

WEB PLACE NAME DICTIONARY IMPLEMENTATION USING TWITTER AS SOURCE TO DEVELOP TRAFFIC INFORMATION SYSTEM

Rachmad Akbar, Ari Wibisono and Petrus Mursanto

Faculty of Computer Science, Universitas Indonesia,
Kampus Baru UI Depok, Jawa Barat, 16424, Indonesia

E-mail: rachmad.akbar@ui.ac.id

Abstract

Congestion is one of the severe problems that occurs in big cities like Jakarta, the center of the Indonesian government. One alternative solution to solve this problem is with a Traffic Information System which is easily accessible by the public. Twitter as a social media can be used as one of information source to develop the Traffic Information System. The most common problem in the processing of information on Twitter, especially in Indonesia, is the use of abbreviations and typographical error (typo). We propose an application of Web Place Name Dictionary in addressing abbreviations and typos in the name of the location contained on Twitter to develop a Traffic Information System. Three cases in find coordinate point in the street that derived from location name in Twitter can be handled through the proposed system. This system produces 90% accuracy rate on 100 data set that obtained via @TMCPoldaMetro. This study only focus on the implementation of Web Place Name Dictionary method so we do not compare our system with the other methods.

Keywords: *Coordinate, Dictionary, Location, Traffic, Twitter*

Abstrak

Kemacetan merupakan salah satu masalah berat yang terjadi pada kota-kota besar di Indonesia, tidak terkecuali di DKI Jakarta yang menjadi pusat pemerintahan Indonesia. Salah satu alternatif solusi untuk mengatasi masalah tersebut adalah dengan adanya Traffic Information System yang mudah diakses oleh masyarakat. Media sosial Twitter dapat dimanfaatkan sebagai salah satu sumber informasi untuk mengembangkan Traffic Information System tersebut. Permasalahan yang sering dihadapi dalam pengolahan informasi pada Twitter, khususnya di Indonesia adalah penggunaan singkatan serta typo. Kami mengajukan penerapan *Web Place Name Dictionary* dalam mengatasi masalah singkatan serta kesalahan ketik (*typo*) pada informasi nama lokasi yang terdapat pada Twitter untuk mengembangkan Traffic Information System. Tiga kasus pencarian titik koordinat pada jalan berdasarkan nama lokasi yang diperoleh dari Twitter dapat ditangani melalui sistem yang diajukan. Sistem ini menghasilkan akurasi 90% pada 100 data uji nama jalan yang diperoleh melalui @TMCPoldaMetro. Studi ini hanya berfokus pada implementasi metode *Web Place Name Dictionary* sehingga perbandingan dengan metode lain tidak dilakukan.

Kata Kunci: *Lalu lintas, Lokasi, Kamus, Koordinat, Twitter*

1. Introduction

Congestion is the impact of the imbalance between capacity of the road and the number of passing vehicles. Currently the length of roads in Jakarta are only 7,208 kilometers [1], while the number of motor vehicles until April 2012 is 13,346,802 [2]. Based on that facts, we will produce a ratio of 1,851 vehicles per kilometer. The number is very far from the ideal ratio. The ideal ratio is less than 100 vehicles per kilometer [3].

The rapid development of technology makes price of mobile devices become more affordable. Hence, the use of social media Twitter increases.

Various informations are distributed to the public through the media, including information of traffic conditions shared by @TMCPoldaMetro, @LewatMana and @SonoraFM92. This phenomenon can be used as an opportunity to develop a Traffic Information System to fulfill the needs of information regarding traffic conditions on the roads that are going to be passed.

The most common problem in the processing of information on Twitter, especially in Indonesia, is the use of abbreviations and typographical error (typo). There are some methods that can be used to solve this problem, such as utilization of geotag, utilization of tagging by tagger, and construction of a web place-name dictionary [4].

Web place-name dictionary is chosen to be a method because Indonesian people seldom use geotagging feature and usually use non-standard sentence that causes tagger tagging process to be difficult to implement.

Similar studies have been carried out by [5]. Endarnoto used Natural Language Processing to identify the name of the location being mentioned in a Tweet. However, the research conducted does not include how to obtain precise location coordinates of a location name that has been obtained through Twitter.

2. Literature Review

2.1 OpenStreetMap

OpenStreetMap (OSM) is a free and editable map that allows user to access the entire dataset on the map (latitude and longitude coordinates) around the world [6]. Disclosure of this dataset can be used as a source of data and information about coordinates of path and name of street in Jakarta. The structure of the node in OSM can be seen in figure 1.

```
<node id="29938999" version="9"
timestamp="2011-01-31T14:07:19Z"
uid="38066" user="ArjanO" changeset
="7145221" lat="-6.1811231" lon=
"106.8227586"/>
```

Figure 1. Data structure of node in OSM

A street is represented as a collection of nodes that have latitude and longitude coordinates, the nodes is ordered from starting point to the end point of the street. The structure of OSM street can be seen in figure 2.

```
<way id="40198102" version="4" timestamp=
"2011-10-20T01:47:36Z" uid="76518" user="Firman Hadi"
changeset="9604949">
<nd ref="484322049"/>
<nd ref="484321975"/>
<nd ref="484321937"/>
<nd ref="1473672608"/>
<nd ref="484326153"/>
<nd ref="363186769"/>
<nd ref="484326063"/>
<tag k="highway" v="primary"/>
<tag k="name" v="Kalibata"/>
<tag k="oneway" v="yes"/>
</way>
```

Figure 2. Data structure of way in OSM

MapQuest is one of the leading providers of online mapping services that support OSM [7]. MapQuest provides several features to support developers to build applications, both desktop applications, web and mobile. JavaScript 6.0 SDK is available to develop applications on the client-

side. While on the server-side developers can develop an application using C++, .NET and Java SDK [8]. In determining the coordinate of a point on the street that represents information from the Twitter, a mechanism is needed to calculate the distance between two coordinate points via existing roads. One of the services that can be used is a web-service which is provided by yournavigation.org.



Figure 3. MapQuest display

Yournavigation.org use Gosmore routing machine as routing engine so that the service is more lightweight, scalable and fast [9]. There are some parameters is used in this web service :

- flat = latitude from start location
- flon = longitude from start location
- tlat = latitude from destination location
- tlon = longitude from destination location
- v = kind of vehicle
- fast = 1 for fastest route, 0 for shortest route
- layer = the selection of Gosmore agencies that used to calculate the route

The yournavigation.org web service returns a kml file that has a structure shaped like xml, so the file is easy to understand. One of the data used in developing this system is the value of the distance from two point coordinates. The sample of yournavigation.org kml file can be seen in figure 4. In fact, there are 3 kinds of Tweet that will be informed by @TMCPoldaMetro. The kinds are Tweet with one location name, two location and three location names. The Tweet that has two or three location names will be processed by utilizing the web service.

2.2 Web Place-name Dictionary

One of the solutions that can be proposed to overcome the problem of finding matching location name by raw location name obtained from Twitter is Web Place-name Dictionary [4]. Web place-name dictionary is made by collecting

the sources of the location name, then using it as a database that will be accessible to search for the location names that are similar to a given word or query.

```
<kml xmlns="
http://earth.google.com/kml/2.0">
<Document>
  <name>KML Samples</name>
  <open>1</open>
  <distance>77.974443</distance>
  <traveltime>3325</traveltime>
  <description>To enable simple
instructions add: 'instructions=1' as
parameter to the URL</description>
  <Folder>
    <name>Paths</name>
    <visibility>0</visibility>
    <description>Examples of paths.
</description>
    <Placemark>
      <name>Tessellate</name>
      <visibility>0</visibility>
      <description><![CDATA[If the
<tessellate> tag has a value of
1, the line will contour to the
underlying terrain]]></description>
      <LineString>
        <tessellate>1</tessellate>
        <coordinates> 7.962893,52.214513
```

Figure 4. kml file example

In developing this system, the name and street coordinates are obtained by utilizing the data sets held by OpenStreetMap data on a particular area of Jakarta. The data used can be obtained through http://downloads.cloudmade.com/asia/southeastern_asia/indonesia/jakarta_raya/jakarta_raya.highway.osm.bz2. Data that will be used in this study is only the primary and secondary data path of the street that is known by tags `<tag k="highway v="primary"/>` and `<tag k="highway v="secondary"/>`.

2.3 String Matching Algorithm

String Matching Algorithm is an algorithm for searching string that have the highest similarity to the given query string. This algorithm is used to solve the typo and abbreviation problem that often happens on Indonesian Tweets.

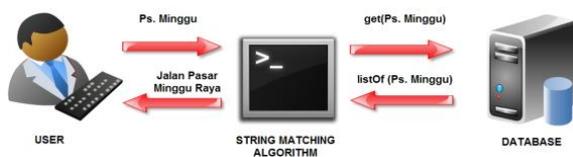


Figure 5. String Matching Algorithm Process

String Matching Algorithm consists of two main processes, search string algorithm to find candidate result and the calculation of similarity between each candidate string and query string. The string matching algorithm flow details can be seen in figure 6.

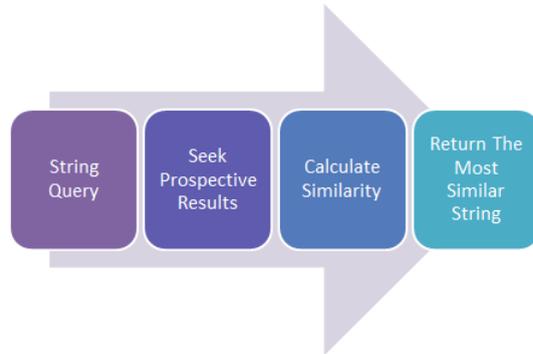


Figure 6. String Matching Algorithm Flow Details

Q-grams method is an algorithm that is often used in search string candidates. Q-grams chops the string into several substrings with each substring consists of q characters. Q-grams method is divided into two types, overlapping and consecutive. In the overlapping Q-Grams, a q-gram will begin at each position in the string. While for consecutive Q-Grams, a q-gram will begin at each position that is multiple of q [10]. On a String "juanda" overlapping 3 - grams will yield {(jua), (uan), (and), (nda)}, whereas if using consecutive 3 - grams will yield {(jua), (nda)}. The value of Q can be changed as needed, but in this study we used Q = 3. This is based on the characteristics of words in Indonesian and street naming in Indonesia.

There are three methods commonly used to calculate the similarity between two strings; Levenshtein Distance, Jaccard, and Cosine. Levenshtein Distance algorithm or often called Edit distance describes the similarity value between two strings based on minimal number of insert, delete, and substitution process to transform a string into another string [11] [12].

Edit distance $ed(x, y)$ between string $x = x_1...x_m$ and $y = y_1...y_n$ where $x, y \in \Sigma^*$ and Σ is alphabetic character. We can construct a matrix with size $M_{1...m+1, 1...n+1}$ to calculate $ed(x, y)$ where the value of edit distance is at $M_{m+1, n+1}$. The calculation of each cell at the position $M_{i, j}$ is calculated based on equation 1, where the value of

$$\delta(a, b) = 0 \text{ if } a = b \text{ and } 1 \text{ if } a \neq b.$$

$$M_{1,1} \leftarrow 0$$

$$M_{i,j} \leftarrow \min \begin{cases} M_{i-1,j}+1 \\ M_{i,j-1}+1 \\ M_{i-1,j-1}+\delta(x_i, y_j) \end{cases} \quad (1)$$

Even if we get the value of the edit distance, we have not been able to assess whether the two strings are similar or not because the value of the edit distance is relative to the length of the string. We need the equation to calculate how similar the two strings. One of the equations that can be used is

$$sim(s_1, s_2) = \frac{\max(|s_1|, |s_2|) - ed(s_1, s_2)}{\max(|s_1|, |s_2|)} \quad (2)$$

Unlike the Levenshtein Distance algorithm, Algorithm Jaccard and Cosine do not count from letter to letter but the grams that are formed. Set X is the result of Q-grams method to first String, and Set Y is the result of Q-grams method to second String. Jaccard algorithm calculates the similarity value by the equation:

$$Jaccard(x, y) = \frac{|X \cap Y|}{|X \cup Y|} \quad (3)$$

Cosine algorithm calculates the similarity value by the equation:

$$Cosine(x, y) = \frac{|X \cap Y|}{\sqrt{|X| \cdot |Y|}} \quad (4)$$

In this study, we only use Levenshtein Distance algorithm because the sequence of the Q-gram made by Q-gram method is an important point.

2.4 Efficient Approximate Candidate String

In developing this kind of system, efficiency is certainly an important factor. One method that can be applied is the Scan Count and Divide Skip, which is an algorithm proposed by Li in 2008 [13].

Scan Count algorithm counts the number of occurrences of each string in the database that contains the elements of a given set of queries grams. Scan Count Algorithm is then developed into Divide Skip. Divide Skip is an algorithm where filters are added into Scan Count algorithm in advance to speed up the calculation. We applied some modifications to suit the environment of the experiments conducted. The pseudo code of Divide Skip algorithm can be seen in figure 6.

In general, Divide Skip algorithm discards non potential data by dividing the entire list into two parts, list with a little number of members (L short) and list with the lot number of members (L

long). The main problem in this algorithm is the selection of the value of L so that the algorithms can run properly and efficiently.

```

Input : lists and threshold T
Output : List of string that appears at least T times
1   Initiate an empty list R
2   Create a L long list which is a collection of L
    longest list of the input list
3   Create a L short list which is a collection of the
    remaining list (short lists)
    Use Scan Count Algorithm for L short to find a
4   collection of records that appears at least T-L
    times
5   FOR (each record r that found) {
6       FOR (each lists in L long )
7           Check whether r appears in the list
8           IF (r appears >= T times on the whole list)
9               Add r to R;
10  }
11  RETURN R;
    
```

Figure 6. Pseudocode of Divide Skip Algorithm

Li stated that the value of L can be found using the equation 5:

$$L = \frac{T}{\mu \log M + 1} \quad (5)$$

M is the length of the longest list, and μ is an independent coefficients. Based on the results of the experiment that have been carried out by Li, the value of μ is 0.0085.

3. Analysis and Result

3.1 Problem Identification

Problem identification is done by analyzing the pattern of the provision of information via Tweet from the account @TMCPoldaMetro and methods that have been applied to map the location of the street based on information from social media. The methods that have been applied in Indonesia is utilizing Google API with location or street name as an input.

The use of Google's API does not guarantee that the precise location can be found as a result of the presence of the typos and abbreviation of the name of the location/street.

There are three cases that must be solved in this system. The case with one location name, two locations name, and three locations name.



Figure 7. Tweet with one location

Tweet with one location name usually informs the overall road conditions.



Figure 8. Tweet with two location

Tweet with two locations usually informs the certain path on the way.



Figure 9. Tweet with three location

Tweet with three locations usually inform the way in certain path that intersects with some other way.

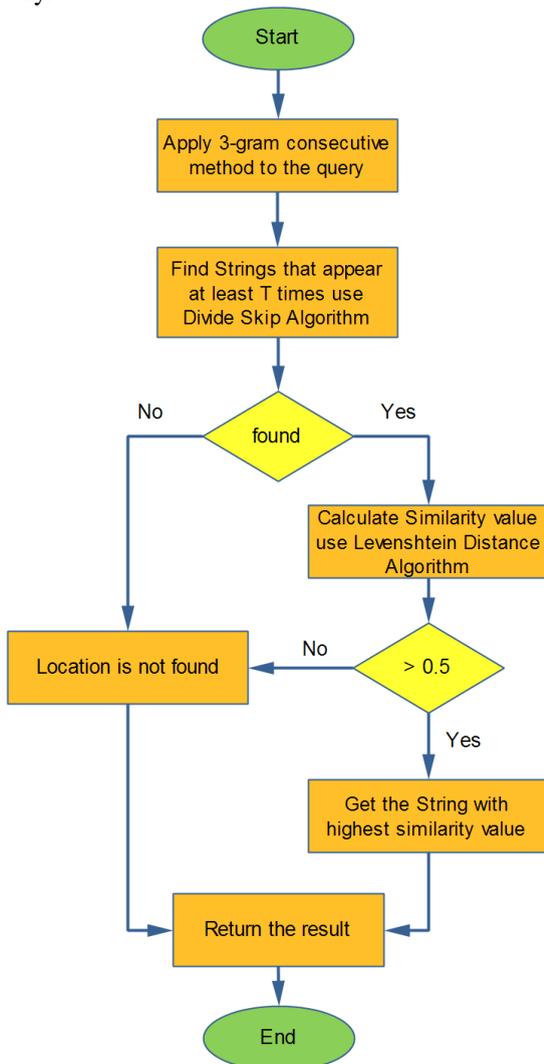


Figure 9. Flow diagram for one input query

3.2 System Requirement

In developing a system that is able to provide accurate coordinate information, a fairly in-depth analysis of the pattern of information provided through Twitter is required. In general, the system is built to meet the following specifications:

- The system is able to search the street name accurately although there are some typos or abbreviation.
- The system is able to locate specific coordinates on the road with one input location name.
- The system is able to locate specific coordinates on the road with two input location name.
- The system is able to locate specific coordinates on the road with three input location name.
- The system is able to visualize on a map the coordinates of the intended path.

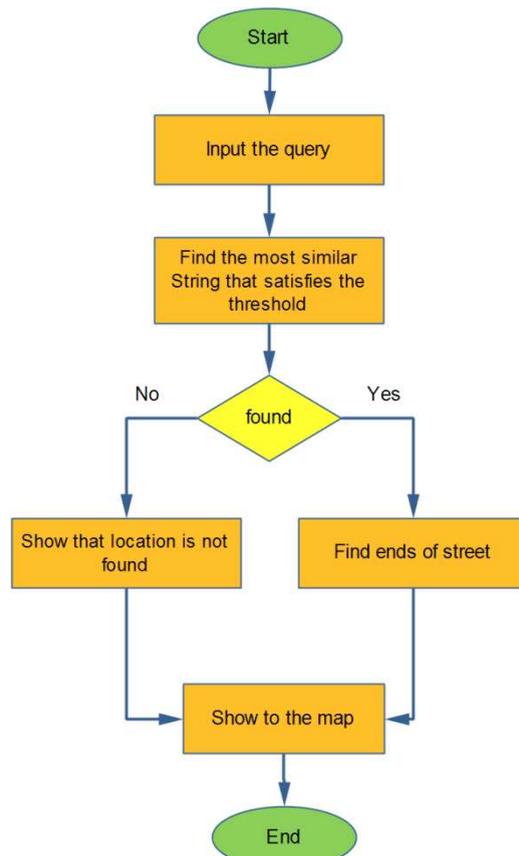


Figure 10. Flow diagram for one input query

3.3 Systems Design

In developing this system, the primary focus is the process of finding a location or street name that corresponds to what is meant by information providers. The flow diagram is shown in figure 9. The value of T that is used in the flow can be

obtained by calculating the result of $(\text{query.length} - 2) / 2$. This means that there are at least a candidate string containing half of gram that are formed from a given query.

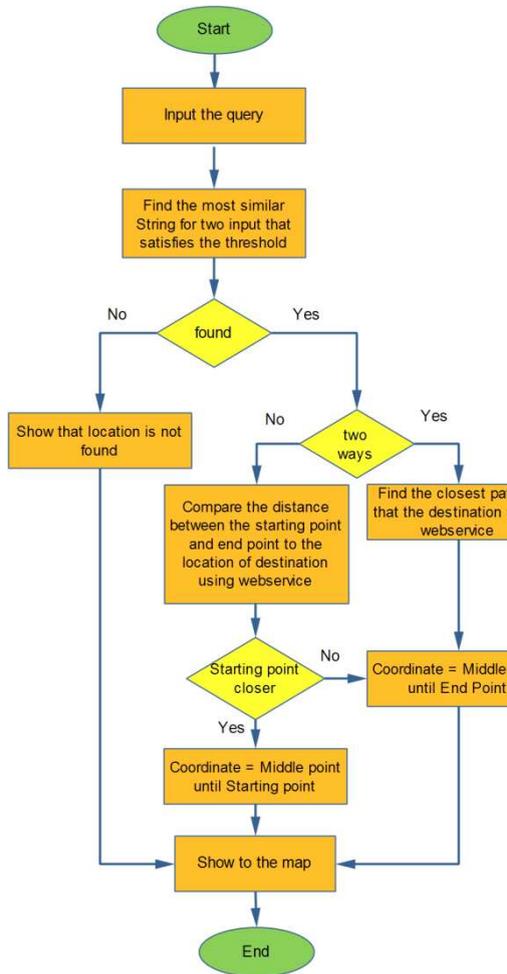


Figure 11. Flow diagram for two input query

The system has four main flows; flow on the input processing, flow on finding coordinate by one location name, flow on finding coordinate by two locations name, and the last is flow on finding coordinate by three locations name. The system will return the coordinates that consist of latitude and longitude which considered as the location referred by a Tweet.

On the flow of processing one input query, the coordinates that will be taken are the start point and end point of a path. This decision is made because the location that is meant usually covers the entire road. If the road which is meant is a two-way street, then the system will do the search process twice to obtain the coordinates of the starting point and end point on each track. Then the coordinates will be visualized into a map.

Contrast to one input query processing, two input query processing utilizes yournavigation.org web service to determine the most appropriate location that is provided by a Tweet. There are two possibilities in this case; first the case where the location is from middle point to starting point or the second case is from middle point to end point on the road.

The determination of the case is based on the reference point given by the second query. The first query is a street that will be searched for its coordinates. If the road on the first query is a road with two ways, we must search which lane that is closest to the reference point on the second query.

The case of two input queries is similar to the case of three input queries. We use the services of yournavigation.org to determine the specific coordinates on the way. In the case of three input query, the system will search the node that has shortest distance to the first reference (query 2) and the node with the shortest distance to the second reference point (query 3).

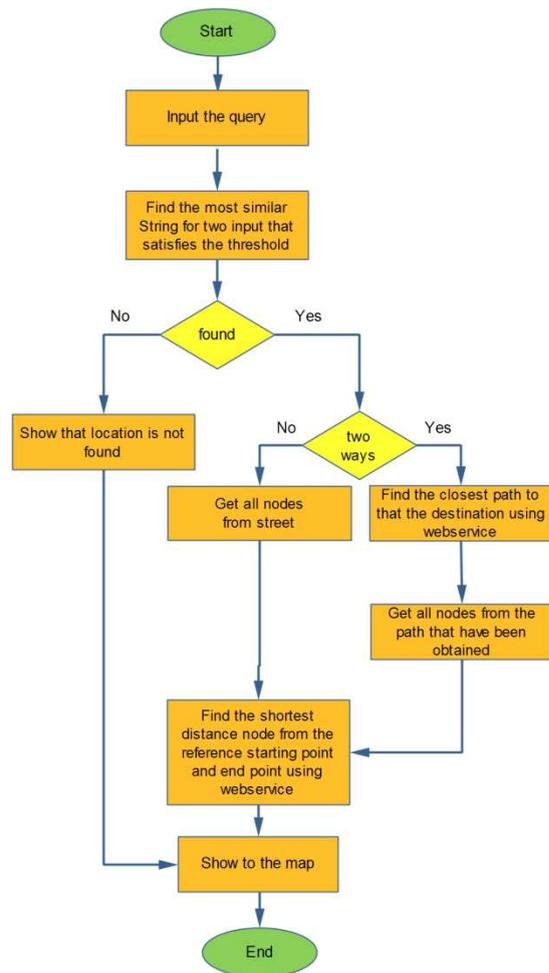


Figure 12. Flow diagram for three input query

3.4 Implementation

There are some applications that are used in building this system such as Tomcat server, Xampp and MySQL. Implementation of the code is done by using the Java language. Figure 13 is an overall system view.



Figure 13. Prototype of system

The prototype consists of two component; input menu for one input, two inputs, three inputs (figure 14) and the map to visualize the coordinate (figure 15).

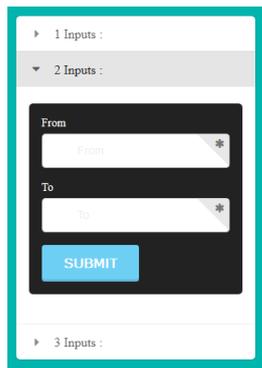


Figure 14. Input Query Field

The given input on the field will be sent to the server and processed to produce the corresponding location as shown in figure 15.

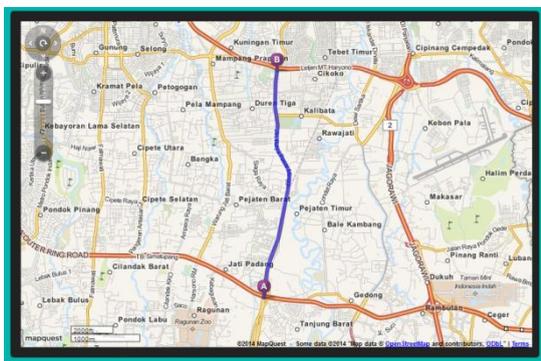


Figure 15. Location visualization

3.5 Testing

Testing of the prototype system is done by conducting a blackbox testing against the street name / location name as input. There are 100 street name / location name that are taken randomly from the data that ever provided in Tweets by @TMCPoldaMetro from May to August 2013.

TABLE I
EXPERIMENT RESULT SAMPLE

Input	Output
JI H R Rsauna Said	data tidak ditemukan
Pos Tanung Priok	taman stasiun tanjung priok
JlThamrin	mh thamrin
jl Danau Sunter Utara	danau sunter utara
jalan Sunter Utara	danau sunter utara
Slipi Jaya	Slipi Jaya
JI Matraman	matraman raya
Tomang	tomang raya
JIPos Pengumben	pos pengumben
JIGatot Subroto	gatot soebroto
Kalibat	kalibata
Cililitan	cililitan besar
Tomang	tomang raya
JIRasuna Said	rasuna said
JJatibaru	jatibaru
JI Gajah Mada	gajah mada
Slipi	Slipi Jaya
JI Panglima Polim	panglima polim raya
JI Sudirman	sudirman
JI I Gusti Nghrah Rai	i gusti ngurah rai
Mangga Besar	mangga besar
jl n raya bogor	raya bogor
JI Margonda Raya	jalan raya margonda
Fatmawati	fatmawati
Pasar Jumat	pasar jumat
Pancoran	pancoran
JIPramuka Raya	pramuka
JlImam Bonjol	imam bonjol
jalan Raya Kalimalang	kalimalang
jalan Raya Pesanggrahan	pesanggrahan

Experiments show that the prototype system can recognize the input with an accuracy of 90% on 100 test cases.

4. Conclusion

Development of traffic information system using web place-name dictionary to determine the location mentioned in a tweet is able to handle the typos and abbreviations that are often done by Indonesian society. The system is able to reduce the inaccuracy that often occurs as a result of existing APIs such as Google Maps that could not solve the typos and abbreviations used by Indonesian people. In black box testing, the system is able to produce an accuracy of 90% with 100 data set that taken randomly.

Utilization of web place-name dictionary can be developed in a variety of other systems which utilize the name of the location mentioned through Twitter to be processed. In this study, the system still has weaknesses because there are no location name (non street name) on OpenStreetName so the problem must be resolved.

5. Acknowledgment

This research is partly supported by Masterplan Percepatan dan Perluasan Pembangunan Ekonomi Indonesia (MP3EI) and Hibah Riset PUPT UI (Awal) 2014. The grants have been made possible by Indonesia's Ministry of Education and Culture No:3715/H2.R12/HKP.05.00/2012 and Universitas Indonesia.

Reference

- [1] PU, "Panjang Jalan Jakarta Baru Penuhi 60 Persen Kebutuhan," [Online]. Available: www.pu.go.id/main/view_pdf/7843. [Accessed 17 November 2013].
- [2] BPS, "Perkembangan jumlah kendaraan bermotor menurut jenis tahun 1987-2011," [Online]. Available: http://www.bps.go.id/tab_sub/view.php?tabel=1&id_subyek=17¬ab=. [Accessed 17 November 2013].
- [3] D. J. Widiyanto, "Kebijakan dan Strategi Penanganan Kemacetan Lalulintas di Perkotaan," *tataruang*, vol. 2, no. 1, pp. 65--70, March 2008.
- [4] T. a. M. Y. a. Y. T. a. C. N. a. N. K. Sakaki, "Real-time event extraction for driving information from social sensors," in *Cyber Technology in Automation, Control, and Intelligent Systems (CYBER), 2012 IEEE International Conference on*, 2012.
- [5] S. Endarnoto, S. Pradipta, A. Nugroho and J. Purnama, "Traffic Condition Information Extraction and Visualization from Social Media Twitter for Android Mobile Application," in *Electrical Engineering and Informatics (ICEEI), 2011 International Conference on*, 2011.
- [6] OpenStreetMap, "About," [Online]. Available: <http://wiki.openstreetmap.org/wiki/About>. [Accessed 18 November 2013].
- [7] OpenStreetMap, "MapQuest," [Online]. Available: <http://wiki.openstreetmap.org/wiki/MapQuest>. [Accessed 18 November 2013].
- [8] MapQuest, "MapQuest," [Online]. Available: <http://company.mapquest.com>. [Accessed 18 November 2013].
- [9] I. A. Hamid, "YOURS navigation routing system - Lambertus Ijsselstein," *Society of Cartographer*, no. November, pp. 20-22, 2009.
- [10] C. Miao, G. Chang and X. Wang, "Filtering Based Multiple String Matching Algorithm Combining q-Grams and BNDM," in *Genetic and Evolutionary Computing (ICGEC), 2010 Fourth International Conference on*, 2010.
- [11] V. I. Levenshtein, "{Binary codes capable of correcting deletions, insertions, and reversals}," *Soviet Physics Doklady*, vol. 10, no. 8, pp. 707--710, 1966.
- [12] G. Navarro and M. Raffinot, *Flexible Pattern Matching in Strings: Practical Online Search Algorithms for Texts and Biological Sequences*, New York, NY, USA: Cambridge University Press, 2002.
- [13] Chen Li, Jiaheng Lu and Yiming Lu, "Efficient Merging and Filtering Algorithms for Approximate String Searches," in *Data Engineering, 2008. ICDE 2008. IEEE 24th International Conference on*, 2008.