

A Mobility Analytical Framework for Big Mobile Data in Densely Populated Area

Yuanyuan Qiao, *Member, IEEE*, Yihang Cheng, Jie Yang,
Jiajia Liu, *Senior Member, IEEE*, and Nei Kato, *Fellow, IEEE*

Abstract—Due to the pervasiveness of mobile devices, a vast amount of geolocated data is generated, which allows us to gain deep insight into human behavior. Among other data sources, the analysis of data traffic from mobile Internet enables the study of mobile subscribers' movements over long time periods at large scales, which is paramount to research over a wide range of disciplines, e.g., sociology, transportation, epidemiology, networking, etc. However, to efficiently analyze the massive data traffic from the view of user mobility, several technical challenges have to be tackled before releasing the full potential of such data sources, including data collection, trajectory construction, data noise removing, data storage, and methods for analyzing user mobility. This paper introduces a mobility analytical framework for big mobile data, based on real data traffic collected from second-, third- and fourth-generation networks, which covered nearly 7 million people. To construct a user's history trajectories, we apply different rules to extract users' locations from different data sources and reduce oscillations between the cell towers. The comparison of mobility characteristics between our mobile data and other existing data sources shows the large potential of mobile Internet data traffic to study human mobility. In addition, our experiments discover the changing of city hotspots, the movement patterns during peak hours, and people with similar history trajectories, which uncover the common rules that exist among huge populations in a city.

Index Terms—Big mobile data, human mobility, mobile internet, mobility analytical framework.

I. INTRODUCTION

THE STUDY of human mobility can yield insight into a variety of social issues on geographical scales, such as urban planning [1], population distribution [2], and the spread of disease [3]. As we know from our daily lives, the movement

Manuscript received November 5, 2015; revised February 19, 2016; accepted April 7, 2016. Date of publication April 12, 2016; date of current version February 10, 2017. This work was supported in part by the Fundamental Research Funds for the Central Universities under Grant 2015RC11, by the Director Foundation Project 2015BKL-NSAC-ZJ-01, by the 111 Project of China under Grant B08004, and by the EU Seventh Framework Program (FP7) International Research Staff Exchange Scheme MobileCloud Project under Grant 612212. The review of this paper was coordinated by Prof. Y.-B. Lin.

Y. Qiao, Y. Cheng, and J. Yang are with the Center for Data Science, the Beijing Key Laboratory of Network System Architecture and Convergence, and the Beijing Laboratory of Advanced Information Network, the School of Information and Communication Engineering, Beijing University of Posts and Telecommunications, Beijing 100876, China (e-mail: yyqiao@bupt.edu.cn; mrcheng@bupt.edu.cn; janyang@bupt.edu.cn).

J. Liu is with the School of Cyber Engineering, Xidian University, Xi'an 710071, China (e-mail: liujiajia@xidian.edu.cn).

N. Kato is with the Graduate School of Information Sciences, Tohoku University, Sendai 980-8579, Japan (e-mail: kato@it.ecei.tohoku.ac.jp).

Color versions of one or more of the figures in this paper are available online at <http://ieeexplore.ieee.org>.

Digital Object Identifier 10.1109/TVT.2016.2553182

of people in space is far from random. However, finding the formulation of quantitative laws explaining human mobility, which is essential to uncover the mechanisms governing human activities [4]–[9], still remains as an open question. There are essentially two ways for studying the nature of mobility: synthetic models and traces. Synthetic models attempt to represent the human behaviors by sets of mathematical equations, such as random mobility models (Lévy walk [10], random walk [11]), models with temporal dependence (Gauss–Markov [12], smooth random [13]), models with spatial dependence (probabilistic random walk [14]), models with geographic restriction (pathway or city section [15], obstacle [16]), and so on [17]–[19]. The aforementioned models are tractable, scalable, easy to deploy, and particularly useful in the field of *ad hoc* networks if traces have yet to be created [20]–[23]. Although synthetic models can easily reproduce human mobility patterns up to an acceptable degree of accuracy, it is still quite difficult to assess to what extent they map reality, and the generated trajectories are different from those observed in real scenarios [24]. On the other side, real human traces that provide accurate information, which traditionally are restricted by expensive data-collection methods, are of theoretical and practical significance in the area of mobility analytics.

Nowadays, smart devices bring us the ubiquitous mobile Internet access. People's movements could be sensed and easily collected by mobile phone, generating large volumes of mobility data, such as call detail records (CDRs), and global positioning system (GPS) tracks. CDRs provide the time that a phone placed a voice call or received a text message, as well as the identity of the cell tower with which the phone was associated at the time [6], [8], [25]–[27]. However, they are sparse in time and coarse in space, which limits the scope of their application to study human mobility. As for GPS tracks, the movements of individuals in latitude and longitude along with time stamps are recorded [28]. However, GPS signals may easily become unavailable in indoor or underground environments, GPS devices may get interferences in environments with high building density, and users are becoming more reluctant to share locations because continuously collecting GPS data may consume devices' energy quickly or make people uncomfortable, considering privacy issues [29]. Due to the aforementioned matters, up to now, there does not exist any GPS data source covering citywide population.

Recently, researchers have found that data traffic from second-, third-, and fourth-generation (2G/3G/4G) data networks is extremely useful for studying human dynamics [30]–[32]. Passively collecting human movement trajectories, while he/she is

accessing to mobile Internet, has lots of advantages: high cost efficiency, low energy consumption, covering a wide range and a large number of people, and with fine time granularity (people tend to surf mobile Internet frequently while moving, and many applications may send or receive network traffic packets periodically when running in the background). Collected trajectories are coarse in space because they record locations only at the granularity of a cellular antenna (with an average error of 175 m [33], and the density of cell tower is much larger in urban areas than in the suburbs or rural areas due to the human population density). This error range is tolerable, and the analytical results are convincing enough to find fundamental laws in human dynamics [6], [34], to build individual mobility models [35] or aggregated mobility models [7], even enough to get a dynamic understanding of the population, activities, and environment [36]–[41].

According to the prediction in [42], monthly global mobile data traffic reached 2.5 exabytes at the end of 2014, and this will surpass 24.3 exabytes by 2019. The explosion in data traffic amount brings many opportunities to obtain data sources with rich information. However, existing methods are not prepared to deal with such huge volume of traffic data. New methods to solve great challenge for data collection, storage, and analysis of big mobile data are needed urgently. As motivated by such observation, in this paper, our goal is to present a framework for efficiently analyzing massive data traffic from the view of user mobility in densely populated areas. The contributions of our work are summarized as follows.

- To the best of our knowledge, we are the first to present a cloud-computing-based analytical framework to analyze big mobile data, from the view of user mobility, covering the mobile networks of 2G/3G/4G and the scale for nearly 7 million people. Big data technologies and analytical algorithms are used to store and process massive data traffic. In particular, our framework is developed for analyzing user mobility patterns based on real mobile Internet data collected from 2G/3G/4G networks.
- Our framework is suitable for human trajectories consisting of a series of positions of cell towers, including CDRs. Since there is noise in raw data, different rules are required to construct human trajectory from different data sources. Toward this end, we define raw data processing rules for constructing human trajectory from different interfaces of 2G/3G/4G networks, and we remove data noise by reducing the oscillation between cell towers. In addition, to ensure the effectiveness of our data set, we calculate three widely accepted mobility indicators, i.e., the trip distance distribution, the radius of gyration, and the number of visited locations over time. Our results show that, for the same indicator, different data sources follow similar models but with different values of parameters.
- We further use our framework to explore human movement behavior in densely populated areas. We employ a parameter-free method to identify city hotspots from the view of population, apply a modified version of the Apriori algorithm to mine maximal sequential pattern, discover similar users based on their history trajectories,

and predict users' future movements from both temporal and spatial perspectives. These functionalities are of significant meaning for improving the user experience of location-based service (LBS), for optimizing network resources, and for advising city planning.

The remainder of this paper is organized as follows. In Section II, related works in the field of mobility analytics are introduced. Section III provides the overall structure of our mobility analytical framework, including data-collection methods for 2G/3G/4G networks, rules for constructing human trajectories, the design of database, and algorithms for mobility analysis. Section IV gives experimental results from our framework based on real data traffic. Conclusions are drawn in Section V.

II. AVAILABLE WORKS RELATED TO MOBILITY ANALYTICAL FRAMEWORKS

To study the inherent properties of human mobility in an efficient way, a mobility analytical framework, aiming to analyze big mobile data by providing data collection, data storage, data preprocessing function, and mobility functionalities, is essential. Mobility Profiler [38] is a complete framework for discovering mobility profiles from raw cell tower connection data. It removes the cell tower oscillation, and constructs a cell tower topology to discover a user's movement pattern. The framework "Jyotish" [43] constructs a predictive model by exploiting the regularity of people movement found in real joint WiFi/Bluetooth traces. With the rising of social–location–mobile (SoLoMo), based on big data platforms of IBM, Cao *et al.* presented a unified SoLoMo analysis approach from a system-oriented view [44], which is designed to process the vast amount of data generated in the telecommunications area every day. Zhang [30] proposed a systematic analysis methodology that considered inaccuracy from cellular data networks. Although the aforementioned frameworks have solved some issues regarding mobility analytics, none of them covered the whole procedure for data traffic analysis (i.e., data collection, data storage, noise removing, trajectory construction, and data analysis from the view of user mobility), considered the methods for big data storage and analysis, or gave detailed experiments with real-world data traffic collected from citywide areas.

With the rising of mobile Internet, research institutions and enterprises pay more and more attention to LBS and try to apply their mobility analytical frameworks to practical applications. IBM developerWorks introduces an Advanced Analysis Platform (AAP) for analyzing location to discover mobility pattern. They want to apply AAP to all kinds of location data, such as GPS data from cars, planes, or other equipment; the use of credit cards or public transportation cards; CDR; and deep packet inspection from operators. Microsoft Research developed "GeoLife" [45] to analyze GPS trajectory and provide location-based social-networking service. GeoLife is based on a framework, i.e., hierarchical-graph-based similarity measurement, to uniformly model each individual's location history and effectively measure the similarity among people. Different data sources usually have different features that require distinct

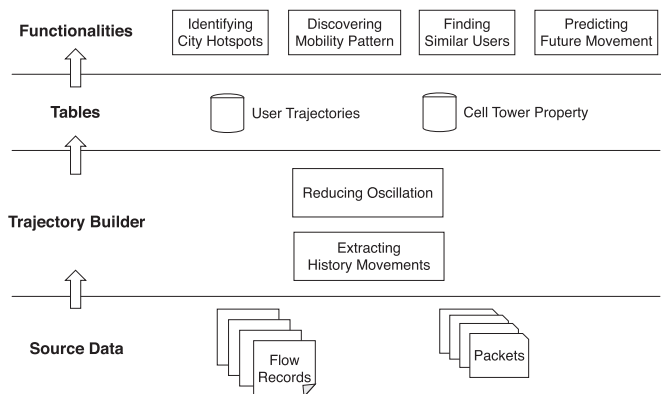


Fig. 1. Architecture of the mobility analytical framework.

analysis methods. For example, GPS trajectory is usually generated with occasional outliers or some noisy points caused by the poor signal of location positioning systems [46]. Raw cell tower connection data usually have “cell tower oscillation,” where, even when the user is static, he/she may be assigned to a number of neighboring cell towers because of load-balancing issues or changes in the ambient radio-frequency environment [38]. Hence, for different kinds of spatial trajectories, specific rules must be considered.

Different from the previous work, this paper aims at developing a framework for user mobility analytics based on massive mobile Internet data traffic, which integrates 1) the techniques for big data collection, storage, and preprocessing; 2) the rules for extracting location data and for constructing people trajectories; 3) the methods for solving data noise (i.e., cell tower oscillation); and 4) the algorithms for discovering common mobility patterns in densely populated areas. The real mobile Internet data traffic collected from 2G/3G/4G networks covering millions of people is used to verify the effectiveness of our framework. Our framework could be applied to analyzing human mobility with mobile Internet data traffic, and it is particularly useful for efficiently processing big mobile data.

III. METHODOLOGY

Here, we introduce a mobility analytical framework to analyze massive data traffic from mobile Internet. As shown in Fig. 1, we have two kinds of data sources, i.e., flow records and packets. Here, we define “flow” as the bidirectional data transmission at the usual 5-tuple source Internet protocol (IP), destination IP, source port, destination port, and transport protocol within a certain period of 64 s. Our framework is based on the cloud computing platform, which is the best tool to handle big data at present. It consists of a trajectory builder (removing data noise and extracting user history movements from cell tower ID sequence), a database (storing users’ trajectories and cell tower property), and functionalities (mobility analytics based on users’ trajectories).

A. Data Collection

By deploying our self-developed traffic monitoring system (TMS) at the core network connecting to the 2G/3G/4G

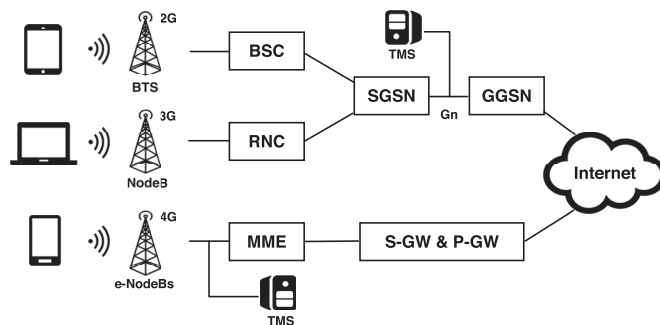


Fig. 2. Mobile Internet network architecture and the deployment of TMS.

network interfaces, data traffic generated by user equipment (UE), such as smartphones, tablets, laptop computers equipped with mobile broadband adapter, or any other devices that access to the Internet through 2G/3G/4G networks, is collected. As shown in Fig. 2, in 2G or 3G networks, a UE communicates with a base transceiver station (BTS) or Node B, which transmits its network traffic to a base station controller (BSC) or radio network controller (RNC). The controllers (BSC/RNC) then deliver the network traffic to a serving GPRS support node (SGSN) that establishes a tunnel on Gn interface (interface between the GGSN and the SGSN) with a gateway GPRS support node (GGSN) through which the data enters the Internet (GPRS represents “General Packet Radio Service”). In 4G networks, evolved Node B (eNodeB) establishes the connection between UE and mobility management entity (MME). Users’ data traffic goes into the Internet through the Serving GateWay (S-GW) and Packet Data Network (PDN) GateWay (P-GW).

In 2G/3G networks, we collect mobile Internet traffic from Gn interface and store the data traffic as flow records. To get people’s location information from 4G networks, we collect Long-Term Evolution (LTE) control-plane packets between eNodeBs and MME, which contain the whole signaling procedures, such as connection establishment, release procedure, or handover procedure. We can get a sequence of time-stamped records, each of which contains the current service eNodeB ID, signaling procedure code, user ID, etc.

As network applications become increasingly complex and heterogeneous, there is an increasing need for application-oriented traffic analysis. Most of the existing network TMSs only equip physical probes to capture and store raw packets because software-based traffic monitoring techniques are inadequate to achieve real-time monitoring. However, the TMS, which is based on a combined software/hardware architecture with flexibility to cope with the modification and addition of monitoring requirements, as well as future rate increase, can conduct application-oriented traffic analysis for a 10-Gb/s network line in real time using an eight-core machine [47].

B. Big Mobile Data Processing Platform

To store and process the massive data collected from a big city covering large population, a cloud-computing-based big mobile data platform with high storage capacity and computing power is essential. The mobility analytical framework is built on a

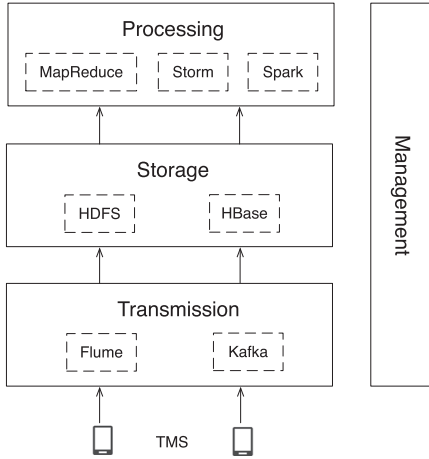


Fig. 3. System architecture of the big mobile data processing platform.

cloud computing platform based on Hadoop (an open-source software for reliable, scalable, and distributed computing) [48], which provides functions of data transmission, storage, processing, and management [47]. The system architecture of our platform is shown in Fig. 3.

1) *Transmission Module*: The data traffic collected by the TMS is uploaded to the cloud computing platform through a transmission module, which provides real-time and stable transmission by using Flume [49] and Kafka [50].

2) *Storage Module*: Hadoop Distributed File System (HDFS) [51] and HBase [52] are used to store massive data traffic in the form of flow records, or packets. All files in the platform are replicated for fault tolerance. The storage space of the platform can be easily extended by adding disks or new machines.

3) *Processing Module*: We use MapReduce [53], Spark [54], and Storm [55] to process the massive data traffic. MapReduce is a programming model and an associated implementation for processing and generating large data sets. Spark supports cyclic data flow and in-memory computing, which is very efficient for iterative and matrix computation. Storm is very useful to deal with real-time analysis.

4) *Management Module*: To monitor the whole platform, we developed a management module to monitor the status of all machines, equipment, software, and modules. All monitoring data are collected by Flume and stored in a database. If the value of a monitoring item is over a set threshold, specific alarm information is sent to the administrator via a short message, an e-mail, and a web interface, immediately. In addition, we use ZooKeeper [56], to modify the configuration parameters of each machine and equipment (i.e., enabling/disabling the machines, equipment, software on machines, and modules of each software), and change the value of alarm threshold.

C. Trajectory Builder

To construct a user's history trajectories, some rules should be applied, while extracting the user's location from data traffic. Meanwhile, data noise must be removed before analyzing data.

1) *Extracting History Movements*: We extract a user's trajectories by 4-tuple {user ID, cell tower ID, time stamp, duration}.

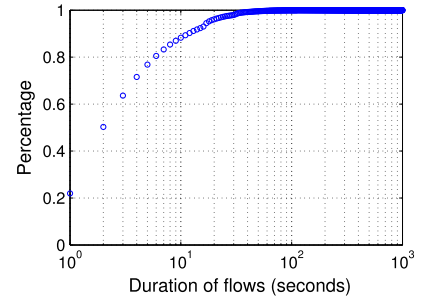


Fig. 4. CDF of flow durations for data traffic collected from 2G/3G networks.

A trajectory is constructed by a sequence of stays, and a stay is defined as

$$S = (U, L, T, D)$$

where U is the user ID; L is the cell tower ID; T is the start time of this stay, which is stored in Coordinated Universal Time (UTC); and D is the duration in seconds that a user accesses with a cell tower. Here, we have the user's trajectory as

$$Traj = \langle S_1, S_2, \dots, S_n \rangle$$

where $S_k = (U, L_k, T_k, D_k)$, $1 \leq k \leq n$.

For different data sources (flow records or packets), we apply different rules to construct the user's trajectories.

(a) *Flow records as data source*: For flow records, each flow will generate a stay, T is the start time of a flow, and D is the duration of that flow. Note that one flow only records the cell tower that a user accessed to when the flow started. If the user switches to another cell tower before the current flow ends, this transition could not be captured. To examine this deviation, we further draw cumulative distribution functions (CDFs) of flow duration of 2G/3G traffic data, i.e., we draw the value of flow duration on the x -axis and the cumulative percentage of each observed flow duration value on the y -axis. We can clearly observe in Fig. 4 that over 80% of flows last less than 6 s, and the duration of 90% of flows is less than 10 s. Since crossing over the coverage of a cell tower within 6 s is nearly impossible, we believe our data set can capture a user's movements perfectly. Note that, for those who move around the border area of a cell tower, we may not capture their next location correctly, if he/she switches the cell tower in 6 s.

Since the changing of users' locations is not recorded in flows, to reduce the deviation, we use two rules to construct a user's trajectories.

Rule 1 (merging overlapping flows with same location): If $L_k = L_{k+1}$ and $T_{k+1} < T_k + D_k < T_{k+1} + D_{k+1}$, remove $S_{a_{k+1}}$; we have $S_{a_k} = (U_a, L_k, T_k, T_{k+1} - T_k + D_{k+1})$.

If user a generates two flows at the same cell tower and the second flow starts before the first flow ends, two stays extracted from these two flows should be merged into one stay.

Rule 2 (identifying transitions from overlapping flows with different locations): If $L_k \neq L_{k+1}$ and $T_{k+1} < T_k + D_k$, then $S_{a_k} = (U_a, L_k, T_k, T_{k+1} - T_k)$, $S_{a_{k+1}} = (U_a, L_{k+1}, T_{k+1}, D_{k+1})$.

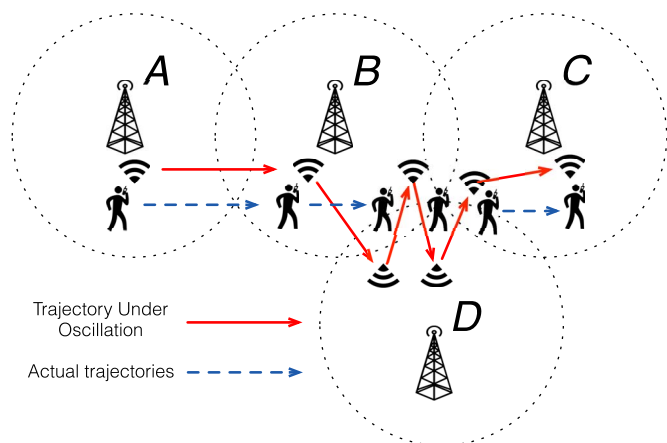


Fig. 5. Example of cell towers in mobile Internet.

If user a generates two continuous flows at different cell towers and the second flow starts before the first flow ends, the stay extracted from the first flow should end when the second flow starts.

(b) *Packets as data source*: When a UE connects to LTE network, it will be either in active state or in idle state. For different states, we apply two rules to construct the UE's trajectories.

In active state, LTE networks are aware of the ID of cell tower to which the UE is currently connected. Every time the UE switches to a cell tower, a stay is generated. For one user, there is no time interval between consecutive stays, and all movements (cell tower switching) can be captured.

Rule 3 (identifying transitions when a UE is in active state in 4G networks): In active state, for two stays generated by user a , if T_k and T_{k+1} are the time that the UE attaches to cell tower L_k and L_{k+1} , respectively, we have $T_k + D_k = T_{k+1}$.

While in idle state, the UE shall initiate the tracking area updating procedure by sending a tracking area update (TAU) request every 12 min (periodic tracking area updating is used to periodically notify the availability of the UE to the network).

Rule 4 (identifying transitions when a UE is in idle state in 4G networks): In idle state, for each user, we have a stay $S_k = (U, L_k, T_k, 0)$ for every 12 min, and movements between each TAU is lost.

2) *Reducing Oscillation*: Mobile phone may switch between different cells even when a user is not mobile. It usually happens when a user is in the overlapping area of two or more cells. This phenomenon is called "cell oscillation" or "ping-pong effect." For example, as shown in Fig. 5, if the user's real movement is $A \rightarrow B \rightarrow C$, when oscillations happened, the trajectory extracted from data traffic would be $A \rightarrow B \rightarrow D \rightarrow B \rightarrow D \rightarrow B \rightarrow C$ (this kind of oscillations is easy to identify) or $A \rightarrow B \rightarrow D \rightarrow C$ (this kind of oscillations is very hard to identify).

Some studies handle the oscillation problem by clustering cell towers, which will reduce the position accuracy [38], [57]. A recent study has proposed an algorithm framework called DECREASE (Detect, Expand, Check, Remove) [58], which resolves oscillation by selecting a cell tower to approximate the mobile device's actual location.

To reduce the oscillations effectively under a big mobile data environment, we only consider two features of cell tower oscillation: 1) It happens between adjacent cells; 2) the duration for oscillations is quite short (if the duration between switching is long, we cannot tell whether it is an oscillation or a real movement). Hence, we apply a simple method to handle the oscillation problem.

Rule 5 (reducing oscillations): Calculate the average displacement (switching) time d_t for the given data set. If a user changes location during d_t , oscillations may happen. Replace the locations that the user connects to during d_t with the one that has the longest accessing time during d_t .

Because we only capture the user's movements when his/her smartphone generates data traffic, we only take consecutive stays, between which there is no time interval, into account when calculating the average displacement time.

D. Data Storage

After removing data noise, we design two tables to store the data of users' movement trajectories and information of cell towers.

1) *Table 1 (User Trajectories)*: We store users' trajectories as 4-tuple {user ID, cell tower ID, time stamp, duration}, to draw a user's history movements from spatial and temporal dimensions.

2) *Table 2 (Cell Tower Property)*: For each cell tower, we store cell tower's property as 6-tuple {cell tower ID, network type, longitude, latitude, a list of adjacent cell towers, a list of semantic location tags}. Network type includes 2G/3G/4G networks; longitude and latitude are the geographical location of cell tower. Semantic location tag is the regional characteristic of cell tower (such as "xx shopping mall" and "xx station"), and this usually includes the name of station, commercial/residential area, educational/industrial/government building, etc. A cell tower may have many semantic location tags.

Table 1 stores the most basic information of a user's trajectory, which could be used to analyze mobility pattern, to discover similar users, or to predict the user's next location. After combining data in Tables 1 and 2, advanced semantic information can be illustrated. We could predict the user's future movement in a spatiotemporal scene and discover the most popular/hot/crowded place in a city.

E. Functionalities

Here, we introduce the methods used by four mobility functionalities, including identifying city hotspots, discovering mobility pattern, finding similar users based on history movement path, and predicting user's future movements.

1) *Identifying City Hotspots*: By collecting the mobile Internet data traffic generated by users, we could have a glance at the structure and dynamics properties of a city. In particular, it is very important to identify the "heart" of the city, which is also called "city hotspot." If there are some abnormal changes of city hotspots, it may imply that unexpected events are happening or a big event will happen soon.

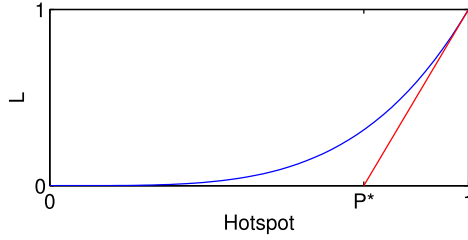


Fig. 6. Criteria selection on the Lorenz curve. The threshold corresponds to the value of $1 - P^*$.

City hotspots, i.e., the most significant locations along the human's trajectories, are made of the geographical area covered by one or many cell towers. Depending on the properties that we want to focus on, many different kinds of city hotspots could be identified. Previous work identified city hotspots in the view of density of population [36], subscribers' mobile data [37], the residence time [59], and semantic location [46]. The aforementioned properties are characterized by a specified parameter. If the parameter value of a place is bigger than a threshold, the place could be identified as a city hotspot. Therefore, identifying hotspots is an issue of exploring an efficient threshold for a specified parameter. We employ a parameter-free method proposed by Louail [36] to select the threshold. The employed method is based on the Lorenz curve.

For a specified time period, we can obtain the value of an indicator $r(j)$, such as population (number of users), from candidate hotspots. Here, j is the sequence number of a candidate hotspot. We sort candidate hotspots in an increasing order and then denote them by $r(1) < r(2) < \dots < r(n)$, where n is the total number of locations. The Lorenz curve is constructed in the following way:

- On the x -axis, draw the proportion of candidate hotspots $P = j/n$, where $j = 1, 2, \dots, n$.
- On the y -axis, plot the corresponding proportion of the number of users with

$$L(j) = \frac{\sum_{i=1}^j r(i)}{\sum_{i=1}^n r(i)}. \quad (1)$$

The method employs the natural way to identify the typical scale of the number of hotspots, which is to take the intersection point P^* between the tangent of $L(P)$ at point $P = 1$ and the horizontal axis $L = 0$ (see Fig. 6). Then, the method gives $1 - P^*$ as the threshold for identifying the hotspots. This method is inspired from the classic scale determination for an exponential decay: if the decay from $F = 1$ is an exponential of the form $e^{(F-1)/a}$, where a is the typical scale that we want to extract, this method would give $1 - F^* = a$.

2) *Discovering Mobility Pattern*: In the view of 2G/3G/4G networks, users' history movements are made of a series of stays, i.e., $Traj = \langle S_1, S_2, \dots, S_n \rangle$, where $S_k = (U, L_k, T_k, D_k)$, $1 \leq k \leq n$. Discovering the mobility pattern of individuals or groups is the matter of finding a maximal continuous trajectory. Therefore, we apply a modified version of the Apriori algorithm to discover maximal sequential pattern. A mobility pattern is identified only if the support value of discovered maximal continuous trajectory is larger than the support

TABLE I
NOTATION DESCRIPTION OF MODIFIED APRIORI ALGORITHM

Notation	Description
δ	Minimum support threshold
$Traj_S$	Set of mobility trajectories
H	Set of hotspots
C_k	Set of $length(k)$ candidate mobility patterns
P_k	Set of $length(k)$ mobility patterns
P	Set of all mobility patterns

threshold [60]. In our case, for a set of trajectories $Traj_S = \{Traj_1, Traj_2, \dots, Traj_N\}$, the support value of pattern p is the ratio of the number of pattern p that appeared in trajectories to the number of trajectories, which is defined as

$$supp(p) = \frac{|\{Traj_i | p \subset Traj_i \text{ and } 1 \leq i \leq N\}|}{N}. \quad (2)$$

For example, if a $Traj_S$ contains 10 trajectories and 6 trajectories contain the mobility pattern p , $supp(p)$ equals 0.6. Given a minimum support threshold δ , the mobility pattern p is defined as a mobility pattern if and only if p has a support value satisfying $supp(p) \geq \delta$.

In this paper, we discover the mobility pattern only between city hotspots. The main notations used in our method are listed in Table I, and the pseudocode of our algorithm is shown in Algorithm 1. A mobility pattern $p = \langle a_1, a_2, \dots, a_n \rangle$ is a candidate mobility pattern only if its subpattern $q = \langle a_1, a_2, \dots, a_{n-1} \rangle$ is discovered as a mobility pattern. For example, if $q = \langle a, b, c \rangle$ is a mobility pattern, $p = \langle a, b, c, d \rangle$ is a candidate mobility pattern. The main idea of this algorithm is discovering a continuous trajectory, the support value of which is larger than δ . We first calculate each hotspot's support value, and the set of $length(1)$ mobility patterns are generated. Then, the mobility patterns with $length(k)$ are generated through mobility patterns with $length(k-1)$. The iteration is ended when the set of $length(k)$ is \emptyset .

Algorithm 1 Discovering Mobility Patterns

INPUT: Support threshold δ

Set of mobility trajectories T

Set of hotspots H

OUTPUT: Set of mobility patterns P

1: $k = 1$

2: $C_k = \{h | h \in H\}$

3: $P_k = \{h | h \in H \wedge supp(h) > \delta\}$

4: $P = \{ \}$

5: **repeat**

6: $k = k + 1$

7: **for all** mobility pattern $p_{k-1} \in P_{k-1}$ **do**

8: **for all** frequent pattern $p_1 \in P_1$ **do**

9: $C_k = \{c_k | c_k = p_{k-1} \cup p_1\}$

10: **end for**

11: **end for**

12: **for all** trajectory $Traj_s \in Traj_S$ **do**

13: $C_t = \text{subset}(C_k, Traj_s)$

14: **for all** candidate $c \in C_t$ **do**

15: $\text{count}(c) = \text{count}(c) + 1$

```

16:   end for
17: end for
18:  $P_k = \{c | c \in C_k \wedge \text{supp}(c) > \delta\}$ 
19:  $P = \cup P_k$ 
20: until  $P_k = \emptyset$ 
21: return  $P$ 

```

3) *Finding Similar Users Based on Path*: The history movements of users may reflect their relationship. If two different users have similar moving path every day, they may know each other or have the potential to be friends. The more unpopular locations (the locations that people seldom visit) that they visit at the same time interval, the more likely that they share similar interests. In [61], user similarity is mined, based on GPS data collected from mobile phones, to recommend friends or discover a community. According to the features of our data set, users' trajectories are extracted even when users are in a big shopping mall or subway (GPS signal is not available in indoor places, underground, and the area of intensive buildings), which ensure that users' daily movements in urban areas are captured.

First, we apply the improved Apriori algorithm to find the maximum similar sequence (MSS) from two users' moving paths (path_1 and path_2).

Second, we calculate the "Inverse Document Frequency (IDF)" [62] for all the locations

$$\text{idf}(s) = \log \frac{N}{n} \quad (3)$$

where N is the total number of users in the data set, and n is the number of users visiting locations s . That is to say, if a lot of people visit location s , the value of $\text{idf}(s)$ will be very small.

Third, the IDF of the i th MSS for two paths are calculated as

$$\text{IDF}(MSS_i) = 2^{|MSS_i|-1} \times \sum_{i=1}^{|MSS_i|} \text{idf}(s_i). \quad (4)$$

Here, $|MSS_i|$ refers to the number of locations in the i th MSS.

Finally, we have the "Similar Score" for two paths as follows:

$$\text{SimScore}(\text{path}_1, \text{path}_2) = \frac{\sum_{j=1}^m \text{IDF}(MSS_j)}{|\text{path}_1| \times |\text{path}_2|} \quad (5)$$

where m is the number of MSS for path_1 and path_2 , and $|\text{path}_1|$ is the number of distinct locations in path_1 [61].

Based on the different time intervals in a day that we focus on, an existing or potential relationship would be found. For example, paths of colleagues or family members tend to have high Similar Score during work time or night, respectively. If two paths collected from two different phone numbers have very high Similar Score during weeks, we may assume that these two phone numbers belong to the same person.

4) *Predicting User's Future Location*: Predicting users' future positions allows us to be ready for their movement and to react in advance. To identify user groups according to their temporal and spatial characteristics, we discretize a day into 24 time segments, i.e., each segment lasts an hour long, as shown in Table II.

TABLE II
TIME SEGMENT OF THE CORRESPONDING TIME INTERVAL

Time segment	Time interval
0	0:00 ~ 0:59
1	1:00~1:59
2	2:00~2:59
...	...
22	22:00~22:59
23	23:00~23:59

We use entropy to measure the activity of users and capture the degree of predictability, which is defined as follows:

$$H(X) = \sum_{i=1}^n (p(x_i)I(x_i)) = - \sum_{i=1}^n p(x_i) \log_b^{p(x_i)} \quad (6)$$

where n is the number of different locations that a user visited in one time segment, i represents the location index that the user visited, and b equals to 2. $p(x_i)$ is the probability of a user staying in a certain place in one time segment. The bigger the entropy value, the more locations that the user visits in the current time segment. For each user, we build two entropy vectors: the entropy value for each segment in weekdays and on weekends, respectively, i.e.,

$$\begin{aligned} E_{\text{weekday}} &= [e_{\text{weekday}/\text{segment}_0}, \dots, e_{\text{weekday}/\text{segment}_{23}}] \\ E_{\text{weekend}} &= [e_{\text{weekend}/\text{segment}_0}, \dots, e_{\text{weekend}/\text{segment}_{23}}]. \end{aligned}$$

Considering the correlation between location and time, we group users into two groups (group 1 and group 2) with distinct mobility patterns by clustering them with k -means clustering.

(a) *Group 1: Users who have regular repeated patterns of movements*: For those who visit very limited locations every day and follow regular pattern in different days, such as white collar workers who usually go to work around 8:00 A.M. and go home around 6:00 P.M., as shown in Fig. 7, we apply "Intelligent Time Divisions (ITD)" [63] to predict not only their future movement but also the time that they may arrive. The method ITD takes spatial probability distribution as a significant characteristic to predict a user's future mobility pattern. The spatial probability distribution of a user shows the probability that a user is at a specific point in the space and is defined as

$$P(X, Y) = \text{prob}(x = X, y = Y)$$

where (x, y) represents the location of a user. If we take time factor into consideration, we can define the spatial probability distribution as

$$P_t(X, Y) = \text{prob}(x(t) = X \ \& \ y(t) = Y).$$

If we can get a user's history movements, the spatial probability distribution can be easily estimated.

(b) *Group 2: Users who move randomly*: For those who spend much time moving around the city every day, such as postmen and taxi drivers, as shown in Fig. 8, we use a time-based Markov predictor to predict their next location. Although they travel relative randomly in the city, still, some patterns may be discovered due to personal habits, traffic conditions, and road planning in the city. For users in group 2, we predict not

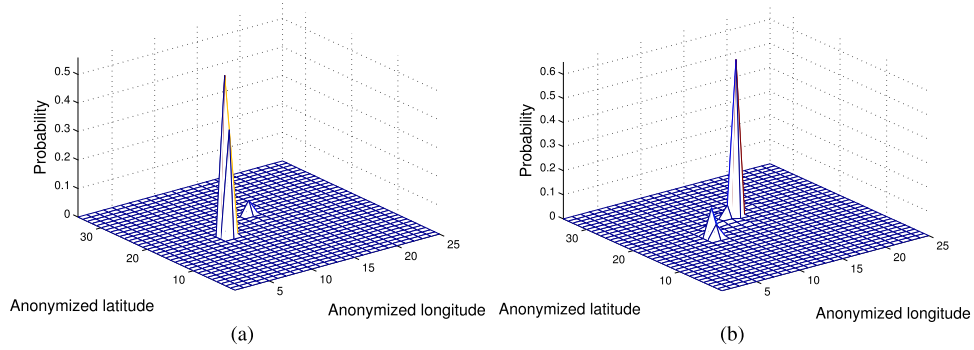


Fig. 7. Spatial probability distribution of a user with regular pattern of movements for workdays. The user basically stays in one place (probably his/her working place) during working hours and stays in another place (probably his/her home) during non-working hours. (a) Working hours. (b) Non-working hours.

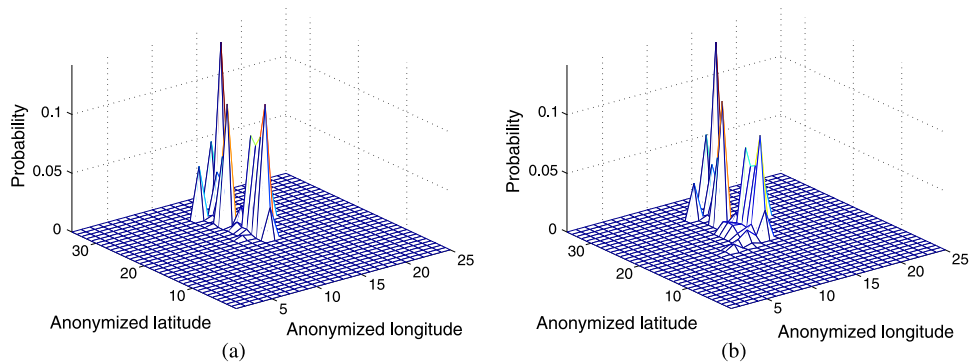


Fig. 8. Spatial probability distribution of a user who moves randomly on workdays. The user moves among many places during working or non-working hours. (a) Working hours. (b) Non-working hours.

only the user's next location but the time interval that he/she stays in this location as well.

For a trajectory $Traj = \langle S_1, S_2, \dots, S_n \rangle$, to predict the next location L_{n+1} , we find all $S_k = (U, L_k, T_k, D_k)$, $1 < k < n$, satisfying $L_k = L_n$, and $getHours(T_n) - T_{inter} < getHours(T_k) < getHours(T_n) + T_{inter}$. Here, the value of T_{inter} depends on the length of the time interval that we predict that he/she may stay at the next location, and $getHours(T)$ equals to current time (24-h format), since T is the start time of stay S , which is stored as UTC time. Finally, the location L_{k+1} that meets the aforementioned conditions and appears most frequently is the predicted next location. For example, a user just passes cell tower A at 9:00 A.M., and we want to know his/her next location, if we have $T_{inter} = 1$, we find all the cell tower A that the user passed by during 8:00 (an hour before 9:00) and 10:00 (an hour after 9:00) in his/her history trajectories. Then, the next cell tower that the user connected to after cell tower A with the highest probability of occurrence in his/her history trajectories is the prediction result.

IV. USER MOBILITY BEHAVIOR

To test the effectiveness of the framework, we collected two data sets from real mobile Internet to do the experiments. In this part, first, we illustrate the basic characteristics of our data sets. Second, three human mobility indicators are calculated to show the mobility feature of our data sets. Then, the experimental results of mobility functionalities are introduced.

A. Data Set

We collected 2G/3G/4G data traffic from real mobile Internet, as shown in Table III. The 2G/3G data traffic is extracted as flow records from July 25 to 31, 2015, which covers nearly seven million people of a big northern city in China. The 4G data traffic is the control-plane packets from October 10 to 31, 2013, with over 3000 people in a big city in southern China.

In our experiments, we use 2G/3G data traffic, which covers nearly seven million people but lasts only seven days, to study the mobility features of large-scale human mobility, and detect the hotspots with large population in the city. The 4G data traffic, which captures the trajectories of thousands of people in 21 days, is more suitable to discover the mobility patterns, to find similar users based on path, and to predict user's future movements.

B. Mobility Features

Nowadays, large-scale human mobility is described by three widely accepted indicators: the trip distance distribution $p(r)$, the radius of gyration $r_g(t)$, and the number of visited locations over time $S(t)$ [64]. These three measures contain the basic ingredients to describe the individual trajectories, in which frequent travels occur among a limited number of places, with less frequent trips to new places outside each individual radius.

1) *Trip Distance Distribution $p(r)$* : The trip distance distribution $p(r)$ quantifies the relative probability of finding a displacement of length r in a short time. By analyzing the

TABLE III
CHARACTERISTICS OF DATA SOURCES

Data sources	Networks	Duration	Number of users	Number of cell towers	Number of flows/packets
Flow records	2G/3G	7/25/2015-7/31/2015 (7 days)	6.90×10^6	85453	2.80×10^{10}
Packets	4G	10/10/2013-31/10/2013 (21 days)	3474	2252	3.60×10^7

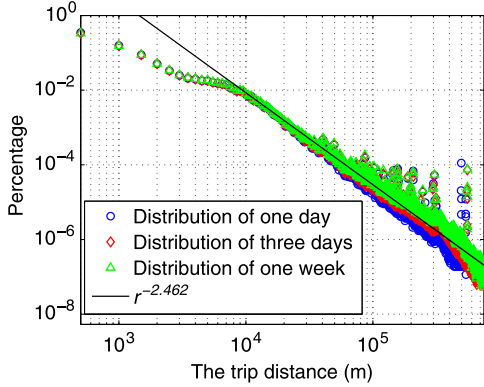


Fig. 9. Trip distance distribution $p(r)$.

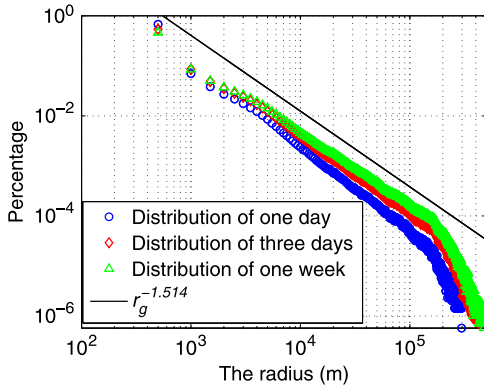


Fig. 10. $p(r_g)$ distribution of the radius of gyration r_g for users.

circulation of bank notes in the United States, a previous study [4] found that the distribution of $p(r)$ decays as a power law, i.e., $p(r) \sim r^{-\beta}$ with $\beta \approx 1.59$. In our case, as shown in Fig. 9, $p(r)$ follows power law with $\beta \approx 2.462$, which implies that the proportion of large trip distance of bank note trajectories is bigger than our data sources. It is because the data source in [4] covers the nationwide area (United States), but our data sources covers a tier-2 city in China.

2) *Radius of Gyration r_g* : The r_g reveals how extensively users move rather than capture the practical distance. Visiting the same sequence of locations in a circle continuously does not increase the value of radius of gyration, whereas a straight-line movement does [6]. r_g is defined as follows:

$$r_g = \sqrt{\frac{1}{n} \sum_{i=1}^n (\vec{r}_i - \vec{r}_{cm})^2} \quad (7)$$

where \vec{r}_i is as i th location in a user's history trajectories, $i = (1, 2, \dots, n)$, and $\vec{r}_{cm} = (1/n) \sum_{i=1}^n \vec{r}_i$ is the center of a trajectory. As shown in Fig. 10, the distribution of $p(r_g)$ for users in seven days follows power law $p(r_g) \sim r_g^{-\beta}$ with $\beta \approx 1.514$. In [35], power law is observed with $\beta \approx 1.55$ for CDRs

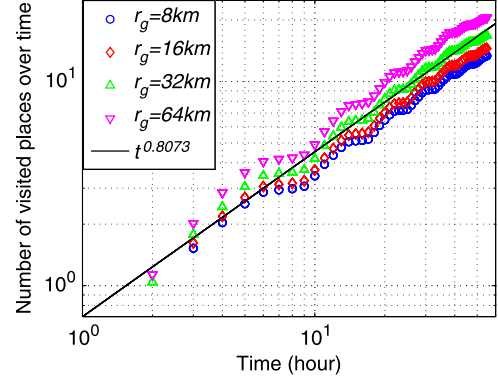


Fig. 11. Number of visited distinct locations $S(t)$ versus time.

of 3 million users in one year. It implies that the movement range of mobile phone users in [35] is smaller than that of users in our data traffic collected from 2G/3G networks. A reasonable guess is that the geographical scope covered by our data set (about 53 840 km²) is bigger than that covered by the data set in [35].

3) *Number of Visited Distinct Locations Over Time $S(t)$* : The number of visited distinct locations over time describes how frequently a user visits new places, which is expected to follow $S(t) \sim t^\mu$ (see Fig. 11). $\mu < 1$ indicates a slowdown at large-time cases, which implies a decreasing tendency of the user to visit previously unvisited locations. In our case, $S(t)$ grows as t^μ with $\mu = 0.807$.

In summary, the aforementioned results show that, for the same indicator, different data sources follow similar models with different values of parameters. We can conclude that human trajectories extracting from cell towers accessing records are able to capture the basic characteristics of human movements.

C. Hotspot Detection

In this part, we will identify the hotspots from the view of population. Intuitively, the population at a place is directly proportional to the importance that is attributed to it by the users. Places with large population, such as a big shopping mall, a residential area, a traffic hub, or the places for group activities, have significance for the city. However, hotspots always change with time, which shows the movements of population in different regions in the city, as shown in Fig. 12. We detect hotspots for each hour in one day. Over 24.4% of hotspots appear once in a day, and only 7.5% of hotspots last more than 12 hours (like 24-hour eating areas, traffic hubs, and universities).

D. Mobility Pattern of Individuals and Groups

Understanding the mobility pattern of groups in the city reveals the population stream among specific locations at a

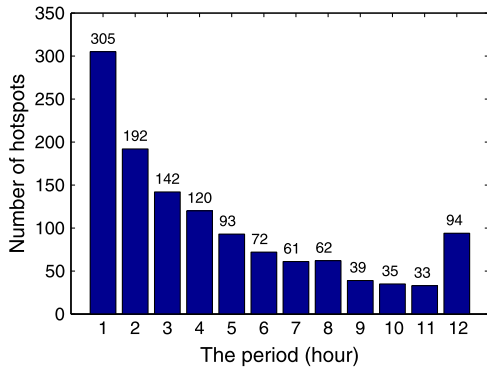


Fig. 12. Duration distribution of hotspots in one day.

TABLE IV
TOP FIVE MOBILITY PATTERNS OF GROUPS IN THE CITY

Pattern number	Pattern	Support value
1	(transportation hub, residential area)	0.049
2	(residential area, transportation hub)	0.037
3	(government building, economic center, university)	0.034
4	(university, government building, economic center)	0.033
5	(megamall, mall, food street, residential area)	0.032

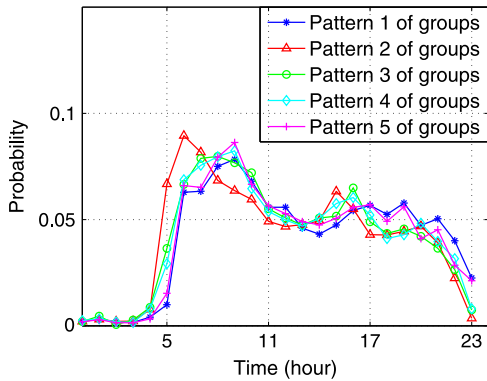


Fig. 13. Probability of the occurrences of patterns for groups varies with time.

specified time, which has important practical applications to making better urban planning, such as set new bus or subway routes and increase cell towers in a more efficient way. Among all the mobility patterns for groups that we evaluated, five patterns are most common in the city, as shown in Table IV.

In our experiments, the support threshold is set to 0.01 to make sure that the average mobility pattern length achieves the longest. The most frequent patterns are a roundtrip between a transportation hub and a residential area. It implies that there are a lot of people living in this residential area and that they usually go to a transportation hub when they need to travel in the city. As shown in Fig. 13, all patterns in Table IV start to occur at 4:00 A.M. Most of them happen between 6:00 A.M. and 9:00 A.M. in the morning, or between 3:00 P.M. and 7:00 P.M. in the afternoon. It indicates that some popular patterns appear during commuting time in the city.

In addition, finding individual mobility patterns provides important information about personal habit and interest, which is of practical significance for a service provider (SP), particu-

TABLE V
TOP FOUR MOBILITY PATTERNS OF USER x IN THE CITY

Pattern number	Pattern	Support value
1	(residential area, road 1, road 2, century mansion)	0.500
2	(century mansion, information mansion)	0.909
3	(information mansion, century mansion)	0.909
4	(century mansion, road 2, road 1, residential area)	0.519

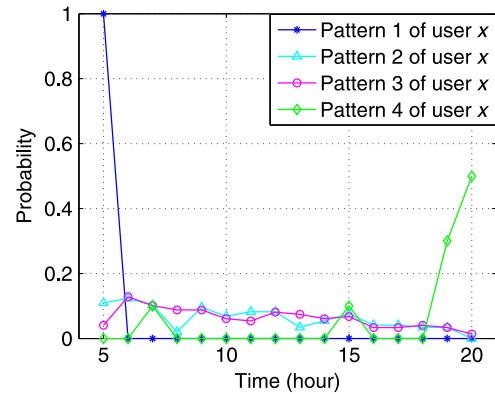


Fig. 14. Probability of the occurrence of patterns for user x varies with time.

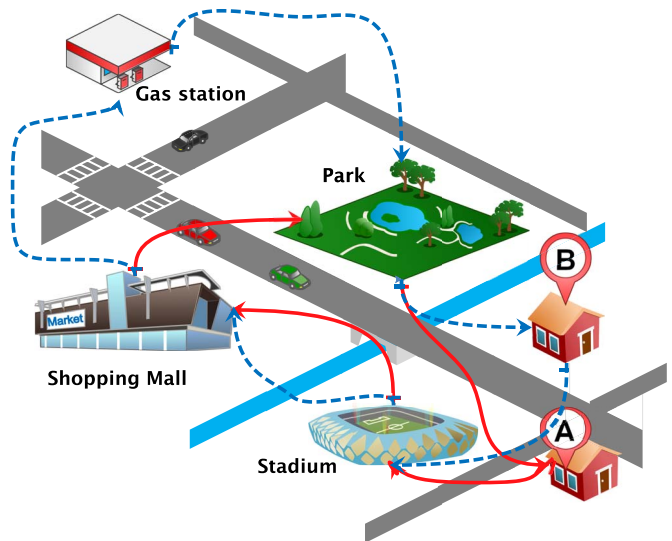


Fig. 15. Semantic path of two sampled users.

larly for location-based SP. We apply the same algorithm to a sampled user x and discover his/her daily patterns, as shown in Table V. The support threshold is set to 0.5, which means that supported patterns emerge at least 10 days over 21 days.

As shown in Fig. 14, pattern 1 happened around 5:00 A.M. for most of the time, and pattern 4 usually happened between 7:00 P.M. and 8:00 P.M. Apparently, pattern 1 and pattern 4 are commute routes for user x . In addition, the century mansion and information mansion (the office building) may be his/her workplaces, and the residential area is his/her home.

E. Relationship Among Users

A user's history trajectories uncover his/her daily interests and living habits. For example, sports fans tend to visit the gym

TABLE VI
ACCURACY OF PREDICTION ALGORITHMS FOR APPLYING DIFFERENT ALGORITHMS TO GROUPS WITH DISTINCT MOBILITY PATTERN

Groups	Number of users	Prediction accuracy of Intelligent Time Divisions	Prediction accuracy of time-based Markov	Markov
Users with regular pattern of movements	795	75.2%	22.5%	22.8%
Users move randomly	1409	26.3%	48.6%	46.9%
All users	2204	43.9%	39.2%	38.2%

and stadium more often, and fashion girls like going shopping during leisure time. The similarity of historical moving path among users draws a potential or existed relationship. The more the users come to places that other users seldom visit, the closer they tend to be. It brings new strategy for SPs to find target users and even helps authorities to locate suspicious people who may have close relationship with target persons.

By applying the method in Section III (Finding Similar Users Based on Path), the Similar Score value of every two users is calculated. Fig. 15 shows the semantic path of two sampled users that get 0.87 normalized Similar Score during a period of time in the weekend. User *A* and user *B* came to a stadium from two different places (probably their home). After a while, they went to the same shopping mall. Then, user *B* visited a filling station and appeared in a park where *A* had already arrived. After about an hour, they returned respectively to their “home.” We can easily conclude that user *A* and user *B* may be friends or have same interests.

F. Mobility Prediction

By applying ITD [63], time-based Markov, and Markov to different groups of users, we have experimental results, as shown in Table VI. We use data traffic of 4G networks, which captures users’ trajectories in 21 days. We select 2204 users from our data, who generate more than 500 packets in 21 days (if user generates very few packets, the movements extracted from packets are not enough to do the experiments). By quantifying the activity of users, we identify 795 users with regular pattern of movements and 1409 users as randomly moving people. For each trajectory, we use the previous $n(n > 0)$ continuous movements to predict the $n + 1$ movements. Then, the prediction accuracy of each user is the proportion of correct predictions for his/her trajectories. The average value of prediction accuracy of all users is the prediction results.

We predict not only a user’s next location but also the time that he/she will arrive, and we achieve better prediction accuracy than the benchmark (Markov algorithm), as shown in Table VI. ITD beats time-based Markov, while predicting the next locations for users with regular pattern of movements. For predicting the future movements of randomly moving users, time-based Markov achieves the highest prediction accuracy when compared with ITD and Markov. It implies that we should employ different prediction algorithms for distinct groups with different movement patterns.

V. CONCLUSION

In this paper, we have proposed a framework to analyze user mobility using big mobile data in densely populated areas. The

whole framework is based on a cloud computing platform, which provides data collection, preprocessing, storage, and analysis function. We further introduce the rules for constructing users’ trajectories from different data sources; the methods for reducing data noise; and the algorithms for identifying hotspots, discovering mobility pattern of groups and individuals, finding similar users based on path, and predicting user’s future movements.

Some interesting findings have come out after the experiments. First, comparing with other studies based on trajectories extracted from bank notes, CDRs, and cellular networks, the three indicators (i.e., the distribution of trip distance, the radius of gyration, and the number of visited distinct locations over time) calculated by our data sets follow similar models with different values of parameters, which vary with the duration, the covered area, and the population of data. It indicates that users’ trajectories extracted from data traffic of mobile Internet are very suitable for analyzing users’ mobility in a big city. Second, by applying our methods, collected data traffic reveals some interesting phenomenon in the city, such as the following: more than half of the hotspots last less than 3 hours; a large crowd moving between a transportation hub and a residential area during morning and evening peak hours; users with similar interests could be easily identified from their history trajectories. All the analysis results uncover the common rules that exist among huge populations in a city, which are of theoretical and practical significance for urban planning, traffic control, mobile network resource optimization, etc. Third, people in the city usually have distinct mobility patterns. Considering the mobility pattern while predicting user’s future movements could improve prediction accuracy.

In the future, we will apply the data stream algorithms and get the real-time analysis result. In this way, we could make some applications more practical, such as predicting and monitoring large-scale events. In addition, for each mobility application, applying the most suitable algorithm to our data set and improving existing methods are our future work.

REFERENCES

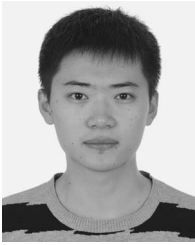
- [1] A. Noulas *et al.*, “A tale of many cities: Universal patterns in human urban mobility,” *PloS One*, vol. 7, no. 5, May 2012, Art. no. e37027.
- [2] F. Giannotti, L. Pappalardo, D. Pedreschi, and D. Wang, “A complexity science perspective on human mobility,” in *Mobility Data: Model, Manage., Understand.* New York, NY, USA: Cambridge Univ. Press, Oct. 2012, pp. 297–314.
- [3] D. Balcan *et al.*, “Multiscale mobility networks and the spatial spreading of infectious diseases,” *Natl. Acad. Sci.*, vol. 106, no. 51, pp. 21 484–21 489, Dec. 2009.
- [4] D. Brockmann, L. Hufnagel, and T. Geisel, “The scaling laws of human travel,” *Nature*, vol. 439, no. 7075, pp. 462–465, Jan. 2006.
- [5] A. L. Barabasi, “The origin of bursts and heavy tails in human dynamics,” *Nature*, vol. 435, no. 7039, pp. 207–211, May 2005.

- [6] M. C. Gonzalez, C. A. Hidalgo, and A. L. Barabasi, "Understanding individual human mobility patterns," *Nature*, vol. 453, no. 7196, pp. 779–782, Jun. 2008.
- [7] F. Simini, M. C. González, A. Maritan, and A. L. Barabási, "A universal model for mobility and migration patterns," *Nature*, vol. 484, no. 7392, pp. 96–100, Apr. 2012.
- [8] C. Song, Z. Qu, N. Blumm, and A. L. Barabási, "Limits of predictability in human mobility," *Science*, vol. 327, no. 5968, pp. 1018–1021, Feb. 2010.
- [9] X. Lu, E. Wetter, N. Bharti, A. J. Tatem, and L. Bengtsson, "Approaching the limit of predictability in human mobility," *Scient. Rep.*, vol. 3, no. 2923, pp. 1–9, Oct. 2013.
- [10] I. Rhee *et al.*, "On the Levy-walk nature of human mobility," *IEEE/ACM Trans. Netw.*, vol. 19, no. 3, pp. 630–643, Jun. 2011.
- [11] D. Aldous and J. Fill, Reversible Markov chains and random walks on graphs, 2002. [Online]. Available: <http://www.stat.berkeley.edu/~aldous/RWG/book.html>
- [12] B. Liang and Z. J. Haas, "Predictive distance-based mobility management for PCS networks," in *Proc. IEEE 18th Ann. Joint Conf. INFOCOM*, New York, NY, USA, Mar. 1999, pp. 1377–1384.
- [13] C. Bettstetter, "Smooth is better than sharp: A random mobility model for simulation of wireless networks," in *Proc. 4th ACM Int. Workshop Model*, New York, NY, USA, Jul. 2001, pp. 19–27.
- [14] C.-C. Chiang, "Wireless network multicasting," Ph.D. dissertation, Univ. California, Los Angeles, CA, USA, 1998.
- [15] V. A. Davies, "Evaluating mobility models within an ad hoc network," M.S. thesis, Colorado Sch. Mines, Golden, CO, USA, 2000.
- [16] A. Jardosh, E. M. Belding-Royer, K. C. Almeroth, and S. Suri, "Towards realistic mobility models for mobile ad hoc networks," in *Proc. 9th Ann. Int. Conf. Mobile Comput. Netw.*, New York, NY, USA, Sep. 2003, pp. 217–229.
- [17] S. Batabyal and P. Bhanuik, "Mobility models, traces and impact of mobility on opportunistic routing algorithms: A survey," *IEEE Commun. Surveys Tuts.*, vol. 17, no. 3, pp. 1679–1707, Apr. 2015.
- [18] H. Nishiyama, M. Ito, and N. Kato, "Relay-by-smartphone: Realizing multihop device-to-device communications," *IEEE Commun. Mag.*, vol. 52, no. 4, pp. 56–65, Apr. 2014.
- [19] Y.-B. Lin *et al.*, "EasyConnect: A management system for IoT devices and its applications for interactive design and art," *IEEE Internet Things J.*, vol. 2, no. 6, pp. 551–561, Dec. 2015.
- [20] H. Cai and D. Y. Eun, "Crossing over the bounded domain: From exponential to power-law intermeeting time in mobile ad hoc networks," *IEEE/ACM Trans. Netw.*, vol. 17, no. 5, pp. 1578–1591, Oct. 2009.
- [21] J. Liu, S. Zhang, N. Kato, H. Ujikawa, and K. Suzuki, "Device-to-device communications for enhancing quality of experience in software defined multi-tier LTE-A networks," *IEEE Netw.*, vol. 29, no. 4, pp. 46–52, Jul./Aug. 2015.
- [22] K. Zheng *et al.*, "Soft-defined heterogeneous vehicular network: Architecture and challenges," *arXiv preprint arXiv:1510.06579*, vol. 30, no. 4, pp. 72–80, Jul.–Aug. 2016.
- [23] K. Zheng *et al.*, "Big data-driven optimization for mobile networks toward 5G," *IEEE Netw.*, vol. 30, no. 1, pp. 44–51, Jan./Feb. 2016.
- [24] J. Yeo, D. Kotz, and T. Henderson, "CRAWDAD: A community resource for archiving wireless data at Dartmouth," *SIGCOMM Comput. Commun. Rev.*, vol. 36, no. 2, pp. 21–22, Apr. 2006.
- [25] R. Becker *et al.*, "Human mobility characterization from cellular network data," *ACM Commun.*, vol. 56, no. 1, pp. 74–82, Jan. 2013.
- [26] Z. Smoreda, A. M. Olteanu Raimond, and T. Couronné, "Spatio-temporal data from mobile phones for personal mobility assessment," in *Transport Survey Methods: Best Practice for Decision Making*. Bingley, U.K.: Emerald Group, Nov. 2013, pp. 745–766.
- [27] B. C. Csáji *et al.*, "Exploring the mobility of mobile phone users," *Phys. A, Stat. Mech. Appl.*, vol. 392, no. 6, pp. 1459–1473, Mar. 2013.
- [28] F. Giannotti, M. Nanni, F. Pinelli, and D. Pedreschi, "Trajectory pattern mining," in *Proc. 13th ACM SIGKDD Int. Conf.*, San Jose, CA, USA, Aug. 2007, pp. 330–339.
- [29] M. Lin and W. Hsu, "Mining GPS data for mobility patterns: A survey," *Perv. Mobile Comput.*, vol. 12, pp. 1–16, Jun. 2014.
- [30] Y. Zhang, "User mobility from the view of cellular data networks," in *Proc. IEEE INFOCOM*, Toronto, ON, Canada, Apr. 2014, pp. 1348–1356.
- [31] J. Yang *et al.*, "Characterizing user behavior in mobile Internet," *IEEE Trans. Emerging Top. Comput.*, vol. 3, no. 1, pp. 95–106, Mar. 2015.
- [32] Y.-B. Lin and P.-K. Huang, "Prefetching for mobile web album," *Wireless Commun. Mobile Comput.*, vol. 16, no. 1, pp. 18–28, 2016.
- [33] A. Thiagarajan, L. Ravindranath, and H. Balakrishnan, "Accurate, low-energy trajectory mapping for mobile devices," in *Proc. NSDI*, Boston, MA, USA, Mar. 2011, pp. 267–280.
- [34] X. Zhou *et al.*, "Human mobility patterns in cellular networks," *IEEE Commun. Lett.*, vol. 17, no. 10, pp. 1877–1880, Oct. 2013.
- [35] C. Song, T. Koren, P. Wang, and A. L. Barabási, "Modelling the scaling properties of human mobility," *Nature Phys.*, vol. 6, no. 10, pp. 818–823, Sep. 2010.
- [36] T. Louail *et al.*, "From mobile phone data to the spatial structure of cities," *Sci. Rep.*, vol. 4, no. 5276, pp. 1–12, Jun. 2014.
- [37] H. Klessig, V. Suryaprakash, O. Blume, A. Fehske, and G. Fettweis, "A framework enabling spatial analysis of mobile traffic hot spots," *IEEE Wireless Commun. Lett.*, vol. 3, no. 5, pp. 537–540, Aug. 2014.
- [38] M. A. Bayir, M. Demirbas, and N. Eagle, "Mobility Profiler: A framework for discovering mobility profiles of cell phone users," *Perv. Mobile Comput.*, vol. 6, no. 4, pp. 435–454, Aug. 2010.
- [39] J. Yang *et al.*, "Global and individual mobility pattern discovery based on hotspots," in *Proc. ICC*, London, U.K., Jun. 2015, pp. 1–6.
- [40] X. Yan, X. Han, B. Wang, and T. Zhou, "Diversity of individual mobility patterns and emergence of aggregated scaling laws," *Sci. Rep.*, vol. 3, no. 2678, pp. 1–5, Sep. 2013.
- [41] "Big data challenge 2015," Telecomitalia, Rome, Italy, 2015. [Online]. Available: <https://www.telecomitalia.com/tit/en/bigdatachallenge.html>
- [42] T. Cisco, "Global mobile data traffic forecast update, 2014–2019 white paper," Cisco, San Jose, CA, USA, Cisco Public Inf., Feb. 2015.
- [43] L. Vu, Q. Do, and K. Nahrstedt, "Jyotish: A novel framework for constructing predictive model of people movement from joint Wifi/Bluetooth trace," in *Proc. IEEE Int. Conf. PerCom*, Seattle, WA, USA, Mar. 2011, pp. 54–62.
- [44] H. Cao *et al.*, "SoLoMo analytics for telco Big Data monetization," *IBM J. Res. Dev.*, vol. 58, no. 5/6, pp. 1–9, Nov. 2014.
- [45] Y. Zheng, X. Xie, and W. Ma, "GeoLife: A collaborative social networking service among user, location and trajectory," *IEEE Data Eng. Bull.*, vol. 33, no. 2, pp. 32–39, Jun. 2010.
- [46] Y. Zheng and X. Zhou, *Computing With Spatial Trajectories*. New York, NY, USA: Springer-Verlag, 2011.
- [47] J. Liu, F. Liu, and N. Ansari, "Monitoring and analyzing big traffic data of a large-scale cellular network with Hadoop," *IEEE Netw.*, vol. 28, no. 4, pp. 32–39, Jul. 2014.
- [48] T. White, *Hadoop: The Definitive Guide*. Sebastopol, CA, USA: O'Reilly, 2012.
- [49] Apache Flume. [Online]. Available: <https://flume.apache.org/>
- [50] Apache Kafka. [Online]. Available: <http://kafka.apache.org/>
- [51] K. Shvachko, H. Kuang, S. Radia, and R. Chansler, "The Hadoop distributed file system," in *Proc. IEEE 26th Symp. Mass Storage Syst. Technol.*, Washington, DC, USA, May 2010, pp. 1–10.
- [52] L. George, *HBase: The Definitive Guide*. Sebastopol, CA, USA: O'Reilly, 2011.
- [53] J. Dean and S. Ghemawat, "MapReduce: Simplified data processing on large clusters," *ACM Commun.*, vol. 51, no. 1, pp. 107–113, Jan. 2008.
- [54] Apache Spark. [Online]. Available: <https://spark.apache.org/>
- [55] Apache Storm. [Online]. Available: <https://storm.apache.org/>
- [56] P. Hunt, M. Konar, F. P. Junqueira, and B. Reed, "ZooKeeper: Wait-free coordination for Internet-scale systems," in *Proc. USENIX Annu. Technol. Conf.*, Berkeley, CA, USA, Jun. 2010, pp. 11–11.
- [57] S. A. Shad, E. Chen, and T. Bao, "Cell oscillation resolution in mobility profile building," *Int. J. Comput. Sci.*, vol. 9, no. 3, pp. 205–213, Jun. 2012.
- [58] W. Wu *et al.*, "Oscillation resolution for mobile phone cellular tower data to enable mobility modelling," in *Proc. MDM*, Brisbane, Qld., Australia, Jul. 2014, pp. 321–328.
- [59] S. Scellato, M. Musolesi, C. Mascolo, V. Latora, and A. T. Campbell, "NextPlace: A spatio-temporal prediction framework for pervasive systems," in *Pervasive Computing*. San Francisco, CA, USA: Springer-Verlag, Jun. 2011, pp. 152–169.
- [60] T. Pangning, M. Steinbach, and V. Kumar, *Introduction to Data Mining*. Boston, MA, USA: Addison-Wesley, 2006.
- [61] Y. Zheng, "Location-Based Social Networks: Users," in *Computing With Spatial Trajectories*. New York, NY, USA: Springer-Verlag, Oct. 2011, pp. 243–276.
- [62] K. Sparck Jones, "A statistical interpretation of term specificity and its application in retrieval," *J. Doc.*, vol. 28, no. 1, pp. 11–21, 1972.
- [63] S. Gatzmir Motahari, H. Zang, and P. Reuther, "Time-clustering-based place prediction for wireless subscribers," *IEEE/ACM Trans. Netw.*, vol. 21, no. 5, pp. 1436–1446, Dec. 2013.
- [64] C. M. Schneider, V. Belik, T. Couronné, Z. Smoreda, and M. C. González, "Unravelling daily human mobility motifs," *J. Roy. Soc. Interface*, vol. 10, no. 84, p. 38, Aug. 2013.



Yuanyuan Qiao (M'15) received the B.E. degree from Xidian University, Xi'an, China, in 2009 and the Ph.D. degree from the Beijing University of Posts and Telecommunications (BUPT), Beijing, China, in 2014.

He is currently a Lecturer with the School of Information and Communication Engineering, BUPT. Her research focuses on traffic measurement and classification, mobile Internet traffic analysis, and big data analytics.



Yihang Cheng received the B.E. degree in communication engineering in 2014 from the Beijing University of Posts and Telecommunications, Beijing, China, where he is currently working toward the M.E. degree with the School of Information and Communication Engineering.

He is engaged in the research of users' relationship discovery and mobility prediction.



Jie Yang received the B.E., M.E., and Ph.D. degrees from the Beijing University of Posts and Telecommunications (BUPT), Beijing, China in 1993, 1999, and 2007, respectively.

She is currently a Professor and Deputy Dean of the School of Information and Communication Engineering, BUPT. Her current research interests include broadband network traffic monitoring, user behavior analysis, big data analysis in Internet and telecommunications, etc. She has published several papers on international magazines and conferences,

including the *IEEE JOURNAL ON SELECTED AREAS IN COMMUNICATIONS*, the *IEEE TRANSACTIONS ON WIRELESS COMMUNICATIONS*, and the *IEEE TRANSACTIONS ON PARALLEL AND DISTRIBUTED SYSTEMS*.

Dr. Yang was the Vice Program Committee Cochair of the IEEE International Conference on Network Infrastructure and Digital Content in 2012 and 2014.



Jiajia Liu (SM'15) is currently a Full Professor with the School of Cyber Engineering, Xidian University, Xi'an, China. His research interests include wireless and mobile *ad hoc* networks, performance modeling and evaluation, stochastic network optimization, Long-Term Evolution Advanced (LTE-A), and fifth-generation (5G) networks.

Prof. Liu received the Yasujiro Niwa Outstanding Paper Award in 2012 and the Best Paper Awards from the IEEE Wireless Communications and Networking Conference in 2012 and 2014. He also re-

ceived the Chinese Government Award for Outstanding Ph.D. Students Abroad in 2011, the Tohoku University Research Institute of Electrical Communication Student Award, and the Tohoku University Professor Genkuro Fujino Award in 2012, as well as the prestigious Dean Award and President Award of Tohoku University in 2013.



Nei Kato (F'13) is a Professor with the Graduate School of Information Sciences, Tohoku University, Sendai, Japan.

Prof. Kato is a Fellow of the Institute of Electronics, Information and Communication Engineers. He is currently the Chair of the IEEE Ad Hoc and Sensor Networks Technical Committee, the Chair of the IEEE Communications Society Sendai Chapter, an Associate Editor-in-Chief of the *IEEE TRANSACTIONS ON PARALLEL AND DISTRIBUTED SYSTEMS* and the *IEEE INTERNET OF THINGS JOURNAL*, an Area

Editor of the *IEEE TRANSACTIONS ON VEHICULAR TECHNOLOGY*, and an Editor of *IEEE WIRELESS COMMUNICATIONS* and *IEEE NETWORK*. He is a Strategic Adviser to the President of Tohoku University.