector ( $Q_t$ )
3:   for all  $i \in Geolocs$  do                                        $\triangleright$  For every Candidate Location (CL)
4:      $mag_{i_t} \leftarrow \sqrt{\sum_{\forall j \in kwVector_{i_t}} kwVector_{i_t}[j]^2}$     $\triangleright$  extract the magnitude of its
   kwVector
5:      $Sim_{i_t, Q_t} \leftarrow \frac{\sum_j kwVector_{i_t}[j] * kwVector_{Q_t}[j]}{mag_{i_t} * mag_{Q_t}}, \forall j \in kwVector_{Q_t} \cap kwVector_{i_t}$     $\triangleright$ 
   Calculate the similarity between  $Q_t$  and CL
6:     for all  $i \in Geolocs$  do                                        $\triangleright$  Get the probability distribution
7:        $Prob_{i_t, Q_t} \leftarrow \frac{Sim_{i_t, Q_t}}{\sum Sim_{Q_t}}$ 
8:     SortDescending  $Prob_{i_t, Q_t}$ 
9: return  $i$  with highest  $Prob_{i_t, Q_t}$ 

```

---

terval. We call the algorithm that uses this method for generating the keyword vectors *TG-TI (Tweet Geotagging Tf-Idf)*.

In order to create the keyword vector for the non-geotagged tweet,  $Q$ , we wish to geolocalise, we follow the same process as before.

#### 4.2 Similarity Calculation and Best Match Extraction

Our next target is to calculate the similarity between the keyword vector of  $Q$  and the keyword vector of each one of the candidate geolocations.

We follow the steps presented in Algorithm 2. The magnitude,  $mag$  is the Euclidean Norm, computed over all the keywords that appear in the vector. We calculate the magnitude of the  $Q$  vector,  $mag_{Q_t}$ , and of each one of the candidate geolocations  $i$ ,  $mag_{i_t}$ , for a given time interval  $t$ . We denote with  $kwVector[j]$  the weight of the  $j$ -th term of the vector. The similarity is computed using the formula shown in line 5 (over all the keywords that appear in both the vector  $Q$  and the vector of the geolocation  $i$ ). The algorithm stores in a sorted list the similarity values for each candidate geolocation. It then normalizes these values over the sum of all similarities, giving us the probability that each candidate geolocation produced  $Q$ . Transforming these values into a probability distribution gives us more flexibility: for example, as we discuss next, we can readily combine this similarity measure with similarities computed using other methods. Furthermore, we can use the probability values in order to produce geolocation predictions only in the cases where we are confident (i.e., these probabilities are high). At the end, the algorithm returns the geolocation(s) with the highest probability(ies).

In our approach, this similarity calculation happens in two phases (using for both the same general method presented above). First, we determine the city with the highest probability for having generated  $Q$ , and then the neighborhood (i.e., square with side 1000 meters) within that city, with the highest



probability. This solution has the added benefit that it can be effective even in the presence of small training datasets (i.e., few geotagged posts), which would not normally be adequate to directly train models with a very large number of candidate geolocations, as in our problem.

Therefore, when we change granularity, from the city to the neighborhood level, we add one more step to our method:  $Prob_{i_t, Q_t} = Prob_{city_{j_t}, Q_t} * Prob_{j_t, Q_t}$ , where  $city_{j_t}$  ranges over all the candidate cities, and  $j_t$  ranges over all the candidate neighborhoods within a city,  $city_{j_t}$ . The probability that a candidate neighborhood in a specific city ( $Prob_{i_t, Q_t}$ ) is the correct geolocation for  $Q_t$  is computed by multiplying the probability that a candidate city is the correct one ( $Prob_{city_{j_t}, Q_t}$ ) by the probability that a given neighborhood within that city is the correct location ( $Prob_{j_t, Q_t}$ ).

### 4.3 Similarity Based on Correlation of Activity Time Series

The similarity measure discussed earlier is based entirely on the contents of the relevant posts, but ignores other useful characteristics of the data. In what follows, we describe a method that exploits the time-evolution behavior in order to derive an additional similarity measure.

This method is based on the activity time series, which record the number of posts generated by a given geolocation over time. We call these series *local activity* time series. We also compute the *global activity* time series, where we record the sum of the number of posts for all geolocations over time. The similarity is then expressed as the correlation value between the local activity of a candidate geolocation with the global activity. The intuition is that posts about an important event will significantly change the local activity and influence in the same way the global activity.

A straightforward idea is to compute the Pearson’s correlation between the local and the global activity. Since we are only interested in similar behavior between local and global activity, we can only keep the positive correlations, and then normalize them over the sum of all correlation values to produce a probability distribution. More specifically, we can construct the global activity time series,  $Gts_i$ , for a coarse-grain geolocation (e.g., a city  $i$ ), as well as the local activity time series,  $Lts_j$ , for all the fine-grain geolocations within the coarse-grain one (e.g., the neighborhoods  $j$  inside the city  $i$ ). Finally, we can compute the correlation between these time series using the Pearson’s correlation.

The above idea proved to be somewhat useful, but with limited benefits (refer to algorithms *TG-C* and *TG-TI-C*, described in [26]). The reason is that this method employs the correlation measure irrespective of the trend exhibited by the local and global activities. For example, these activities can be positively correlated, but have a negative trend (i.e., activity is diminishing). Evidently, in such cases the correlation does not help, and should not be taken into account.

We now describe a new technique that addresses this problem. More specifically, we consider a location as a candidate location only if both the local and the global activity increase. As we demonstrate later, this modification on the usage of the correlation measure leads to a significantly better result.

This new correlation-based technique is shown in Algorithm 3.

Initially, we construct the global activity time series,  $Gts_i$ , for a coarse-grain geolocation (e.g., a city  $g$ ), as well as the local activity time series,  $Lts_{CL}$ , for all the fine-grain geolocations within the coarse-grain one (e.g., the neighborhoods  $j$  inside the city  $i$ ). Since we are only interested in similar behavior between local and global activities only in the case where we have an increasing trend, we use the linear regression line in order to test this trend. In particular, we use the  $\lambda$  parameter of the equation representing the linear regression line,  $y = \lambda * x + b$  (refer to line 4). If  $\lambda$  is positive, we assume that the time-series has a positive slope (lines 5 – 6 and 12 – 13). In this process, we use smaller sliding sub-windows of size  $n/2$  (lines 3 and 10), sliding them across the original window. As a result, we have a sub-window that slides  $n/2$  times on the original  $n$ -timeslot window, counting the number of the slides that result into positive linear regressions for both time-series describing the local and the global activity.

After having calculated all the  $\lambda$  for each candidate locations, we calculate the Pearson correlation between the time-series describing the local and the global activity, and add 1 to this value, in order to shift the range of values between  $[0,2]$  (line 14). This has the desirable effect that we avoid negative similarities (that would result from negative correlations). Note that candidate locations that correspond to positive correlation receive a bonus (they get multiplied by a number in the range  $(1,2]$ ), while those that correspond to a negative correlation get penalized (they get multiplied by a number in  $[0,1)$ ). Finally, we set a threshold  $th_{LR}$ , and we check if the number of the sliding windows for each location that have positive  $\lambda$  is greater than  $th_{LR}$  (line 7 and 16).

If the number of the sliding sub-windows that have positive  $\lambda$  exceeds  $th_{LR}$ , then this location is considered as a candidate location, and we assign to the location its correlation and the value *True* for exceeding the threshold (line 17), otherwise we assign to the location its correlation and the value *False* (line 19). Finally, the algorithm returns the final set of Candidate Locations, CL, which includes for each location its correlation value and the attribute that indicates if the location exceeds the  $th_{LR}$  threshold (line 22).

We can then combine this method with the TG algorithm, by multiplying the two similarity measures (concordance similarity and correlation), to obtain the *TG-CLR* (*Tweet Geotagging with activity Correlation with Linear Regression*). When we do the same with the TG-TI algorithm, we get the *TG-TI-CLR* (*Tweet Geotagging with Tf-Idf and activity Correlation with Linear regression*) algorithm. If the candidate location that has the greatest similarity with the non-geotagged tweet  $Q$  does not exceed the  $th_{LR}$  threshold (i.e., it has been assigned the value *False*), then we do not match  $Q$  to any location.

**Algorithm 3** Activity Correlation

---

```

1: procedure CORRELATIONSIM(global  $Gts$  and local  $Lts_{CL}$  activity time series, threshold
    $th_{LR}$ , Candidate Locations  $CL$ , Window time intervals  $[t_1, t_2]$ )
2:    $counter_g \leftarrow 0$ 
3:   for all  $subWindow_i \in Window$  do ▷ For how many subWindows Global
time-series have positive  $\lambda$ 
4:      $\lambda_g \leftarrow \frac{\Sigma((x-\bar{x})(y-\bar{y}))}{\Sigma(x-\bar{x})^2 \Sigma(y-\bar{y})^2}$ 
5:     if  $\lambda_g \geq 0$  then
6:        $counter_g \leftarrow counter_g + 1$ 
7:   if  $counter_g > th_{LR}$  then ▷ If Global time-series have at least  $th_{LR}$  subWindows
with positive  $\lambda$ 
8:     for all  $loc \in CL$  do ▷ check Local time-series of all Candidate Locations ( $loc$ )
9:        $counter_{loc} \leftarrow 0$ 
10:      for all  $subWindow_i \in Window$  do
11:         $\lambda_{loc_t} \leftarrow \frac{\Sigma((x-\bar{x})(y-\bar{y}))}{\Sigma(x-\bar{x})^2 \Sigma(y-\bar{y})^2}$ 
12:        if  $\lambda_{loc_t} \geq 0$  then
13:           $counter_{loc} \leftarrow counter_{loc} + 1$ 
14:       $corr_{g,loc} \leftarrow \frac{\Sigma_{t=t_1}^{t_2} (Gts_{gt} - \bar{Gts}_g)(Lts_{loc_t} - \bar{Lts}_{loc})}{\sqrt{\Sigma_{t=t_1}^{t_2} (Gts_{gt} - \bar{Gts}_g)^2 \Sigma_{t=t_1}^{t_2} (Lts_{loc_t} - \bar{Lts}_{loc})^2}} + 1$  ▷ Calculate the
correlation between
15:        ▷ the time-series of the  $loc$  and the global time-series
16:        if  $counter_{loc} > th_{LR}$  then ▷ Check if  $loc$  exceeds the  $th_{LR}$ 
17:          Assign to  $loc$  its correlation value, and True (for exceeding the  $th_{LR}$ 
threshold)
18:        else
19:          Assign to  $loc$  its correlation value, and False (for not exceeding the  $th_{LR}$ 
threshold)
20:        else
21:          Assign to all  $loc$  in  $CL$  the values 0 (for the correlation) and False
22: return  $CL$ 

```

---

## 4.4 Sliding Windows

We observe that previous methods use all past data in order to build their models. Methods such as [18] and [9] start building their models taking into consideration all available data. However, this may lead to situations where some local events may be mishandled. For example, consider the case, where a concert takes place in a city, followed by a second concert the following day. Then, a model that is based on all the data (and in the absence of specific and detailed keywords) is likely to assign the tweets relevant to the second concert to the location of the first concert, for which more data are available.

In order to avoid similar problems, we can use a tumbling window model [26]. Although this helps to address the problem mentioned above, tumbling windows may still mis-assign tweets that are generated at the beginning, or at the end of the window, and are connected to an event that is outside the window period.

A better idea is to use sliding windows, which we exploit in this work. In this case, a particular timeslot can be part of  $n-1$  windows, where  $n$  is the length of the window. If a timeslot is at the beginning of an event (the latest

in the window), the new timeslots to be inserted later are going to be more relevant. As a result, the timeslot is going to be in  $n-1$  windows, the majority of which will be relevant.

Using the sliding window idea, we can now take advantage of the already extracted models of each location and incrementally update them for every slide, reducing dramatically the time needed for the contraction of the keyword vectors. In order to achieve this, we do not recalculate the concordance of each word for each location across the window. Instead, we extract the concordance across the window only for the first model created, and for every slide we update the concordances of each word by subtracting the concordance of the words in the data removed and adding those in the data added to our dataset. We can see the steps of the incremental update of the vectors in Algorithm 4.

Furthermore, due to the incremental update that we achieve at our concordance *kwVectors*, we prove that our method can be applied in streaming manner. Unfortunately, the incremental update is not straight applicable on the Tf-Idf *kwVectors* but still Tf-Idf methods get advantage on the incremental update of the concordances.

---

#### Algorithm 4 Incremental Update of kwVector

---

```

1: procedure UPDATE OF KWVECTOR(all kwVectorst, geotagged tweets from location i for
   time intervals  $t - 1$  and  $t + 1$ )
2:   for all kwVectorit ∈ {kwVectorst} do
3:     for all word ∈ {kwVectorit} do
4:       concit-1 ← concordance in i at  $t - 1$ 
5:       concit+1 ← concordance in i at  $t + 1$ 
6:       conci ← concit − concit-1 + concit+1
7: return kwVectorst

```

---

## 5 Experimental Evaluation

**Experimental Setup.** We performed the experiments on a server running on Ubuntu 14.04.2 LTS, with 64GB RAM, and an Intel(R) Xeon(R) CPU E5506 @ 2.13GHz processor. For the implementation of our methods and the reimplementations of the QL and KL we used Python 2.7.

**Datasets.** For the evaluation of our approach, we use 3 datasets containing geotagged<sup>3</sup> posts from Twitter, generated in Italy, Germany and the Netherlands. In particular, we have data from 6 of the largest Italian cities, namely, Rome, Milan, Naples, Bologna, Venice and Turin, and from the capital of Germany, Berlin, and the capital of Netherlands, Amsterdam. The tweets from Italy were generated between June 20 and July 23, 2014, while the tweets from Germany and the Netherlands were generated between August 10 and

---

<sup>3</sup> Earlier studies have shown that techniques and models built for geotagged data indeed generalize to non-geotagged data, since geotagged and non-geotagged tweets have similar data characteristics [15].

September 11, 2014. The granularity of the neighborhood level we use for every city is a square with side of 1000 meters. The number of tweets is 543.295 for Italy (219.681 originated from Rome, 137.622 from Milan, 60.065 from Naples, 49.434 from Bologna, 46.982 from Turin, and 29.511 from Venice), 77.179 for Berlin and 136.189 for Amsterdam. The time windows we use have a duration of 4 hours (which can effectively capture an important event, as well as the start and the aftermath of this event), while also keeping the detailed aggregated information for every 15min time interval. As mentioned in Section 4.4, we use the sliding window model. We experimented sliding the window by 1 and by 2 time intervals, getting almost the same results; thus, we chose to slide our window by 2 time intervals per slide (30-minutes), which led to faster execution times. Finally, the default grid we use in this study is 20 by 20 squares.

**Algorithms.** We experimentally evaluate the six algorithms we described in Section 4, namely, TG, TG-TI, TG-C, TG-TI-C, TG-CLR and TG-TI-CLR (the last two only for the neighborhood level). As baselines, we implemented the QL and KL methods [18], which aim to solve a similar problem. In order to choose the value for the  $\mu$  parameter, we followed the same methodology as in the original paper [18]: we experimented with several values for the  $\mu$  parameter, in the range [100,10000]), and verified that  $\mu = 10000$  gave the best results in our setting, as well.

**Evaluation Measures.** We study the time performance, as well as the effectiveness of each approach using the precision and recall measures:  $Precision = \frac{cgTweets}{gTweets}$  and  $Recall = \frac{cgTweets}{aTweets}$ , where  $cgTweets$  is the number of the correctly geolocated tweets,  $gTweets$  is the number of tweets we geolocated, and  $aTweets$  is the number of all tweets in the test set. In the case of the city level, where we predict the geolocation for all the tweets in the test set, the above precision and recall measures coincide, and we use the term *accuracy* instead. For the neighborhood level though, we do not predict the geolocation of tweets for which all our candidate locations have a similarity of 0. Thus, we report results for both precision and recall. We also report the balanced F1 measure,  $F1 = 2 * \frac{Precision * Recall}{Precision + Recall}$ . Following previous work [18], we report the results when we consider the top-1 (@Top1), top-3 (@Top3), and top-5 (@Top5; only for neighborhood level) predicted geolocations, as well as the results when considering as correct the prediction of the exact geolocation (@0-Step), or of any geolocation at distance 1 (@1-Step; exact and its eight immediate neighbors), or 2 (@2-Step; exact and its 24 closest neighbors) from the exact. In all our experiments, we randomly divided the dataset in 80%training and 20% testing, repeated each experiment 30 times, and reported the mean values in the results.

## 5.1 City-Level Results

We start our analysis by running our method on city-level. We extract the geotagged tweets from the 6 cities, removing the duplicated posts in order to

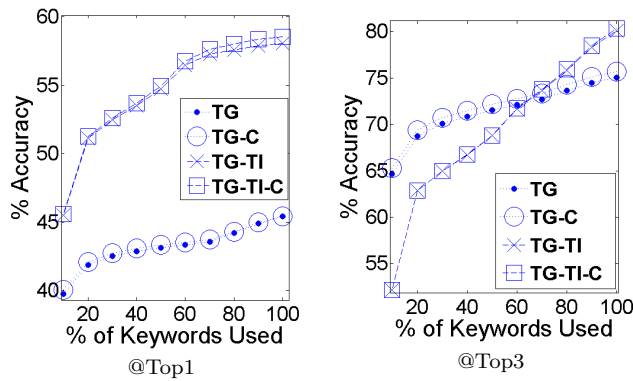


Fig. 2: Accuracy for city level when using TG, TG-C, TG-TI, and TG-TI-C (@0-Step).

avoid spam. We record the activity every 15 minutes, and we consider time intervals of 4 hours, leading to 181 timeslots (due to technical problems some of the timeslots were empty, and we excluded those from our analysis).

In this case we extracted the similarities between the test tweets and the 6 cities, and we also evaluated our approach using the correlation of the activity time series (refer to Section 4.3): we use the Pearson’s correlation between the activity time-series of the 6 cities and the activity time series of Italy. The results (@Top1 and @0-Step) are presented in Figure 2(a). As we can see in this plot, the accuracy for the city level is increased compared to the accuracy before the correlation. More precisely, we get the maximum number of matches in all four cases when we keep 100% of the keywords. The accuracy of TG and TG-C is almost identical, at 45%. For TG-TI we get 58% accuracy, while when using TG-TI-C we get 59% (though, our t-test analysis revealed that this difference is not statistically significant). After further analyzing the results of these two algorithms, we found that for 134 windows TG-TI-C has better accuracy, for 4 windows TG-TI and TG-TI-C have the same accuracy, and for the rest 43 windows TG-TI performs better. We note that the accuracy of the random algorithm is 17%.

After evaluating the algorithms using the most similar candidate geolocation (@Top1), we also evaluated them using the 3 most similar candidates (@Top3). As we can see in Figure 2(b), when using only a small percentage of the keywords we get better results with the TG and TG-C algorithms. In contrast, when we use more than 70% of the keywords, the Tf-Idf based algorithms, TG-TI and TG-TI-C, result in better accuracy. The accuracy is increasing when the percentage of the keywords used increases.

## 5.2 Neighbourhood-Level Results

In this subsection, we present the results for the neighborhood level evaluation, for which we used data from four different European cities: Milan, Rome,

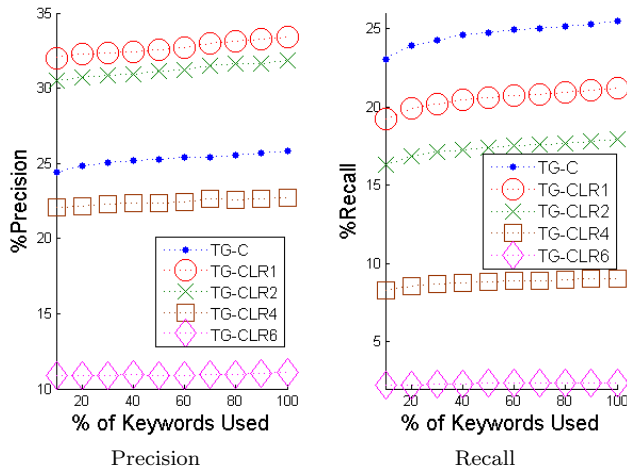


Fig. 3: TG-CLR for Different LR Parameters (@Top1 and @0-Step).

Berlin and Amsterdam. As we have already mentioned, we created a grid of 400 squares (20 by 20) for each city. For the city of Rome, we additionally ran some experiments using a grid of 900 squares (30 by 30).

### Setting the Parameters

We first identify the best threshold to use for the LR parameter. As we mentioned at the beginning of this section, we use a window of 16 timeslots, and sub-windows of size  $n_{sub-window} = n_{window}/2 = 8$  (refer to Section 4.3). Furthermore, the maximum LR equals to the number of slides, which is 8, as well. We experimented by setting the LR-threshold equal to  $\{1, 2, 4, 6\}$ , and depict the results in Figures 3 (precision and recall for algorithms not using Tf-Idf), 4 (precision and recall for algorithms using Tf-Idf), and 5 (F1 measure for all algorithms). For brevity, we only report the results for the city of Milan; results for the other cities are similar.

In this experiment, we had 3264 15-min timeslots, resulting into 1624 window slides. For each method, we extracted the mean precision, recall and F1 scores among all windows, while varying the percentage of the keywords used. We observe that the best mean precision is 48%, which is achieved by TG-TI-CLR1 when using 100% of the keywords (Figure 4(a)), while the maximum recall for this method is 32%, when using 40% of the keywords (Figure 4(b)). Note that the same method without the use of the trends, that is, TG-TI-C, has maximum precision and recall 40% and 39%, respectively. Regarding the TG-CLR1 algorithm, we get maximum precision 33% and maximum recall 21%, both when using 100% of the keywords. Due to these, and after finding out that the F1 score of the *CLR1* methods isn't too different compared to those not using linear regression, we concluded that for the rest of the experimental part we are going to use only the *CLR1* methods.

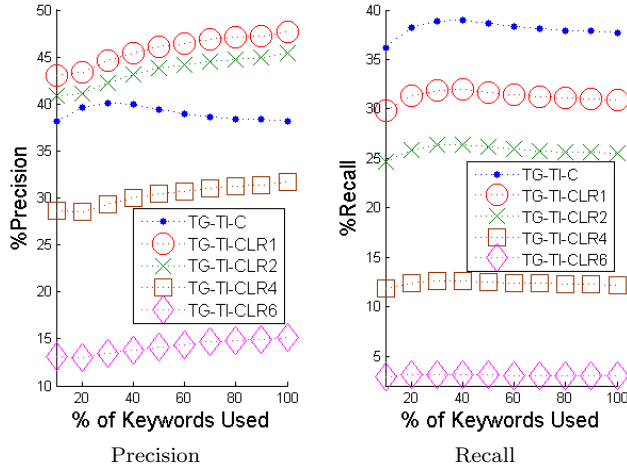


Fig. 4: TG-TI-CLR for Different LR Parameters (@Top1 and @0-Step).

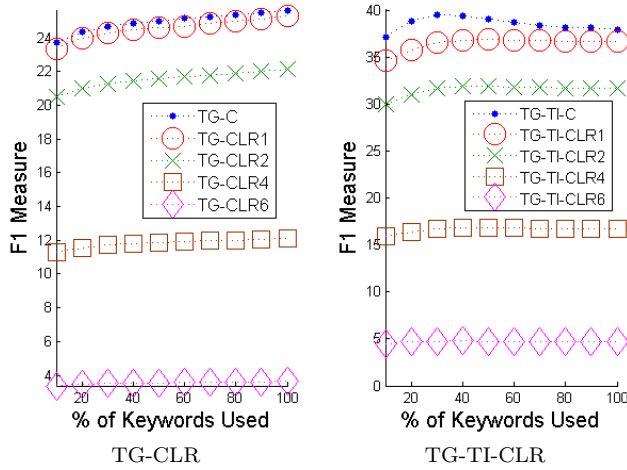


Fig. 5: F1 for TG-CLR and TG-TI-CLR for Different LR Parameters (@Top1 and @0-Step).

### Evaluating the Correlation-Based Methods

In the following experiments, we compare the CLR methods to those that do not use correlation. In Figure 6, we present the mean precision and recall that our algorithms have for the city of Milan among all windows, when varying the percentage of the keywords used. As before, we only consider the first answer given by each algorithm (i.e., @Top1). The best precision is 48% and is achieved by TG-TI-CLR1 using 100% of the keywords. The maximum recall is 38% achieved by TG-TI when using 30% of the keywords. According to the F1 measure, TG-TI achieves its best using 30% of the keywords, with F1



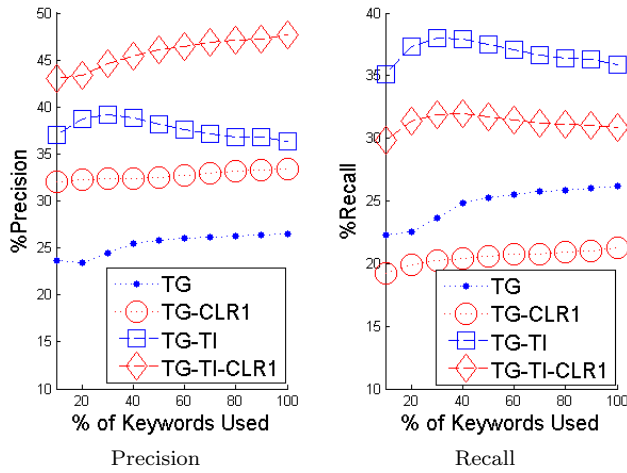


Fig. 6: Trade-off Between Precision and Recall for Neighbourhood Level (Milan, @Top1 and @0-Step).

equal to 39%. TG-TI-CLR1 achieves best F1 score 37% when using 50% of the keywords. The second best precision is 39%, achieved by TG-TI when using 30%. The best precision achieved by TG is 27%, while its best recall is 26%. TG-CLR1 reached up to 33% precision and 21% recall (both achieved when using 100% of the keywords). The mean accuracy of the random algorithm, which was choosing one square at random only among that had data at the train datasets, was less than 2%.

We note that the best precision is always observed when we use the TG-TI-CLR1 algorithm. This means that the correlation between the city and square activities is beneficial, when using the linear regression parameter that prunes activities with negative trends. As a result, we do not estimate the location of tweets that would probably be wrongly predicted, leading to a small penalty in recall, but increased precision.

We now report the results of the same experiment for the cities of Rome (in Figure 7), Berlin (in Figure 8), and Amsterdam (in Figure 9). The best precision we observed for the city of Rome was 48% and was achieved by TG-TI-CLR1, using 100% of keywords, while the best recall was achieved when using TG-TI method, using 40% of keywords. The same methods also resulted in the highest precision and recall for Berlin. In particular, TG-TI-CLR1 achieved a precision of 58%, for a recall of 40%. The best recall for Berlin was 47%, achieved by TG-TI, which also led to the second best precision, 51%. Regarding the city of Amsterdam, we achieve the highest precision of 44% with TG-TI-CLR1, while the best recall of 38% is achieved by TG-TI.

The results show that the behavior of the algorithms is similar across cities, while their relative performance remains the same. An interesting observation is the fact that the precision and recall for Berlin are much higher than the

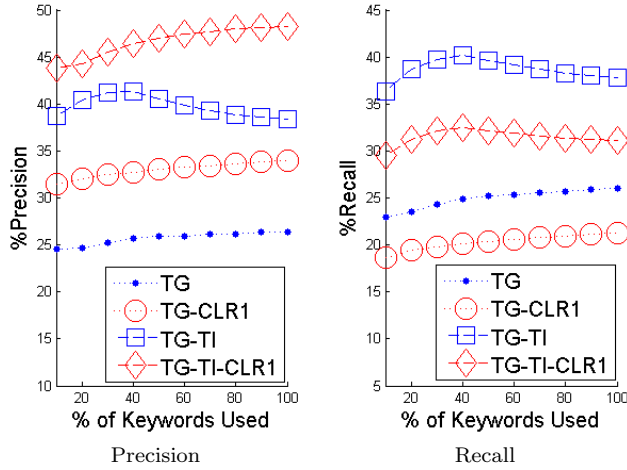


Fig. 7: Precision and Recall for the City of Rome (@Top1 and @0-Step).

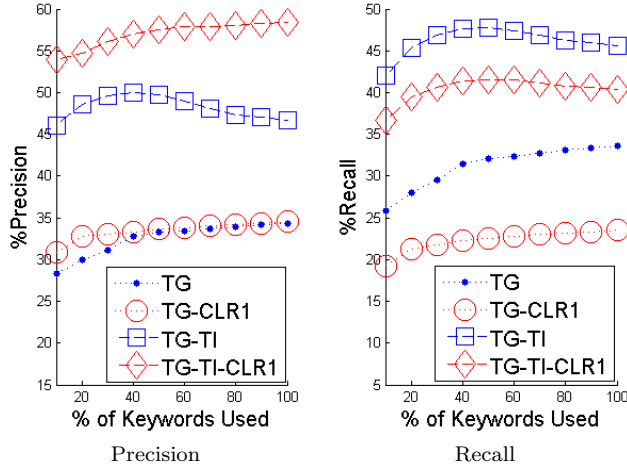


Fig. 8: Precision and Recall for the City of Berlin (@Top1 and @0-Step).

rest of the cities. This is due to the distribution of the keywords among the squares, which resulted into more representative keyword sets for each square.

### Comparing to Baselines

In this set of experiments, we compare our approach to the QL and KL baseline algorithms. We use the same spatial and temporal granularities for all algorithms. Similarly to our methods, we only consider tweets for which there exists at least one candidate location with similarity greater than 0. The results of this comparison are illustrated in Figure 10(a).

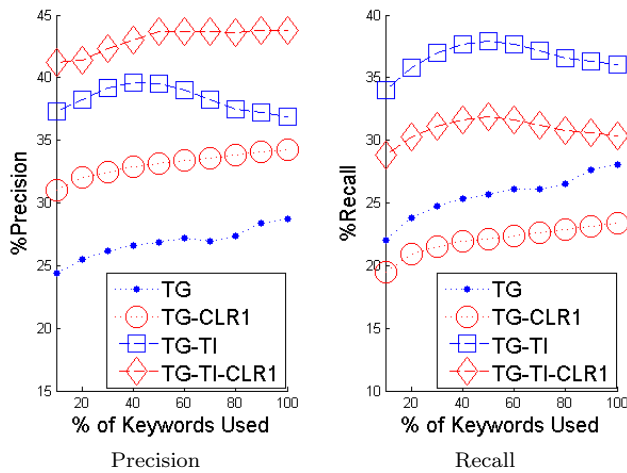


Fig. 9: Precision and Recall for the City of Amsterdam (@Top1 and @0-Step).

We observe that TG-TI-CLR1 achieves up to 18% better recall than the QL algorithm<sup>4</sup>, and up to 22% better F1 score. This difference in performance can be explained by the different focus of the QL algorithm, which was developed to operate at much bigger spatial (in the order of zipcodes, or cities) and temporal granularities (in the order of weeks, or months) [18]. We also note that (for the same reasons) the results between QL and KL are almost the same. Therefore, in our plots we only report the F1 score for QL.

In terms of time performance, we measured the mean execution time needed per 4-hour window for the entire process: training the models, and extracting the similarities between the query-tweets and the candidate locations. Figure 10(b) depicts the execution time needed for each algorithm.

As we can see in the graph, TG is the fastest algorithm. This is natural, since this algorithm does not spend time calculating the Tf-Idf, the correlations, or the linear regressions. The QL algorithm has a consistently high execution time of around 90sec, independent of the number of keywords considered. TG-TI-CLR1 performs in the middle. The interesting point is that although this algorithm has to calculate the Tf-Idf, the correlations and the linear regressions, the total time needed for each square, when using 10-70% of the keywords is smaller than the time needed for TG-CLR1. The reason is the search space pruning. When compared to TG-CLR1, the TG-TI-CLR1 algorithm prunes stopwords, and thus, eliminates the candidate locations that do not share any keyword with the tweet under examination. We also observe

<sup>4</sup> We note that the QL results reported here are much better than those reported in our earlier study [26]. This is due to the different experimental setup (i.e., sliding windows) that we now use for all algorithms, which resulted in an increased number of windows with a high number of tweets, leading to higher execution times and better models.

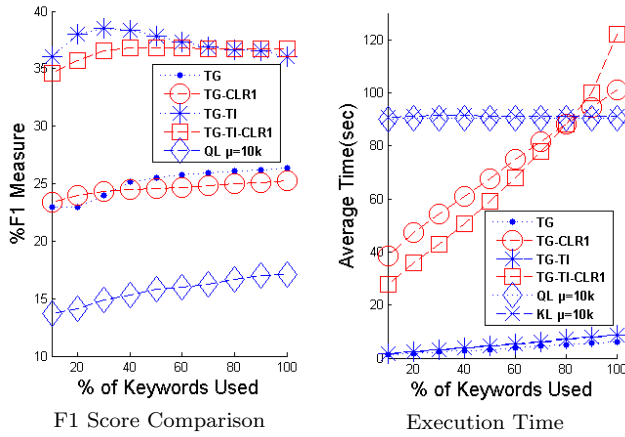


Fig. 10: Trade-off Between F1 Score and Execution-Time for the City of Milan (@Top1 and @0-Step).

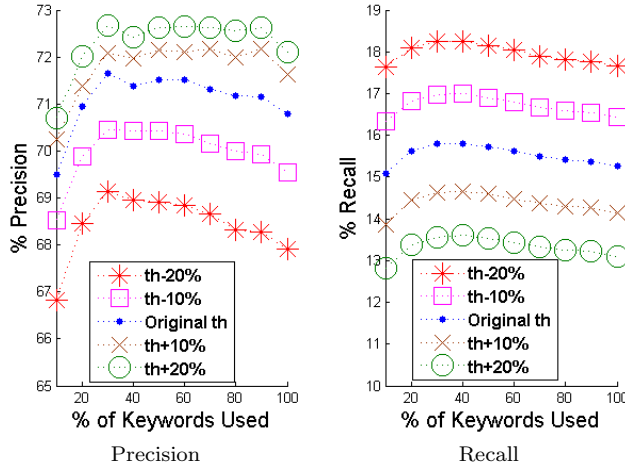


Fig. 11: Precision and Recall on Neighbourhood Level for TG-TI when using dynamic thresholds (Th) (@Top1 and @0-Step).

that when TG-TI-CLR1 achieves its best F1 score, i.e., when using 50% of the keywords, it is significantly faster than the QL algorithm.

Finally, we note that the KL algorithm performs very similar to QL, but requiring at all cases a bit higher time when compared to QL (around 0.8 secs more).

### Focusing on Precision

We now examine the behavior of our algorithms when we want to achieve high precision, which is useful for several applications.

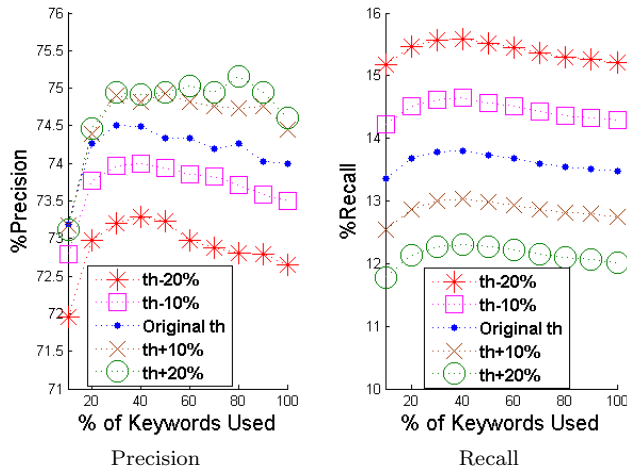


Fig. 12: Precision and Recall on Neighbourhood Level for TG-TI-CLR1 when using dynamic thresholds (Th) (@Top1 and @0-Step).

In the first set of experiments, we employ a dynamic similarity threshold that determines whether the algorithm will make a prediction for the geolocation. The thresholds we use are automatically set, based on the results of the same timeslots of the previous days: they are computed as the mean of the similarities of the correctly identified geolocations, averaged over the corresponding timeslots of the previous days. We have 48 (dynamic) thresholds, one per half-hour slide. Evidently, these thresholds lead to fewer predictions of tweet geolocations, reducing the recall, but increasing the precision.

In Figure 11, we present the precisions and the recalls after the introduction of the thresholds for the method TG-TI, while in Figure 12 we present the precision and recall for TG-TI-CLR1. We run experiments by using the exact dynamic threshold, the exact threshold  $\pm 10\%$  and the exact threshold  $\pm 20\%$ . Furthermore, in order to evaluate the results, we use again the balanced F1 score. The F1 score for the two methods presented before is depicted in Figure 13.

After evaluating our methods using the first most similar answer, we analyzed the results when taking under consideration the first 3 (Top3) and the first 5 (Top5) most similar candidates. In Figures 14 and 15, we can see depicted the mean precisions and recalls when using 10-100% of the keywords for both cases. The results show that both precision and recall are benefiting, with the F1 scores increasing from around 35% to around 55% (Figure 16).

Finally, we study the performance of our methods in the case where we relax the definition of the correct answer to include answers that are 1 square (1 - Step), or 2 squares (2 - Steps) away from the exact answer. That is, we consider the near neighbors of the exact answer to be correct answers, as well.

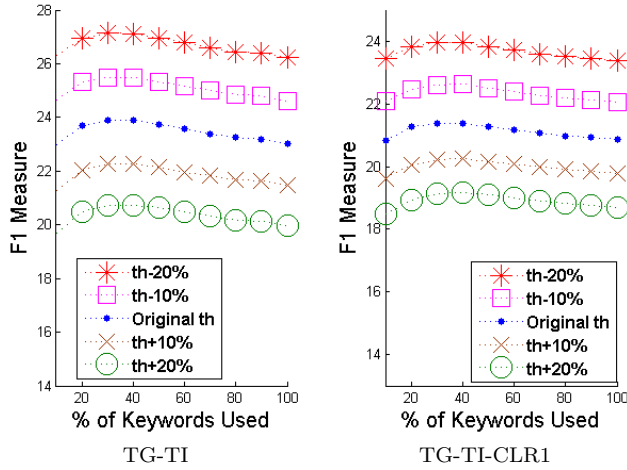


Fig. 13: F1 measure for Neighbourhood Level with threshold (Th) (@Top1 and @0-Step).

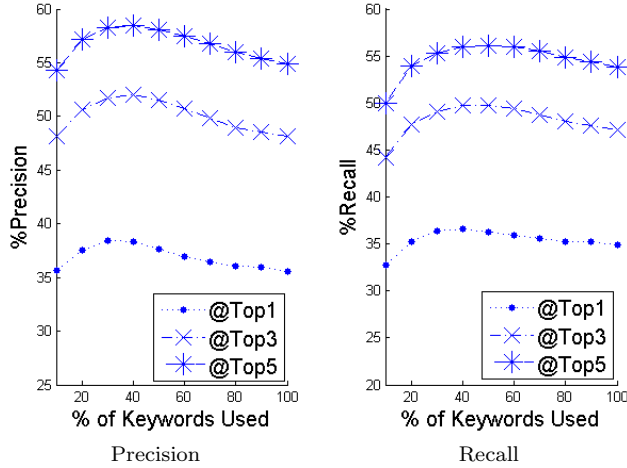


Fig. 14: Precision and Recall for TG-TI (@0-Step).

The results of this evaluation are depicted in Figures 17 and 18 (we report the results for the city of Milan).

When using the 1 – Step evaluation, we observe an increase of up to 6% for precision, and up to 4% for recall. The additional benefit for 2 – Steps is diminishing, exhibiting an increase of up to 4% for precision and up to 2% for recall. This effect of diminishing returns is due to the fact that immediately neighboring squares tend to share the same topic, while the topic dilutes and differs more when we move further away. In all cases, TG-TI-CLR1 accounts

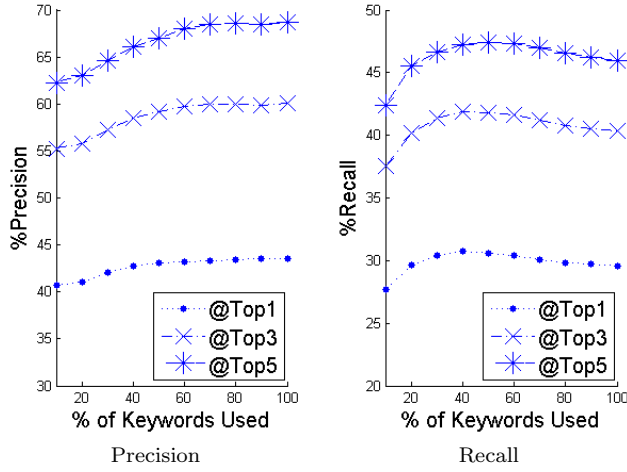


Fig. 15: Precision and Recall for TG-TI-CLR1 (@0-Step).

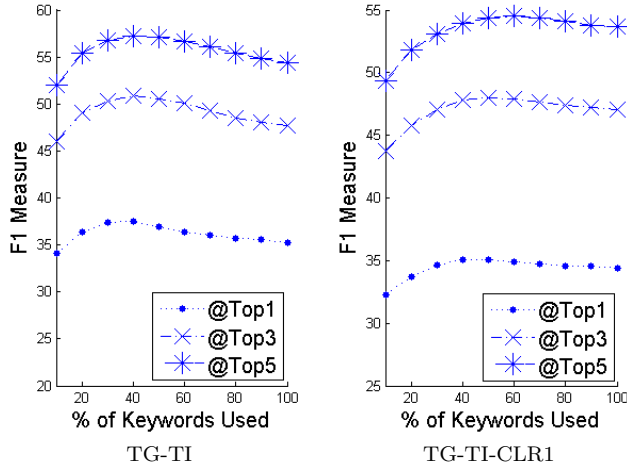


Fig. 16: F1 Score for TG-TI and TG-TI-CLR1 (@0-Step).

for the best mean precision. The second best precision is the one achieved by TG-TI, while the third best is achieved by TG-CLR1.

Finally, we run experiments modifying at the same time all the three parameters presented before, namely the similarity threshold, the @Step and the TopK. In Figure 19, we illustrate the precision and recall of the TG-TI-CLR1 method. The results show that we can achieve a significant increase in precision, but only a modest increase in recall. We note that precision hovers above the 75%, therefore, making the proposed approach attractive for applications that need access to the geolocations of tweets.

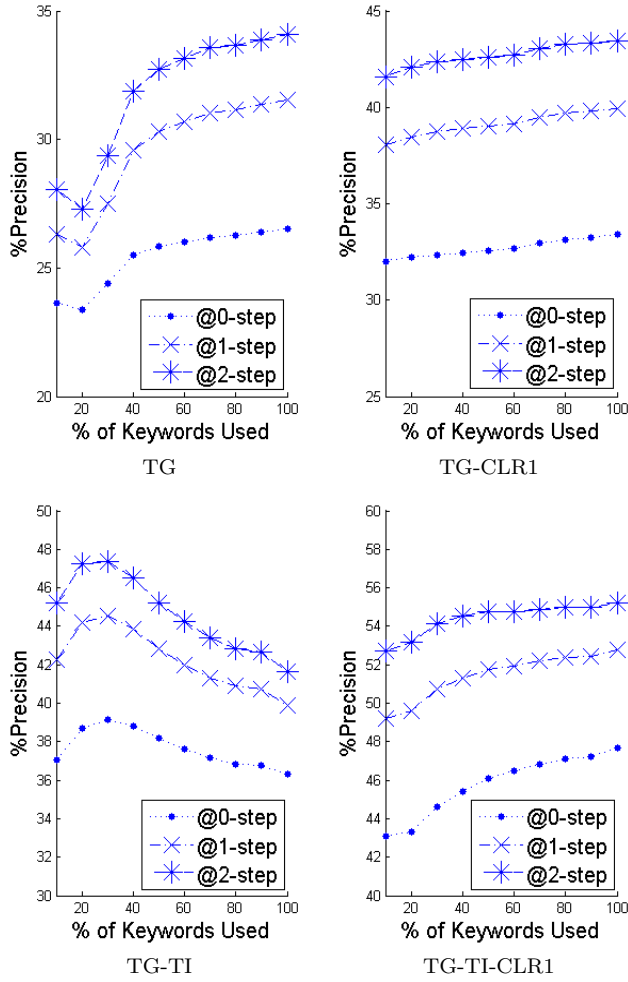


Fig. 17: Precision for Neighbourhood Level (@Top1).

### Size of Search Space

In order to evaluate our method with larger search spaces, we created a bigger grid for the city of Rome, and we ran experiments on this new dataset. In particular, we created a grid of 900 squares (30 by 30), while keeping the rest of the setup parameters the same. The size of each square is the same as before: 1km. In Figure 20, we compare the precision and recall of the Tf-Idf methods for the 20 by 20 and the 30 by 30 grids.

The best precision for the 30 by 30 grid is 45% and is achieved by TG-TI-CLR1 when using 100% of keywords, while the best recall is 37% and achieved by TG-TI when using 40% of the keywords. As expected due to the higher search space, the precision and recall achieved by each method are lower than



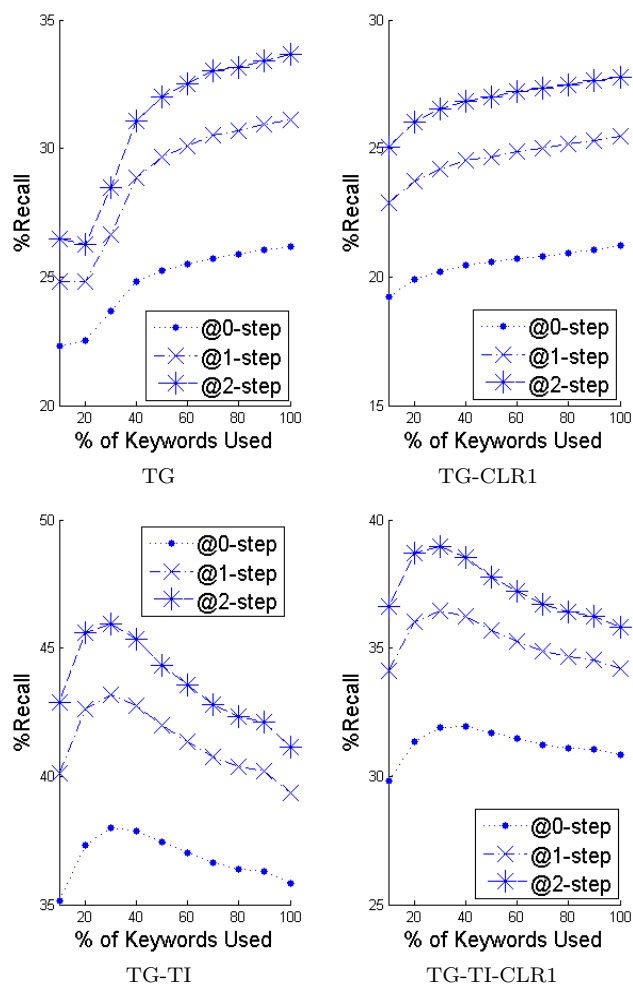


Fig. 18: Recall for Neighbourhood Level (@Top1).

those for the smaller grid: they were up to 4% lower for both algorithms, when the search space increased by 225%. These results demonstrate that the effect of the increase of the search space on the proposed algorithms is relatively small.

### 5.3 Discussion

Overall, our results show that using the correlation between local and global activity has the potential (when properly employed) to lead to significantly better accuracy.

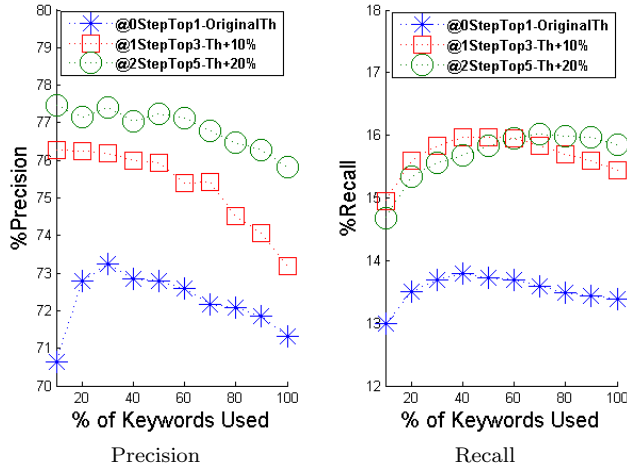


Fig. 19: Precision and Recall for TG-TI-CLR1 for varying similarity threshold, TopK, and @Step.

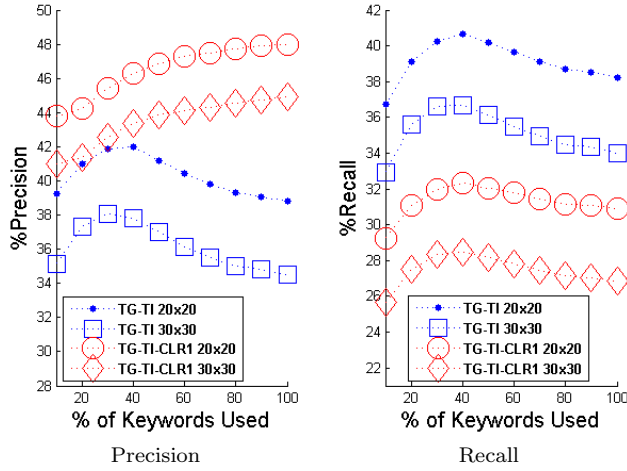


Fig. 20: Precision and Recall Comparison for the City of Rome (grids 20x20 and 30x30, @Top1 and @0-Step).

We also observe that, contrary to previous work, the time needed to train and test our models depends on the percentage of keywords used. This allows us to achieve a trade-off between execution time and accuracy. An interesting point regarding this trade-off is the fact that the increase of the execution time that the *CLR* methods exhibit when using higher percentage of keywords, does not pay off with a proportional increase in precision and/or recall.

For the @Top1 case, the Tf-Idf based algorithms are the winners, providing better results than the simpler algorithms based on concordance. Furthermore,

when using the Tf-Idf based algorithms, the best result is achieved when pruning some of the keywords. This is due to the fact that pruning the keywords with the lowest weight, we primarily remove stopwords, which has a positive impact on accuracy. This is also true for TG-TI-CLR1, when considering the F1 score.

Regarding the difference in precision and recall between our approach and the baselines, we believe that it is due to the very different granularity requirements of the problems, especially the temporal granularity. Even though the baselines provide good results for identifying the characteristic topics of a location (when there are enough data), our approach has an advantage for geolocating tweets referring to time-focused events, especially those with a relatively short time-span (e.g., concerts).

## 6 Conclusions

The extended use of social networks has resulted in an abundance of information on different aspects of everyday social activities, and has led to a proliferation of tools and applications that can help end-users and large-scale event organizers to better plan and manage their activities. Several of these applications are based on the knowledge of the geolocation of the relevant information. However, in Twitter, only a small percentage of the posts are geotagged.

In this work, we address the problem of geolocating non-geotagged tweets. We have proposed a framework that allows the estimation of the location from which a post was generated, by exploiting the similarities in the content between this post and a set of geotagged tweets. Contrary to previous approaches, our framework provides geolocation estimates at a fine grain, thus, supporting a range of applications that require this detailed knowledge. The experimental evaluation with real data demonstrates the efficiency and effectiveness of our approach, which when coupled with the right visualizations [27], can become a powerful analysis tool. In our future work, we plan to study the use of more elaborate models for the representation of the keywords in a tweet. The challenge here is to identify the right abstraction, given the short length of tweets.

## Acknowledgments

This work was supported by a fellowship from Telecom Italia.

## References

1. Facebook, <https://www.facebook.com/>
2. Google+, <https://plus.google.com>
3. Twitter, <https://twitter.com>
4. Abdelhaq, H., Sengstock, C., Gertz, M.: Eventweet: Online localized event detection from twitter. *Proceedings of the VLDB Endowment* **6**(12) (2013)
5. Ajao, O., Hong, J., Liu, W.: A survey of location inference techniques on twitter. *Journal of Information Science* **41**(6), 855–864 (2015)

6. Balduini, M., Bocconi, S., Bozzon, A., Della Valle, E., Huang, Y., Oosterman, J., Palpanas, T., Tsytsarau, M.: A case study of active, continuous and predictive social media analytics for smart city. In: ISWC Workshop on Semantics for Smarter Cities (S4SC)
7. Balduini, M., Della Valle, E., Dell'Aglio, D., Tsytsarau, M., Palpanas, T., Confalonieri, C.: Social listening of city scale events using the streaming linked data framework. In: ISWC (2013)
8. Blei, D.M., Ng, A.Y., Jordan, M.I.: Latent dirichlet allocation. *the Journal of machine Learning research* (2003)
9. Chang, H.w., Lee, D., Eltaher, M., Lee, J.: @ phillies tweeting from philly? predicting twitter user locations with spatial word usage. In: ASONAM (2012)
10. Cheng, Z., Caverlee, J., Lee, K.: You are where you tweet: a content-based approach to geo-locating twitter users. In: CIKM (2010)
11. Crooks, A., Croitoru, A., Stefanidis, A., Radzikowski, J.: # earthquake: Twitter as a distributed sensor system. *Transactions in GIS* **17**(1), 124–147 (2013)
12. Earle, P.S., Bowden, D.C., Guy, M.: Twitter earthquake detection: earthquake monitoring in a social world. *Annals of Geophysics* **54**(6) (2012)
13. Eisenstein, J., O'Connor, B., Smith, N.A., Xing, E.P.: A latent variable model for geographic lexical variation. In: EMNLP (2010)
14. Frias-Martinez, V., Soto, V., Hohwald, H., Frias-Martinez, E.: Characterizing urban landscapes using geolocated tweets. In: SocialCom-PASSAT (2012)
15. Han, B., Cook, P., Baldwin, T.: Text-based twitter user geolocation prediction. *JAIR* (2014)
16. Hossain, N., Hu, T., Feizi, R., Zheng, D., White, A.M., Luo, J., Kautz, H.: Precise localization of homes and activities: Detecting drinking-while-tweeting patterns in communities (2016)
17. Ikawa, Y., Enoki, M., Tatsubori, M.: Location inference using microblog messages. In: Proceedings of the 21st international conference companion on World Wide Web, pp. 687–690. ACM (2012)
18. Kinsella, S., Murdock, V., O'Hare, N.: I'm eating a sandwich in glasgow: modeling locations with tweets. In: SMUC (2011)
19. Leetaru, K., Wang, S., Cao, G., Padmanabhan, A., Shook, E.: Mapping the global twitter heartbeat: The geography of twitter. *First Monday* **18**(5) (2013)
20. Li, C., Sun, A.: Fine-grained location extraction from tweets with temporal awareness. In: SIGIR (2014)
21. Malmi, E., Do, T.M.T., Gatica-Perez, D.: From foursquare to my square: Learning check-in behavior from multiple sources. In: ICWSM (2013)
22. Mathioudakis, M., Koudas, N.: Twittermonitor: trend detection over the twitter stream. In: SIGMOD (2010)
23. Murdock, V.: Your mileage may vary: on the limits of social media. *SIGSPATIAL Special* (2011)
24. Paradesi, S.M.: Geotagging tweets using their content. In: FLAIRS Conference (2011)
25. Paraskevopoulos, P., Dinh, T.C., Dashdorj, Z., Palpanas, T., Serafini, L.: Identification and characterization of human behavior patterns from mobile phone data. *NetMob* (2013)
26. Paraskevopoulos, P., Palpanas, T.: Fine-grained geolocalisation of non-geotagged tweets. In: Proceedings of the 2015 IEEE/ACM International Conference on Advances in Social Networks Analysis and Mining 2015, pp. 105–112. ACM (2015)
27. Paraskevopoulos, P., Pellegrini, G., Palpanas, T.: When a tweet finds its place: Fine-grained tweet geolocalisation. In: International Workshop on Data Science for Social Good (SoGood), in conjunction with the European Conference on Machine Learning and Principles and Practice of Knowledge Discovery (ECML PKDD) (2016)
28. Sakaki, T., Okazaki, M., Matsuo, Y.: Earthquake shakes twitter users: real-time event detection by social sensors. In: WWW (2010)
29. Schulz, A., Hadjakos, A., Paulheim, H., Nachtwey, J., Mühlhäuser, M.: A multi-indicator approach for geolocalization of tweets. In: ICWSM (2013)
30. Serdyukov, P., Murdock, V., Van Zwol, R.: Placing flickr photos on a map. In: SIGIR (2009)
31. Tsytsarau, M., Amer-Yahia, S., Palpanas, T.: Efficient sentiment correlation for large-scale demographics. In: SIGMOD (2013)

32. Tsytsarau, M., Palpanas, T.: Survey on mining subjective data on the web. *Data Mining and Knowledge Discovery* (2012)
33. Tsytsarau, M., Palpanas, T.: Nia: System for news impact analytics. *KDD Workshop on Interactive Data Exploration and Analytics (IDEA)* (2014)
34. Tsytsarau, M., Palpanas, T., Castellanos, M.: Dynamics of news events and social media reaction. In: *SIGKDD* (2014)
35. Tsytsarau, M., Palpanas, T., Denecke, K.: Scalable discovery of contradictions on the web. In: *WWW* (2010)
36. Tsytsarau, M., Palpanas, T., Denecke, K.: Scalable detection of sentiment-based contradictions. *DiversiWeb, WWW* (2011)
37. Van Canneyt, S., Van Laere, O., Schockaert, S., Dhoedt, B.: Using social media to find places of interest: a case study. In: *SIGSPATIAL (GEOCROWD)* (2012)
38. Yuan, Q., Cong, G., Ma, Z., Sun, A., Thalmann, N.M.: Who, where, when and what: discover spatio-temporal topics for twitter users. In: *SIGKDD* (2013)
39. Zafarani, R., Liu, H.: Evaluation without ground truth in social media research. *Communications of the ACM* **58**(6), 54–60 (2015)