

Transportation Research Record

Quantifying Privacy Vulnerability Under Linkage Attack Across Multi-source Individual Mobility Data

--Manuscript Draft--

Full Title:	Quantifying Privacy Vulnerability Under Linkage Attack Across Multi-source Individual Mobility Data
Abstract:	<p>With the advances in detector and sensor technologies, identity detection-based intelligent transportation systems---such as license plate recognition (LPR) system and parking electronic toll collection (ETC) system---have been widely deployed in urban transportation, generating large quantities of multi-source individual-based mobility data set (e.g., LPR data and parking data). Given the high frequency, precision and wide coverage, these individual-based mobility data can be used in many transportation research areas, such as transportation planning, traffic prediction and individual mobility pattern profiling. With the increasing demand for publishing and sharing these individual-based data sets to researchers and practitioners, the privacy issue of data publishing has been a major concern since true identities of individuals can be revealed by linkage attack. In this paper, we quantitatively measure the privacy disclosure risk caused by linkage attack across multi-source individual-based mobility data sets. Taking an example of LPR data and parking data, a traffic-knowledge-driven adversary model is proposed for linkage attack conducting among LPR data and parking data. Two common modes of LPR data publishing are examined and two quantitative criteria are introduced to present the risk of privacy leakage under linkage attack. The experimental results demonstrate that anonymized individual still under high risk of being linked successfully (71.63% under mode 1 and 36.55% under mode 2). This study serves as a wake-up call for relevant agencies and data owners about the privacy vulnerability caused by linkage attack across multi-source individual-based mobility data.</p>
Manuscript Classifications:	Data and Information Technology; Urban Transportation Big Data ABJ30SC; Big Data; Data Analysis; Data Fusion; Statistical Methods ABJ80; Modeling
Manuscript Number:	20-01384
Article Type:	Presentation
Order of Authors:	Jing Gao
	Qinglong Lu
	Ming Cai

1 **QUANTIFYING PRIVACY VULNERABILITY UNDER LINKAGE ATTACK ACROSS**
2 **MULTI-SOURCE INDIVIDUAL MOBILITY DATA**

3 **Jing Gao**

4 Guangdong Provincial Key Laboratory of Intelligent Transportation Systems
5 Research Center of Intelligent Transportation System
6 &
7 School of Intelligent Systems Engineering
8 Sun Yat-Sen University
9 Guangzhou, Guangdong 510006, China
10 Email: gaoj26@gmail.com

11 **Qinglong Lu**

12 Department of Civil, Geo and Environmental Engineering
13 Technical University of Munich
14 Arcisstra. 21, 80333 Munich
15 Email: qinglong.lu@tum.de

16 **Ming Cai, Corresponding Author**

17 Guangdong Provincial Key Laboratory of Intelligent Transportation Systems
18 Research Center of Intelligent Transportation System
19 &
20 School of Intelligent Systems Engineering
21 Sun Yat-Sen University
22 Guangzhou, Guangdong 510006, China
23 Email: caiming@mail.sysu.edu.cn

24 Word count: 5929 words text + 4 tables x 250 words (each) = 6929 words

25 Submission Date: August 1, 2018

1 **ABSTRACT**

2 With the advances in detector and sensor technologies, identity detection-based intelligent trans-
3 portation systems—such as license plate recognition (LPR) system and parking electronic toll
4 collection (ETC) system—have been widely deployed in urban transportation, generating large
5 quantities of multi-source individual-based mobility data set (e.g., LPR data and parking data).
6 Given the high frequency, precision and wide coverage, these individual-based mobility data can
7 be used in many transportation research areas, such as transportation planning, traffic prediction
8 and individual mobility pattern profiling. With the increasing demand for publishing and sharing
9 these individual-based data sets to researchers and practitioners, the privacy issue of data publish-
10 ing has been a major concern since true identities of individuals can be revealed by linkage attack.
11 In this paper, we quantitatively measure the privacy disclosure risk caused by linkage attack across
12 multi-source individual-based mobility data sets. Taking an example of LPR data and parking
13 data, a traffic-knowledge-driven adversary model is proposed for linkage attack conducting among
14 LPR data and parking data. Two common modes of LPR data publishing are examined and two
15 quantitative criteria are introduced to present the risk of privacy leakage under linkage attack. The
16 experimental results demonstrate that anonymized individual still under high risk of being linked
17 successfully (71.63% under mode 1 and 36.55% under mode 2). This study serves as a wake-up
18 call for relevant agencies and data owners about the privacy vulnerability caused by linkage attack
19 across multi-source individual-based mobility data.

20 *Keywords:* Individual-based mobility data, Linkage attack, License plate recognition data, Parking
21 data, Disclosure risk

1 INTRODUCTION

2 With the advances in mobile sensing technologies, individual behaviors are widely captured and
3 recorded. Large quantities of individual-based mobility data sets, such as license plate recognition
4 (LPR) data and parking data, have been widely generated and collected. License plate recognition
5 (LPR) data, generated by the LPR system which takes pictures of every passing vehicle and con-
6 vert images to detailed spatiotemporal records automatically, capturing vehicle history trajectory
7 with high precision and wide coverage. Given the high frequency, great precision, and extensive
8 coverage, LPR data, have been applied to a wide range of transportation research areas, for exam-
9 ple, transportation planning, traffic prediction and individual mobility pattern recognition. With
10 the wide usage of LPR data, the focus has been increasingly centered on the privacy issue of pub-
11 lishing this kind of data. Previous studies adopted the concept of anonymity (*1*) and examined that
12 individuals in anonymized LPR data set still have a high risk of being re-identified (*2, 3*). Specif-
13 ically, Gao et al. (*2*) selected several spatiotemporal records in one's history trajectory to form his
14 quasi-identifier and the anonymity of an individual was defined as the number of occurrence of
15 the quasi-identifier in all individual's history trajectory. An anonymized individual is called re-
16 identified if his anonymity is one, i.e., his quasi-identifier is only contained in his own trajectory.
17 However, an individual with his anonymity equaling 1 doesn't mean that his true identity (vehicle
18 plate number in our case) would be revealed. To get the true identity of the target for obtaining
19 his privacy information, the further linkage should be made between anonymized data and external
20 information (with true identity), which is called linkage attack.

21 This paper aims at quantifying the privacy disclosure risk under the linkage attack across
22 multi-source individual mobility data. Taking an example of LPR data and parking data, we want
23 to examine to what extent the true identity of anonymized individuals can be revealed by linkage
24 attack between anonymized LPR database and parking database. In doing so, we propose a traffic-
25 knowledge-driven adversary model to conduct linkage attack. Two common modes of LPR data
26 publishing are investigated and two quantitative criteria are introduced to measure the privacy
27 disclosure risk under linkage attack. To the best of our knowledge, this is the first study that
28 empirically quantifies the privacy risk of LPR data by conducting linkage attack among LPR data
29 set and parking data set from urban transportation system.

30 The remainder of this paper is structured as follows. Section 2 reviews the related research
31 in privacy risk measurement and privacy protection in individual-based mobility data. Section 3
32 introduces the LPR data set and the parking data set. In addition, two common modes of LPR data
33 publishing are introduced section 3. Then, in section 4, we propose the traffic-knowledge-driven
34 adversary model concerning two different modes of LPR data publishing and empirical experi-
35 ments and results analysis are presented in section 5. Finally, Section 6 presents some concluding
36 remarks and proposes future work of this study.

37 LITERATURE REVIEW

38 With the emerging data in the traffic area and more and more data are available to the public, the
39 latent insecurity of privacy has attracted the attention of many researchers. To prevent attackers
40 from extracting private information from anonymized database, data suppression techniques are
41 presented in (*4*) to reduce such reidentification risks. Concerning the privacy protection of Location
42 Based Service (LBS) users, (*5*) investigate an adversary model to attack the LBS data with long-
43 term pseudonyms and uncover insecurity of users privacy under this releasing pattern. Li et al.
44 (*6*) establish a traffic monitoring scheme to preserve driver privacy by acquiring driving data and

1 integrating it with a weighted proximity graph to filter out false reports uploaded by malicious
2 drivers to intervene the storage system.

3 In (7), 3 months of credit card records are investigated and the result shows that by only
4 leveraging four spatiotemporal points 90% of individuals could be re-identified which illustrates a
5 very serve problem in the metadata. (8) indicates that 10 pieces of side information of a victim are
6 enough to identify the trace of (s)he by a probability of 30% to 50%. Factors affecting anonymity
7 are summarized in (9), and concerning these factors, two typical techniques for improving data
8 security are discussed, i.e. suppression and generalization. To evaluate the effectiveness of the
9 techniques, an adversary model is proposed to conduct attacks. These studies, however, at an
10 angle of revealing data security at a point level, have not noticed the potential consequence of
11 linkage attack.

12 To reach a trade-off both personal privacy and traffic data utility are at a comparatively
13 high level, some studies have proposed valid systems for different applications. (10) built a de-
14 centralized architecture under which the location data of ridesharing users will be released in such
15 a way that guarantees the security of user privacy without sacrificing the service quality. In (11),
16 to protect the privacy of participants of probe vehicles who have to reveal their GPS positions to
17 the traffic monitoring department, a VTLs-based system cooperating with an associated cloaking
18 technique is provided. However, this system cannot protect traffic monitoring accuracy from any
19 active attack. Sun et al. (12) construct a VTLs zone-based system to balance the data needs for
20 general traffic research and privacy security, and to assess the effectiveness of this system, traffic-
21 knowledge-based adversary models are proposed to conduct privacy attack.

22 Those systems and architectures are always focused on the privacy problem from one single
23 data source, while in reality, different sources of traffic data are mixed and have a strong relation-
24 ship with each other. To this end, a traffic-knowledge-driven adversary model is proposed for
25 conducting linkage attack across multi-source individual-based mobility data. And the privacy
26 disclosure risk under linkage attack among LPR data and parking data is quantitatively measured.
27

28 DATA DESCRIPTION

29 License plate recognition data

30 We use LPR data collected by LPR system in Guangzhou, China on 18th July, 2018. The LPR
31 system is essentially a network of cameras that can take pictures of every passing vehicle and
32 transform the image into a detailed spatiotemporal record automatically. The information in a
33 single LPR record includes license plate number, color or the license plate, timestamp when vehicle
34 capture by the detector, detector ID (representing different camera gantries) and driving direction
35 of the vehicle. Table 1 lists the primary fields from an example LPR record. The trajectories of
36 individuals can be tracked by summarizing a series of spatiotemporal records.

37 Considering the great advantage and increasing demand of LPR data, more and more agen-
38 cies choose to publish and share the LPR data set for transportation research purpose. Taking into
39 account the privacy concerns and issues in LPR data, just like the publication of other source of
40 trajectory data set, data owners usually need to construct an anonymized LPR database for data
41 publishing. Generally, there are two common alternative modes.

TABLE 1 : Primary fields in original LPR data set

Field	Example value
Vehicle plate number	@ABC123
Passing timestamp	2018-07-18 10:59:07 (one-second resolution)
Address	Inner Ring Road, Meizhou Building (East to West)
Detector ID	17068
Drive direction	0

1 *Mode 1: full trajectory publishing*

2 The first one is the most common LPR data publishing method. When LPR data is published in
3 mode 1, the vehicle plate number of each vehicle will be replaced with a unique random identifier.
4 Each record in the anonymized LPR database is a sequence of spatiotemporal tuples (in the form
5 of $(T_i, LocR_i)$), which represents the full history trajectory of an individual/vehicle. In other words,
6 under the data publishing of mode 1, complete history trajectories with its corresponding random
7 identifiers will be published. Table 2 shows an example anonymized LPR database published in
8 mode one.

TABLE 2 : Anonymized LPR database published in mode one

Rec_ID. #	Record value
5ac0bd6239f8b9a	(288,17062) → (302,17067) → (71898,14688) → (72058,16800)
5ac0bd6239f8b9b	(308,17063) → (512,17067) → (698,17069)
...	...

9 *Mode 2: segmented trajectory publishing*

10 Compared with mode 1 in which complete vehicle trajectories are published directly, mode 2
11 splits each individual's complete trajectory into segments according to parking activity and allocate
12 a random identifier to each trajectory segment. Table 3 shows the example anonymized LPR
13 database published in mode 2. Since individuals usually park several times for specific activities
14 (e.g, shopping or working) during a day, it is quite common that the full history trajectory actually
15 consists of several segmented trajectories. Consequently, spatio-temporal information within each
16 O-to-D trip is revealed.

TABLE 3 : Anonymized LPR database published in mode two

Rec_ID. #	Record value
5ac0bd6239f8b9a1	(288,17062) → (302,17067)
5ac0bd6239f8b9a2	(71898,14688) → (72058,16800)
5ac0bd6239f8b9b	(308,17063) → (512,17067) → (698,17069)
...	...

1 Parking data

2 The parking data we used was mainly collected by the parking toll collection system in Guangzhou,
3 China on 18th July, 2018. Similar to the LPR system, the camera detectors at the parking gate
4 take pictures of entering/leaving vehicles to get the vehicle plate number, the entry time and the
5 departure time of vehicles. Table 4 lists the primary fields from an example parking record. It
6 is quite natural that the parking actions arise from specific activity (e.g., commuting or shopping)
7 and the parking data records the spatiotemporal information of origin or destination in a trip.

8 Since the parking lots are usually constructed and managed by the parking owners, the
9 parking data set is consequently an accessible database to them. In addition, to get the parking
10 data set, one only need to observe and record the access of vehicles at the parking gates. As
11 a consequence, the parking database can be regarded as a public external database for linkage
12 attack.

TABLE 4 : Primary fields in original Parking data set

Field	Example value
Parking ID	etcp_788
Parking name	GOGO Xintiandi
Address	No. 1 Middle Road, Xiaoguwai Street
Gate name	No. 1 entrance
Vehicle plate number	@ABC123
Intime	2018-07-18 11:08:07
Outtime	2018-07-18 11:32:07

13 So given the anonymized LPR database (without true identity of vehicle) published in two
14 different modes and the external public parking database (with actual vehicle plate number), the
15 adversary aim at inferring the true identity of records in anonymized LPR database by linking to
16 parking database. Figure 1 shows a simple road network with detectors and parking lots, where
17 the characters represents the locations of camera detectors and the serial numbers represent the
18 locations of parking lots. Suppose an individual drove from a to c, then parked in the parking ③
19 for shopping nearby. After shopping, he left the parking and drove from e to f. So the goal of an
20 adversary is to link the parking record in parking ③ to the LPR record (e.g., $a \rightarrow b \rightarrow c \rightarrow e \rightarrow d$
21 $\rightarrow f$ in mode 1) in anonymized LPR database to obtain the vehicle plate number of the individual.

22 TRAFFIC-KNOWLEDGE-DRIVEN ADVERSARY MODEL

23 In this section, we propose a traffic-knowledge-driven adversary model to quantify the privacy
24 disclosure risk under two different modes of publishing LPR data set. Specifically, the adversary
25 model takes into consideration both the temporal (real-time travel time information) and spatial
26 (spatial connectivity) relation to conduct linkage attack between LPR data (anonymized database)
27 and parking data (public external database). The target of our model is to infer the true identities
28 of records, i.e, the vehicle plate number, in the LPR database by linking to parking data. Two
29 different modes of LPR data publishing proposed in Section ?? are examined.

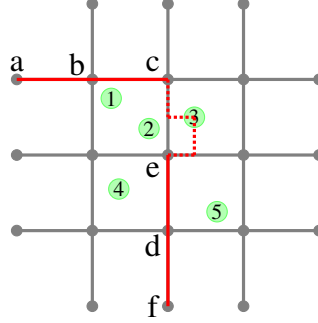


FIGURE 1 : Schematic illustration of road network with camera detectors and parking

1 Linking under full trajectory publishing (Mode 1)

2 As presented in Section 3, in LPR data publishing of mode 1, LPR data is anonymized by replacing
 3 vehicle plate number of each individual with a unique identifier and full history trajectory of each
 4 individual is published. Under this circumstance, each record in an anonymized LPR database rep-
 5 resents the complete history trajectory of an individual and complete spatiotemporal information
 6 (a sequence of spatiotemporal tuples) of individuals is released. Given the parking database, we
 7 aim to link one LPR record to each parking data record to re-identify the true identity (vehicle plate
 8 number) of that anonymized LPR record.

9 Given a vehicle l with its parking record P_k^l in the parking database P_k of parking k ($k =$
 10 $1, \dots, m$), it can be naturally imagined that this vehicle passed through a detector, and then it en-
 11 tered the parking lot k . After that, it left the parking and captured by another detector. In other
 12 words, the parking record can be inserted into a pair of temporal consecutive spatiotemporal tuples
 13 in one's trajectory, i.e., a data record in the anonymized LPR database. Consider an individual j
 14 with an anonymized LPR record R^j , if he used to park in the parking k and since generated the
 15 parking data record P_k^l , without loss of generality, we suppose that parking occurred between i -th
 16 tuple and $(i+1)$ -th tuple in R^j (denoted as $R_i^j = (T_i^j, LocR_i^j)$ and $R_{i+1}^j = (T_{i+1}^j, LocR_{i+1}^j)$, respec-
 17 tively), then we can estimate the timestamp of R_i^j (denoted as T_i^j) as the following Equation 1.

$$\hat{T}_i^j \approx T_{P_k^l}^{in} - t_{i,k}^j \quad (1)$$

18 Here $t_{i,k}^j$ represents the estimated travel time for the linkage between parking record P_k^l and
 19 the spatiotemporal tuple R_i^j , which can be calculated by the distance from $LocR_i^j$ to the location
 20 of parking k (indicated as $LocP_k$) divided by the estimated average speed in this route. To im-
 21 prove the precision of travel time estimation, we request traffic information using navigation API,
 22 by which we can get the distance of planned routes and the navigation duration of routes. Since
 23 traffic flow patterns are usually similar in neighboring links and over the same period, we approx-
 24 imate the average speed of vehicle j in the target route using the average speed in the previous
 25 link it traveled along, which can be calculated by the earlier pair of spatiotemporal tuples in R^j
 26 as $Dist_{i-1,i}^j / (T_i^j - T_{i-1}^j)$, where $Dist_{i-1,i}^j$ is the requested route distance from $LocR_{i-1}^j$ to $LocR_i^j$.
 27 This thus allows us to capture the real-time traffic dynamics (e.g., congestion status or free-flow
 28 status) and intrinsic driving characteristics of vehicles (e.g., heavy truck or private car), which
 29 serves as a key factor for the travel time estimation in the linkage process. Then $t_{i,k}^j$ is calculated

1 by $Dist_{i,k}^j/V_{i,k}^j$, where $Dist_{i,k}^j$ is the requested route distance from $LocR_i^j$ to $LocP_k$. In addition,
 2 taking into account the existence of abnormal estimated travel time arising from traffic anomaly
 3 (e.g., traffic accidents, road constructions and aggressive driving, etc.), we calibrate the estimated
 4 travel time using the requested navigation duration from $LocR_i^j$ to $LocP_k$ (indicated as $Dura_{i,k}^j$) as
 5 equation 2.

$$t_{i,k}^j = \begin{cases} \frac{Dist_{i,k}^j}{V_{i,k}^j}, & \text{if } t_{i,k}^j \in [(1-\alpha)Dura_{i,k}^j, (1+\alpha)Dura_{i,k}^j]; \\ Dura_{i,k}^j, & \text{if } t_{i,k}^j \notin [(1-\alpha)Dura_{i,k}^j, (1+\alpha)Dura_{i,k}^j]. \end{cases} \quad (2)$$

6 Similarly, the travel speed of vehicle j on the route between $LocP_k$ and $LocR_{i+1}^j$ can be
 7 approximated by the calculated average speed on the later link between $LocR_{i+1}^j$ and $LocR_{i+2}^j$, and
 8 the estimated travel time on target route can be calibrated by the requested duration $Dura_{k,i+1}^j$ on
 9 it. Then the travel time of vehicle j from parking k to the following detector can be estimated by
 10 equation 3. If vehicle l and vehicle j are the same individual, then the timestamp that this vehicle
 11 went through the detector (i.e., T_{i+1}^j) can be approximated by $T_{P_k}^{out} + t_{k,i+1}^j$. See equation 4.

$$t_{k,i+1}^j = \begin{cases} \frac{Dist_{k,i+1}^j}{V_{k,i+1}^j}, & \text{if } t_{k,i+1}^j \in [(1-\alpha)Dura_{k,i+1}^j, (1+\alpha)Dura_{k,i+1}^j]; \\ Dura_{k,i+1}^j, & \text{if } t_{k,i+1}^j \notin [(1-\alpha)Dura_{k,i+1}^j, (1+\alpha)Dura_{k,i+1}^j]. \end{cases} \quad (3)$$

$$\hat{T}_{i+1}^j \approx T_{P_k}^{out} + t_{k,i+1}^j \quad (4)$$

12 Adding a fluctuation of the estimated travel time, if the vehicle j leaved the previous detec-
 13 tor within the period $[\hat{T}_i^j - T_f^{i,j}, \hat{T}_i^j + T_f^{i,j}]$ and arrived at the following detector within the period
 14 $[\hat{T}_{i+1}^j - T_f^{i+1,j}, \hat{T}_{i+1}^j + T_f^{i+1,j}]$, then the parking record P_k^l may be linked to the pair of spatiotemporal
 15 tuples R_i^j and R_{i+1}^j in the anonymized LPR record R^j and accordingly they may be the same vehicle.
 16 Differing from a constant time threshold was set for the fluctuation in (12), we allow the estimated
 17 travel time fluctuate within a proportion range, i.e., $T_f^{i,j} = \beta \hat{T}_i^j$, which can screen more reasonable
 18 estimated travel time. For conducting linkage attack, if two continuous tuples $R_i^j = (T_i^j, LocR_i^j)$
 19 and $R_{i+1}^j = (T_{i+1}^j, LocR_{i+1}^j)$ in one's LPR record R^j satisfy $(1-\beta)\hat{T}_i^j \leq T_i^j \leq (1+\beta)\hat{T}_i^j$ and
 20 $(1-\beta)\hat{T}_{i+1}^j \leq T_{i+1}^j \leq (1+\beta)\hat{T}_{i+1}^j$ simultaneously, (i, j) is added into a candidate set for parking
 21 record P_k^l (denoted as C_k^l). If C_k^l is not empty, then a optimal match between R_i^j , R_{i+1}^j and P_k^l can
 22 be selected by equation 5. Equation 5 means that the linkage result of our adversary model is the
 23 one with the minimal total time difference between the actual passing timestamps and estimated
 24 passing timestamps at the previous detector and the following detector. If the match is correct, i.e.,
 25 R_i^j , R_{i+1}^j and P_k^l indeed belong to the same vehicle, the parking record and anonymized LPR record
 26 thus are successfully linked, and the true identity of the anonymized record in LPR database, i.e.,
 27 the vehicle plate number, is successfully inferred, which poses a threat to the individual's privacy.
 28 Applying the model to conduct linkage attack for each parking record in parking database, we can
 29 finally get the match set of the parking database.

$$\hat{i}, \hat{j} = \arg \min_{(i,j) \in C_k^l} |T_i^j - \hat{T}_i^j| + |T_{i+1}^j - \hat{T}_{i+1}^j| \quad (5)$$

1 Linking under segmented trajectory publishing (Mode 2)

In mode 2, traffic data owner separates the complete trajectory to several segments at the nodes where the vehicle is in parking state, and each segment will have a unique identifier. To re-identify the vehicle by linkage attack, for a parking record (e.g. P_k^l), two trajectory segments, e.g. R^j and $R^{j'}$, will be taken into account if between which that parking record can be inserted into, mathematically, this condition can be formulated as Equation 6.

$$T_{end}^j \leq T_{P_k^l}^{in} \leq T_{P_k^l}^{out} \leq T_{st}^{j'} \quad (6)$$

where T_{end}^j and $T_{st}^{j'}$ indicates the passing timestamp of the last tuple of R^j and the first tuple of $R^{j'}$ separately, and $j \neq j'$. For better expound in the following content, variables with the subscription 'end' are for the last tuple of a segment, while variables with the subscription 'st' are for the first tuple of a segment. The adversary model proposed in section 4.1 can be applied to link the parking data to the LPR data released through mode 2 as well. But instead of applying the adversary model to each node, combined nodes, which are defined as the nodes formed by the last tuple of a segment and the first tuple of another segment, e.g. $(R_{end}^j, R_{st}^{j'})$, will be the only attacking objectives, as the parking behavior can only happen in the combined nodes. However, when we have a large LPR dataset, the number of combination of LPR records that can satisfy Equation 6 will be enormous, requiring a stricter restriction. To this end, a threshold for filtering the records swarm into the algorithm is introduced, see Equation 7. Equation 7 reduces the number of possible LPR records for constructing combined nodes by restricting the timestamps of the pre-parking tuple and the pos-parking tuple in a range of the timestamp estimated by leveraging the requested navigation duration. γ_t is a time-unit filtering parameter.

$$\begin{aligned} T_{end}^j &\in [T_{P_k^l}^{in} - (Dura_{end,k}^j + \gamma_t), T_{P_k^l}^{in} - (Dura_{end,k}^j - \gamma_t)] \\ T_{st}^{j'} &\in [T_{P_k^l}^{out} + (Dura_{k,st}^{j'} - \gamma_t), T_{P_k^l}^{out} + (Dura_{k,st}^{j'} + \gamma_t)] \end{aligned} \quad (7)$$

For vehicles which do not have parking records, their complete trajectory will be released directly like mode 1. So in mode 2, the timestamp of the pre-parking tuple and the pos-parking tuple can be estimated by Equation 8.

$$\hat{T}_{end}^j \approx T_{P_k^l}^{in} - t_{end,k}^j \quad (8a)$$

$$\hat{T}_{st}^{j'} \approx T_{P_k^l}^{out} + t_{k,st}^{j'} \quad (8b)$$

And similar to the description in section 4.1, $t_{end,k}^j$ and $t_{k,st}^j$ can be estimated as Equation 9.

$$t_{end,k}^j = \begin{cases} \frac{Dist_{end,k}^j}{V_{end,k}^j}, & \text{if } t_{end,k}^j \in [(1-\alpha)Dura_{end,k}^j, (1+\alpha)Dura_{end,k}^j]; \\ Dura_{end,k}^j, & \text{if } t_{end,k}^j \notin [(1-\alpha)Dura_{end,k}^j, (1+\alpha)Dura_{end,k}^j]. \end{cases} \quad (9a)$$

$$t_{k,st}^{j'} = \begin{cases} \frac{Dist_{k,st}^{j'}}{V_{k,st}^{j'}}, & \text{if } t_{k,st}^{j'} \in [(1-\alpha)Dura_{k,st}^{j'}, (1+\alpha)Dura_{k,st}^{j'}]; \\ Dura_{k,st}^{j'}, & \text{if } t_{k,st}^{j'} \notin [(1-\alpha)Dura_{k,st}^{j'}, (1+\alpha)Dura_{k,st}^{j'}]. \end{cases} \quad (9b)$$

And $V_{end,k}^j$ and $V_{k,st}^{j'}$ can be calculated as Equation 10.

$$V_{end,k}^j = \frac{Dist_{end-1,end}^j}{T_{end}^j - T_{end-1}^j} \quad (10a)$$

$$V_{k,st}^{j'} = \frac{Dist_{st,st+1}^{j'}}{T_{st+1}^{j'} - T_{st}^{j'}} \quad (10b)$$

Similarly, we set a fluctuation range for the estimated travel time, for the LPR record R^j which can meet $(1-\beta)\hat{T}_{end}^j \leq T_{end}^j \leq (1+\beta)\hat{T}_{end}^j$ and $R^{j'}$ which can meet $(1-\beta)\hat{T}_{st}^{j'} \leq T_{st}^{j'} \leq (1+\beta)\hat{T}_{st}^{j'}$, (j, j') will become a candidate of the parking record R_k^l and be added into C_k^l . The optimum will be selected by Equation 11. It means the candidate whose passing timestamps result into the smallest difference from the estimated timestamps will be the output of the adversary model. If j, j' and l are equal, it will be counted as a success.

$$\hat{j}, \hat{j}' = \arg \min_{(j,j') \in C_k^l} |T_{end}^j - \hat{T}_{end}^j| + |T_{st}^{j'} - \hat{T}_{st}^{j'}| \quad (11)$$

1 EXPERIMENT AND RESULT ANALYSIS

2 In this section, we quantitatively measure the privacy vulnerability due to linkage attack among
3 LPR data and parking data by applying the proposed adversary model to an actual scenario. Two
4 different modes of LPR data publishing are examined and the corresponding results are analyzed.

5 The linkage attack was conducted in a scenario with five parking lots. Due to the wide
6 coverage of camera detector in the LPR system, the accessible detectors of parking lots are not
7 far away from them. Consequently, to improve the efficiency of the adversary model, we only
8 consider detectors within a 5-km radius of one parking as its latent previous detectors and following
9 detectors. In this case, 24 detectors are selected and the network geometry of the detectors and
10 parking lots is shown in Figure 3. The parking data set and LPR data set were collected on the
11 same day (July 18th, 2018). We construct the parking database by selecting the parking records
12 of the five parking lots in our experiment area. After data cleaning (e.g., drop duplicated records),
13 there are 145 parking records in the parking database. As for the LPR database, we select all
14 history trajectory of vehicles that passed through the one or more of the latent 24 detectors. Then
15 for the data publishing of mode 1, we replace each vehicle plate number with a unique random
16 identifier to construct the anonymized LPR database. When it comes to mode 2, the trajectory of

1 vehicles that parked in one of the five parking lots is first segmented according to the actual parking
2 record in parking database. Then a unique random identifier is allocated to each trajectory segment
3 for data anonymization in mode 2.

4 Applying the adversary model proposed in section 4, we conduct linkage attack for each
5 parking record in parking database under two different patterns of LPR data publishing. We set the
6 calibration rate α to 0.35 and the fluctuation rate β to 0.4, respectively. For mode 1, we get a match
7 set of 141 matched records, with 101 records successfully linked (i.e., the target parking record
8 and the matched anonymized LPR record belong to the same vehicle) and 40 records incorrectly
9 matched. As for mode 2, we get a match set with all 145 parking records matched. There are
10 53 records successfully linked to anonymized LPR records. To quantify the privacy risk more
11 comprehensively, we introduce two classical criteria in the field of information retrieval, namely
12 precision and recall. We define precision in our linkage attack scenario as the number of parking
13 records successfully matched in our match set divided by the total number of match results in
14 our match set (i.e., the size of match set). Recall is defined as the number of parking records
15 successfully matched divided by the total number of parking records in parking database. Then the
16 precision and recall under mode 1 are $101/141 = 71.63\%$ and $101/145 = 69.66\%$, respectively.
17 The precision and recall under mode 2 are $53/145 = 36.55\%$ and $53/145 = 36.55\%$, respectively.
18 We can see from the results that individual with anonymized LPR data published in mode 1 have a
19 high probability of being re-identified by our adversary model (over 70% parking records in match
20 set can be linked successfully to anonymized LPR database). In addition, the privacy disclosure
21 risk drops significantly by almost a half (from 71.3% to 36.55%) from mode 1 to mode 2, which
22 provide a practical solution for privacy preservation in LPR data publishing. This is because in
23 mode 1, our model conducts linkage attack between two consecutive spatiotemporal points in
24 one's trajectory. The linkage attack is successful when the matched anonymized LPR record and
25 target parking record belong to the same vehicle. While in mode 2, we need to link the target
26 parking record to two trajectory segments, i.e., two anonymized records in the published database.
27 And only when these three records belong to the same individual, the linkage attack is judged
28 successful.

29 Besides selecting a single candidate (one LPR record under mode 1 and two under mode 2)
30 from the candidate set for each target parking record, an adversary may choose to output multiple
31 candidates to improve the probability of linking successfully to the target parking record. Instead of
32 selecting the candidate with minimal time error, we output the top 3 candidates with minimal error
33 in the candidate set to our match set and compare it with our original results (denoted as top 1).
34 With the selection of top 3 candidates, in mode 1, the precision and recall are $130/257 = 50.58\%$
35 and $130/145 = 89.66\%$, respectively. In mode 2, the precision and recall are $98/362 = 27.07\%$
36 and $98/145 = 67.59\%$, respectively. The overall results are compared in Figure 2. It is clear from
37 the bar graph that, when top 3 candidates are selected, more parking records in parking database
38 are linked successfully. But the precision decline since more candidates are added into match set.
39 This means that an adversary can get more individual re-identified at the cost of precision.

40 CONCLUSION AND DISCUSSION

41 In this paper, we quantitatively measure the privacy disclosure risk under linkage attack across
42 multi-source individual-based mobility data sets. Performing a case study on LPR data and park-
43 ing data, we propose a traffic-knowledge-driven adversary model for linkage attack conducting
44 among LPR data and parking data. Two common modes of LPR data publishing/anonymization

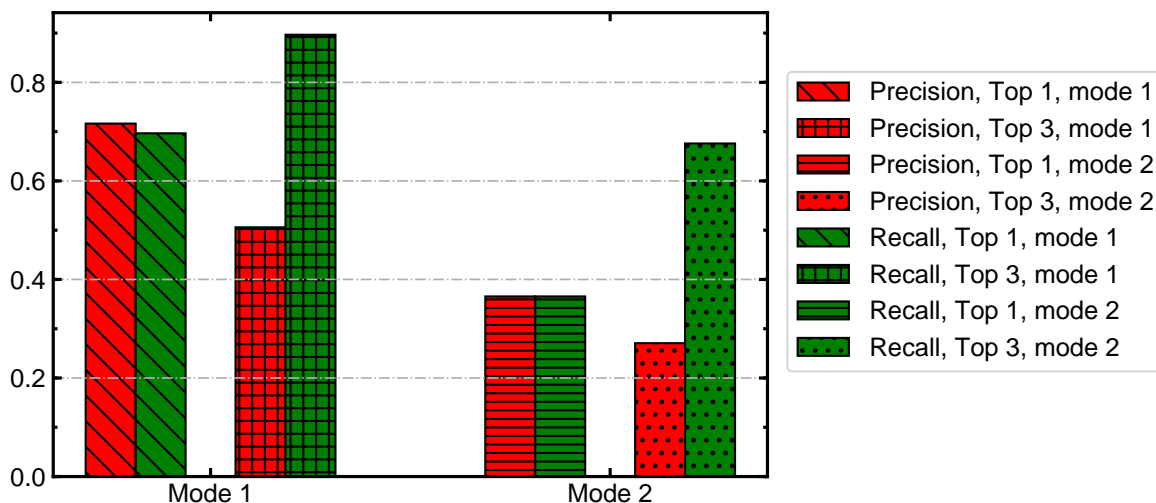


FIGURE 2 : Detectors and parkings considered in the scenario

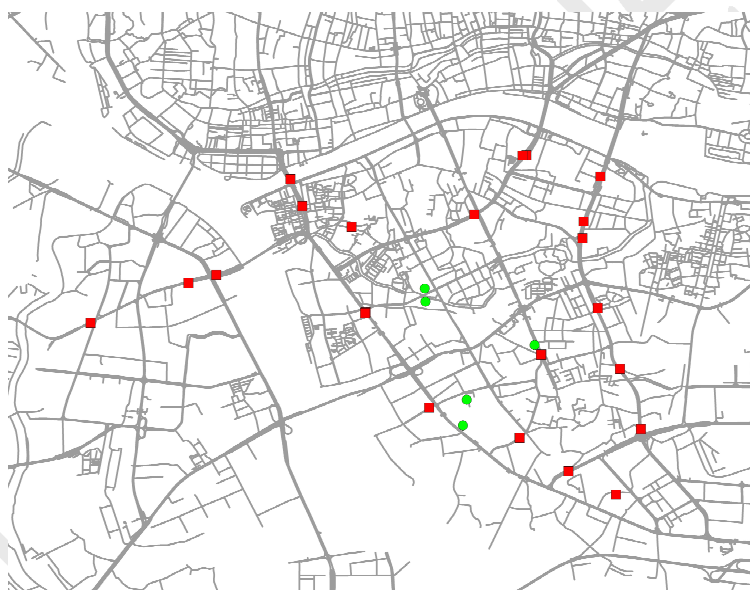


FIGURE 3 : Detectors and parkings considered in the scenario

1 are examined by applying the proposed model. Two quantitative evaluation criteria, precision and
 2 recall are introduced to investigate the privacy disclosure risk under linkage attack. Empirical ex-
 3 periments are conducted in actual scenario in Guangzhou, China. The result shows that when LPR
 4 data is published in mode 1, our adversary model can achieve a precision of 71.61% and a recall
 5 of 69.66% which illustrates that individuals in anonymized LPR database still face a high proba-
 6 bility of being re-identified. However, in mode 2, the precision and recall decline to 36.55%. This
 7 means that LPR data published in mode 2 has a high level of privacy-preserving. When output top
 8 3 candidates, the recall of the adversary model in mode 1 increases to 89.66%, and that of mode
 9 2 increases to 67.59%, respectively. It means, when top 3 candidates are selected, there will have
 10 more parking records been linked successfully to anonymized LPR records.

11 This current work focuses on revealing the great privacy vulnerability caused by linkage

1 attack across multi-source individual-based mobility data sets from a quantitative perspective. Fu-
2 ture extensions will explore the data utility in specific transportation application under two differ-
3 ent modes of LPR data publishing and discuss the privacy-and-utility trace-off of different data
4 publishing modes. In addition, we will investigate possible solutions to preserve privacy against
5 linkage attack across multi-source individual mobility data (e.g., reducing the probability of being
6 linked successfully). We hope this work could stimulate more discussion to address the privacy
7 issues in the context of multi-source individual-based mobility data sets.

8 ACKNOWLEDGEMENT

9 This research is mainly supported by the National Natural Science Foundation of China (No.
10 11574407) and the Science and Technology Planning Project of Guangzhou City, China (No.
11 201704020142).

12 AUTHOR CONTRIBUTIONS

13 J.G., Q.L. and M.C. designed the research; J.G. and Q.L. performed the research; J.G. and Q.L.
14 analyzed the data; J.G., Q.L. and M.C. wrote the paper.

15 REFERENCES

- 16 [1] Sweeney, L., k -anonymity: a model for protecting privacy. *International Journal of Uncer-*
17 *tainty, Fuzziness and Knowledge-Based Systems*, Vol. 10, No. 05, 2002, pp. 557–570.
- 18 [2] Gao, J., L. Sun, and M. Cai, Measuring privacy vulnerability of individual mobility traces: a
19 case study on license plate recognition data. In *Transportation Research Board (TRB) Annual*
20 *Meeting*, 2019.
- 21 [3] Gao, J., L. Sun, and M. Cai, Quantifying privacy vulnerability of individual mobility traces:
22 A case study of license plate recognition data. *Transportation Research Part C Emerging*
23 *Technologies*, Vol. 104, 2019, pp. 78–94.
- 24 [4] Hoh, B., M. Gruteser, H. Xiong, and A. Alrabady, Enhancing security and privacy in traffic-
25 monitoring systems. *IEEE Pervasive Computing*, Vol. 5, No. 4, 2006, pp. 38–46.
- 26 [5] Liu, X., H. Zhao, M. Pan, H. Yue, X. Li, and Y. Fang, Traffic-aware multiple mix zone place-
27 ment for protecting location privacy. In *2012 Proceedings IEEE INFOCOM*, IEEE, 2012, pp.
28 972–980.
- 29 [6] Li, M., L. Zhu, and X. Lin, Privacy-Preserving Traffic Monitoring with False Report Filtering
30 via Fog-assisted Vehicular Crowdsensing. *IEEE Transactions on Services Computing*, 2019.
- 31 [7] de Montjoye, Y. A., L. Radaelli, V. K. Singh, and A. S. Pentland, Identity and privacy. Unique
32 in the shopping mall: on the reidentifiability of credit card metadata. *Science*, Vol. 347, No.
33 6221, 2015, pp. 536–9.
- 34 [8] Ma, C. Y., D. K. Yau, N. K. Yip, and N. S. Rao, Privacy vulnerability of published anonymous
35 mobility traces. *IEEE/ACM transactions on networking (TON)*, Vol. 21, No. 3, 2013, pp. 720–
36 733.

- 1 [9] Gao, J., L. Sun, and M. Cai, Quantifying privacy vulnerability of individual mobility traces:
2 A case study of license plate recognition data. *Transportation Research Part C: Emerging*
3 *Technologies*, Vol. 104, 2019, pp. 78–94.
- 4 [10] Aïvodji, U. M., S. Gambs, M.-J. Huguet, and M.-O. Killijian, Meeting points in ridesharing:
5 A privacy-preserving approach. *Transportation Research Part C: Emerging Technologies*,
6 Vol. 72, 2016, pp. 239–253.
- 7 [11] Hoh, B., M. Gruteser, R. Herring, J. Ban, D. Work, J.-C. Herrera, A. M. Bayen, M. An-
8 navaram, and Q. Jacobson, Virtual trip lines for distributed privacy-preserving traffic moni-
9 toring. In *Proceedings of the 6th international conference on Mobile systems, applications,*
10 *and services*, ACM, 2008, pp. 15–28.
- 11 [12] Sun, Z., B. Zan, X. J. Ban, and M. Gruteser, Privacy protection method for fine-grained ur-
12 ban traffic modeling using mobile sensors. *Transportation Research Part B: Methodological*,
13 Vol. 56, 2013, pp. 50–69.