



HAL
open science

Can We Map-Match Individual Cellular Network Signaling Trajectories in Urban Environments? Data-Driven Study

Loïc Bonnetain, Angelo Furno, Jean Krug, Nour-Eddin El Faouzi

► To cite this version:

Loïc Bonnetain, Angelo Furno, Jean Krug, Nour-Eddin El Faouzi. Can We Map-Match Individual Cellular Network Signaling Trajectories in Urban Environments? Data-Driven Study. *Transportation Research Record*, 2019, 2673 (7), pp74-88. <10.1177/0361198119847472>. <hal-03181077>

HAL Id: hal-03181077

<https://hal.science/hal-03181077v1>

Submitted on 11 Jun 2021

HAL is a multi-disciplinary open access archive for the deposit and dissemination of scientific research documents, whether they are published or not. The documents may come from teaching and research institutions in France or abroad, or from public or private research centers.

L'archive ouverte pluridisciplinaire HAL, est destinée au dépôt et à la diffusion de documents scientifiques de niveau recherche, publiés ou non, émanant des établissements d'enseignement et de recherche français ou étrangers, des laboratoires publics ou privés.



HAL Authorization

1 Can we map-match individual cellular network signaling
2 trajectories in urban environments? A data-driven
3 study.

4
5 Camera-ready version uploaded on: Mars, 15th 2018

6 Loïc Bonnetain,¹,

7 ¹ Univ. Lyon, IFSTTAR, ENTPE, LICIT UMR_T9401, F-69675, Lyon, France
8 loic.bonnetain@entpe.fr

9 Angelo Furno¹,

10 ¹ Univ. Lyon, IFSTTAR, ENTPE, LICIT UMR_T9401, F-69675, Lyon, France
11 angelo.furno@ifsttar.fr

12 Jean Krug¹,

13 ¹ Univ. Lyon, IFSTTAR, ENTPE, LICIT UMR_T9401, F-69675, Lyon, France
14 jean.krug@entpe.fr

15 Nour-Eddin El Faouzi^{1,2}

16 ¹ Univ. Lyon, IFSTTAR, ENTPE, LICIT UMR_T9401, F-69675, Lyon, France

17 ² Queensland University of Technology, STRC, Gardens Point Campus, 2 George Street,
18 G.P.O. Box 2434, Brisbane, Queensland 4001, Australia.

19 nour-eddin.elfaouzi@ifsttar.fr

20 *Accepted for publication in the Transportation Research Record (TRR),*
21 *Journal of the Transportation Research Board*

22
23 **Word Count:**

24 Number of words: 6620

25 Number of tables: 2 (250 words each)

26 Total: 7120

Abstract

Mobile phone data collected by network operators can provide fundamental insights on individual and aggregate mobility of people, at unprecedented spatio-temporal scales. However, traditional Call Detail Records (CDR) have fundamental issues due to low accuracy along both the spatial and the temporal dimensions, which limits their applicability for detailed studies on mobility, especially in urban scenarios. In this paper, we focus on a new generation of mobile phone passive data, individual cellular-network signaling data, characterized by higher spatio-temporal resolutions than traditional CDR data. We design a framework based on unsupervised Hidden Markov Model (HMM) for map-matching such kind of data on multi-modal transportation network, aimed at accurately inferring the complex multi-modal travel itineraries and popular paths people follow in their urban daily mobility. This information, especially if computed at large spatio-temporal scales, can represent a solid basis for studying actual and dynamic travel demand, to properly dimension multi-modal transport systems and even perform anomaly detection and adaptive network control. We evaluate our approach in a case-study based on real cellular traces collected by a major French operator in the city of Lyon, and propose a validation study at both microscopic and macroscopic levels. The results show that our approach can properly handle sparse and noisy cell phone trajectories in urban complex environments. Besides, the results are promising concerning popular paths detection and reconstruction of Origin-Destination matrices.

Keywords: Map-matching, Mobile phone, Hidden Markov Model, Multi-modal transportation network

1. INTRODUCTION

In recent years, the widespread diffusion of mobile devices and the exploding consumption of Internet traffic via 3G and 4G technologies have made mobile phone data a crucial source of information in multiple domains. This is especially true in the field of transportation, as these data, usually including spatio-temporal information related to mobile phone users, can provide fundamental insights on people’s mobility both at individual and aggregate scales.

For instance, Call Detail Records (CDR), also referred to as *mobile phone passive data*, have fed plenty of large-scale studies on human mobility, given the possibility to study urban mobility at unprecedented spatio-temporal scales [1]. Relevant work based on CDR data comprises *i.e.*, modelling the general laws governing human movements [2], reconstructing Origin-Destination (OD) matrices [3], understanding urban land use [4, 5] and inferring population density [6]. Mobile phone passive data are increasingly used also in operational contexts by mobility service providers and traffic authorities, in conjunction with – or even in the place of – more traditional data sources on mobility like census data, local travel surveys and logs from road-side units (e.g., loop detectors, LI-DAR or acoustic sensors, Bluetooth scanners, etc.). In fact, the latter suffer from very high deployment costs, extremely poor spatio-temporal resolutions, and are rarely informative in terms of individual mobility [7, 8].

However, despite significant benefits, CDR still have fundamental issues that need to be addressed due to low accuracy along both the spatial dimension (*i.e.*, user location is only known at the cell sector or base station coverage levels) and the temporal one (*i.e.*, events are recorded only when the user performs a voice call or texts a message), which limits their applicability for detailed studies on mobility, especially in urban settings.

In such scenarios, Global Positioning System (GPS) logs still represent the preferred choice, since they allow for obtaining data with a higher degree of accuracy (*i.e.*, meters) and temporal frequency (*i.e.*, seconds). Such measures can be relatively easily analyzed and mapped to mobility patterns by relying on machine learning techniques [9, 10]. However, a huge overhead exists in collecting detailed GPS datasets at statistically relevant scales, such data being mostly retrieved on voluntary basis or via special agreements involving only a small sample of users or vehicles [7]. Given these limitations, extended variants of CDR (namely, network signaling logs and Internet session reports) are currently collected by network providers and investigated by the research community. Differently from CDR data, network signaling data report on multiple kinds of events besides calls and text messages (e.g., IP protocol message exchanges, hand-overs, location updates, etc.) thus increasing the spatio-temporal sampling frequency of mobile phone passive data. Research on this kind of data is however still at early stages. In this paper, by building on related work from the field of GPS map-matching and CDR analysis, we focus on the possibility of inferring relatively accurate measures of both individual and aggregate mobility flows from cellular network signaling data.

In fact, in the context of next-generation intelligent transportation systems, inferring individual trips with a certain degree of accuracy, even in urban environment, will enable a better and more precise understanding of both microscopic and macroscopic mobility. Such knowledge is expected to be leveraged in many applications such as multi-modal transportation network analysis and optimization, traffic routing and adaptive control.

Map-matching of GPS traces has been widely studied in the literature [11] and state-of-the-art approaches can achieve high accuracy in the presence of large-sampling rate data (e.g., sampling rate of 1 Hz) [12]. Although, it is worth noticing that, in terms of penetration rate and energy consumption, mobile phone data represent much better candidates than GPS data to track users in a large-scale and suitable way [13]. In this paper, we deal with the sparsity (in time and space), the noise and large localization error associated to cell phone trajectories, that make the task of

1 reconstructing trips challenging [14].

2 A methodology based on Hidden Markov Model (HMM) is presented as the core of a map-
3 matching algorithm engineered for cellular network signaling data. The algorithm infers the most
4 likely path of a mobile phone user, given a sequence of network signaling events emitted by her/his
5 smartphone during a trip.

6 The network modeling of the transportation graph and the cellular network, key elements of
7 the proposed approach, are also presented. A study case using the HMM-based map-matching is
8 performed with two different datasets from the city of Lyon (France).

9

10 The key contributions of this work are:

- 11 • The main solution for the challenging problem of mapping cellular trajectories to the multi-
12 modal transportation network instead of only considering the road network.
- 13 • An unsupervised HMM-based map-matching approach allowing to infer trajectories on phys-
14 ical network from any sparse (spatially and temporally) cellular trajectory in dense urban
15 context. This is made possible by a more fine-grained modeling of both transportation and
16 cellular networks, compared to state-of-the art approaches [15].
- 17 • Dataset collection of real-world cellular trajectories related to a group of users in the Lyon
18 metropolitan area, France. The dataset has been collected by Orange, the major French
19 mobile network operator. Despite the sparsity of the available data, we analyze our approach
20 in two case studies, for both macroscopic and microscopic mobility analysis.

21 The rest of the paper is organized as follows. In Sec. 2, we present related work. In Sec. 3,
22 we formulate key definitions to define the map-matching problem. In Sec. 4, network modeling is
23 presented. In Sec. 5, we discuss the methodology of our HMM-based model. In Sec. 6, we evaluate
24 our approach on the considered dataset. We conclude in Sec. 7 by discussing the limits of our
25 approach and future directions.

26 2. RELATED WORK

27 Map-matching is a basic operation for improving positioning accuracy by integrating positioning
28 data with spatial transportation data to identify the correct link on which a mobile object is
29 traveling [16]. Several approaches exist in the literature to solve the problem of map-matching
30 GPS traces to a transportation network. Quddus *et al.* [17] categorize map-matching approaches
31 in four classes.

32 *Geometric approaches* only use the spatial geometry of the network: the most simple and
33 popular map-matching algorithm consists in matching each position point to the closest node in
34 the network [18].

35 *Topological approaches* use geometric information as well as topological information like the
36 existence of connectivity between nodes of the network [19]. Very sensitive to noise and outliers,
37 these approaches are not appropriate to solve map matching problems in presence of highly noisy
38 and sparse data.

39 The third kind of approaches exploits *probabilistic methods*: a confidence region around the
40 location of the moving object is defined. Then, candidate network links are identified as those
41 present in this confidence region. The evaluation of the candidates is based on the geometrical
42 criteria.

43 Finally, *advanced map-matching approaches* use more complex mathematical tools. A non
44 exhaustive list of these methods includes, *i.e.*, the Kalman Filter, its Extended Kalman version [20],

1 Dempster–Shafer theory [19], fuzzy logic models [21], or the application of Bayesian inference [22].
 2 These state-of-the-art algorithms may achieve a quasi-perfect accuracy (location error lower than
 3 10 meters) with high sampling rate GPS data. Newson *et al.* [12] first introduce HMM-based
 4 map-matching dealing with different GPS traces sampling rate. Their approach turned out to be
 5 much more robust and accurate with sparse and noisy trajectories compared to standard advanced
 6 map-matching approaches for high sampling rate data.

7 As a consequence of the growing availability of large-scale mobile phone data collected by
 8 network operators, map-matching cell phone trajectories is recently becoming a challenging task for
 9 researchers. Most of the approaches used with cellular trajectories are based on those traditionally
 10 designed for GPS map-matching. Schulze *et al.* [23] use a probabilistic approach: their solution
 11 restricts the set of admissible routes to a corridor by estimating the area within which a user is
 12 allowed to travel and infers path using the shortest path on candidate routes. With only 55% of
 13 correct matches, this method has been outperformed by a HMM-based approach recently developed
 14 by Jagadeesh *et al.* [24], which reaches 75% of median accuracy.

15 Finally, HMM-based map-matching has become the state-of-the-art approach for noisy and
 16 sparse location data and, *a fortiori*, mobile phone trajectory. Thiagarajan *et al.* [13] and, more
 17 recently, Algizawy *et al.* [25] developed supervised HMM models exhibiting good accuracy (75%
 18 for Thiagarajan *et al.* approach). However, such an approach needs to train the HMM model
 19 with a large amount of labeled cellular trajectories, which are very hard to obtain, especially when
 20 dealing with highly dynamic and irregular environments, such as urban areas. Instead, we prefer to
 21 focus on unsupervised models that do not require collecting and labeling any trajectory. Moreover,
 22 we state that additional information such as signal strength of observation are relatively hard to
 23 obtain from mobile network operators and therefore should not be required by the map-matching
 24 approach, as for example is the case in [13]. Jagadeesh *et al.* [24] proposed an online map-matching
 25 algorithm combining HMM-based map-matching and route choice model.

26 Finally, it is worth noticing that most of the approaches match cellular trajectories only to road
 27 networks, without considering other transportation modes. Among the very few exceptions, it is
 28 necessary to mention the methodology recently proposed by Asgari *et al.* [15]. The authors have
 29 developed an approach, namely CT-Mapper, which has been designed with very similar objectives
 30 to those of our work. CT-Mapper is an unsupervised HMM model which aims at mapping sparse
 31 multi-modal cellular trajectories by using a multilayer transportation network. Yet, CT-Mapper
 32 has some limitations: the multilayer network allows for unrealistic paths (each subway station
 33 is connected to its closest road intersection for simplification matters). In addition, CT-Mapper
 34 requires already cleaned cellular trajectories. Dealing with noisy mobile phone data requires an
 35 advanced cleaning process which is not further specified in CT-Mapper. Finally, Asgari *et al.*
 36 filtered out trajectories whose lengths are shorter than 5 kilometers, only keeping longer trajectories,
 37 with an the average length of 26.5 kilometers. Hence, CT-Mapper has been validated only in
 38 inter-urban mobility scenarios, thus seeming not to handle urban mobility. Our model aims at
 39 investigating and overcoming these limitations, using a more sophisticated approach especially
 40 concerning network modeling.

41 3. PROBLEM STATEMENT

42 The section presents the main definitions, and a formal conceptualization of the problem of map-
 43 matching sparse cell-phone trajectories to the underlying multi-modal transportation network. The
 44 definitions reported in the following are based on those used in strictly related recent work [15, 23]:

45 **Definition 1 (Signaling event)** *A signaling event is defined as any observation resulting from a*
 46 *communication activity between a cell phone and a base station antenna of the telecommunication*

1 network. Each observation o is defined as a tuple (c_i, t) where c_i is the identifier of the antenna
 2 (see definition 5) where the mobile phone event took place and $t \in \mathbb{N}$ is the timestamp of the event.

3 **Definition 2 (Cell phone trajectory)** A cell phone trajectory $T = (o_1, \dots, o_n)$ is defined as a
 4 sequence of network signaling events, ordered by their timestamps and related to the same mobile
 5 phone user.

6 We consider the following as typical kinds of signaling events: i) communication events (i.e., calls
 7 and SMS); ii) handover events (i.e., cell changes during an established communication) and Location
 8 Area (LA) updates; iii) network attachment/detachment events; iv) data/internet connections.

9 **Definition 3 (Multi-Layer Transportation Graph)** A Multi-Layer Transportation Graph is
 10 defined as a directed graph $G = (V, E, L, \Psi)$ where V, E represent the vertices and the edges, re-
 11 spectively, and L is the set of possible layers related to different transportation modes. In our
 12 study, we focus on four layers only: road, bus, tramway and subway. Function Ψ indicates the
 13 layer associated to a given node, i.e., $\Psi : V \rightarrow L$ in G . Transportation Layer $G^l = (V^l, E^l)$ is
 14 a subset of G where $V^l = \{v | v \in V, \Psi(v) = l\}$ and $E^l = \{\langle v_i | v_j \rangle \in E, \Psi(v_i) = \Psi(v_j) = l\}$. Each
 15 node v_i is characterized by its latitude and longitude (i.e., the geographical position $v_i = \langle \text{lat}, \text{lon} \rangle_i$).
 16 CrossLayer edge set $E^{\text{cl}} \subset E$ defines the edges with pair of nodes not belonging to the same layer:
 17 $E^{\text{cl}} = \{\langle v_i | v_j \rangle \in E | \Psi(v_i) \neq \Psi(v_j)\}$.

18 **Definition 4 (Cellular Network)** The cellular network is defined as a set of base station an-
 19 tennas $C = (c_0, c_1 \dots c_p)$, where each antenna $c_p = (\phi, \lambda, z)$ is characterized by its latitude and
 20 longitude in the geographical coordinate system and the direction of the antenna, called azimuth.

21 **Definition 5 (Path)** A path P between two nodes $v, w \in V$ is a sequence of edges $(e_1, \dots, e_n) \in$
 22 E^n such that $e_1 = (v, \cdot)$, $e_n = (\cdot, w)$ and $\forall i \in \llbracket 1, n-1 \rrbracket, \exists u \in V, e_i = (\cdot, u), e_{i+1} = (u, \cdot)$.

23 Finally, using the above definitions, the problem addressed in this paper can be defined as
 24 follows: given a cell phone trajectory T and the Multi-Layer graph G , the aim is to find the most
 25 likely path P in G that leads to the sequence of observed events (o_1, \dots, o_n) in T . This is obviously
 26 a map-matching problem.

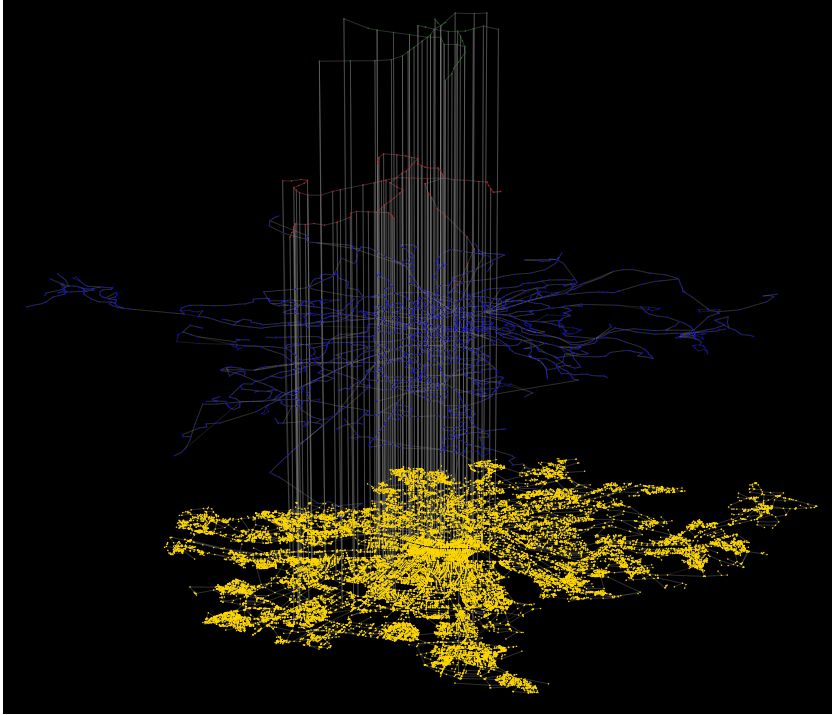
27 4. NETWORK MODELING

28 4.1. Multi-Layer Transportation Graph

29 The transportation network studied in the paper is the multimodal transportation system of the
 30 city of Lyon, France. The network is designed as a multiplex network G composed of four graph
 31 layers representing four transportation modes: road, bus, tramway and subway. The whole Multi-
 32 Layer network and its different layers is shown in Fig. 1a. The Python NetworkX library is used for
 33 multilayer modeling [26]. The graph and its different layers are built using multiple data sources
 34 and programming tools. The road network is generated via OSMnx [27], a Python library which
 35 creates NetworkX graphs from OSM data. Simplification of the road network topology derived
 36 from OSM is integrated as a facility in the library. The resulting road network corresponds to all
 37 drivable routes, representing the finest level of granularity that can be reached in road modeling.

38 Public transport layers have been generated using GTFS (Google Transit Feed Specification)
 39 data. We have performed some preprocessing steps (such as merging same public transport stops,
 40 which are in different directions) to obtain a reliable graph. Finally, cross-layers are added between
 41 layers to obtain the final multiplex graph structure. Between public transport layers, cross layers
 42 are defined as connections at transfer stops between public transport lines (this information is

1 contained in the GTFS transfer file). In Asgari *et al.* [15], each subway station is connected to its
 2 closest road intersection for simplification matters. In light of a more realistic modeling, we prefer
 3 instead to link the road and public transport layers by using parking locations derived from Lyon
 4 OpenData [28]. The closest node of each parking location is thus connected to the closest public
 5 transport node.



(a) Visualization of the Multi-Layer transportation network. Four transportation modes are considered: subway (green nodes, upper layer), tramway (red nodes, mid layer), buses (blue nodes, mid layer) and road (yellow nodes, bottom layer). Cross-Layers (vertical grey edges) connect the different layers. See Fig. 1b for statistics related to each of these layers.

Layer	$ V $	$ E $	$\langle k \rangle$	$\langle l \rangle$ (km)	Source
Multi-Layer	29005	63088	4.39	0.14	OSM/GTFS
Subway	46	80	3.47	0.78	GTFS
Tramway	86	173	4.03	0.60	GTFS
Bus	2023	4495	4.44	0.46	GTFS
Road	26853	58340	4.34	0.11	OSM

(b) Main characteristics of each transportation layer and Multi Layer network: number of nodes $|V|$, number of edges $|E|$, average node degree $\langle k \rangle$ and average edge length in kilometer $\langle l \rangle$.

Figure 1: Lyon multimodal network: graphical representation (a) and main features (b)

6 4.2. Cellular Network

7 The cellular network considered in this study is composed of 13,306 antennas of the Orange mobile
 8 network operator, in the Metropole of Lyon region. Due to the overlapping of 2G, 3G and 3G+
 9 cellular networks, antennas from different layers have the same characteristics (longitude, latitude
 10 and azimuth). After filtering duplicates, the result is a cellular network of 3,706 antennas. How-
 11 ever, there are still antennas with the same spatial location (longitude and latitude) but different

1 azimuths. In order to improve the modeling of the cellular network, we propose a method join-
 2 ing traditional Voronoi tessellation with the azimuth information to remove spatial overlapping.
 3 Specifically, each antenna is translated of an infinitesimal distance in the direction of its azimuth.

4 After applying Voronoi tessellation to the azimuth-corrected set of antennas, we consider the
 5 new location of the antenna as the barycenter of the polygon representing the Voronoi cell, as
 6 visually reported in Fig. 2. Compared to the simple Voronoi tessellation, as applied in [15], our
 7 coverage model is about three times more segmented by taking into account the azimuth of the
 8 antennas (i.e., the area covered by each antenna is on average three times smaller in our approach
 9 than in a traditional Voronoi tessellation).

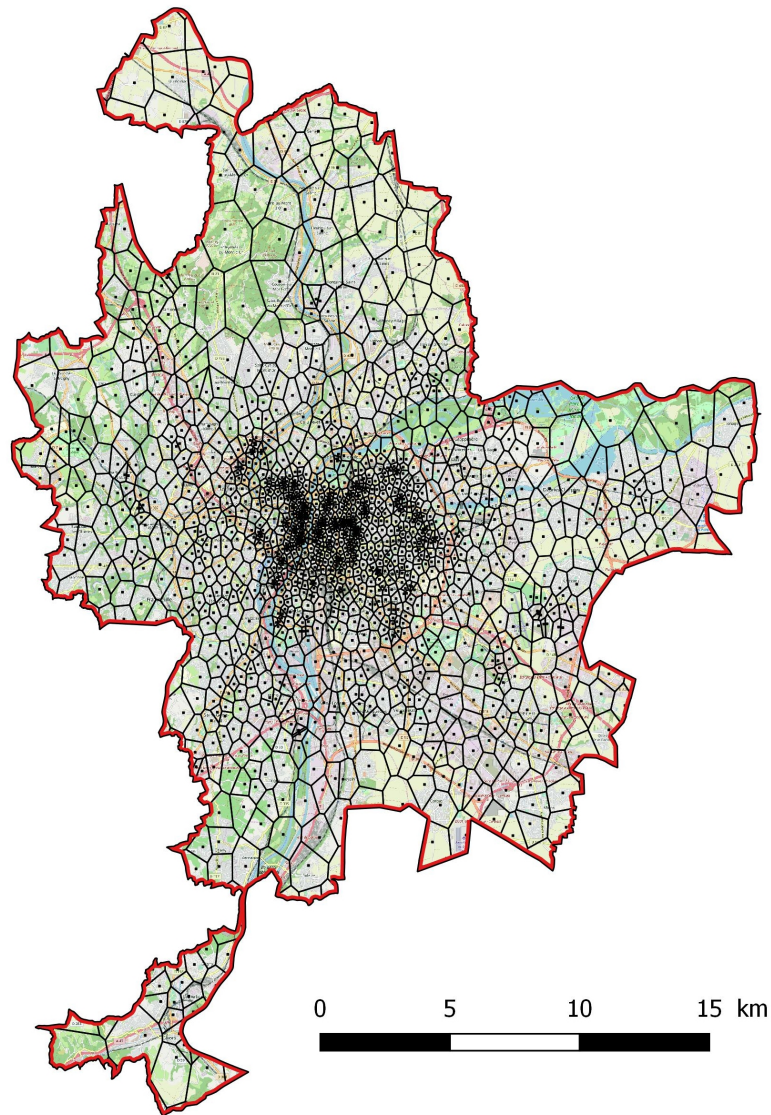


Figure 2: Cellular network of the Grand Lyon

10 5. METHODOLOGY

11 The following section reports on the main methodological background characterizing our solution
 12 to perform map-matching of cellular network trajectories, as issued from individual anonymized

1 network signaling mobile phone passive data.

2 5.1. Hidden Markov Model

3 A Hidden Markov Model can be defined by a five-fold $\langle V, C, \pi, A, B \rangle$, where:

- 4 • $V = \{v_0, \dots, v_{N-1}\}$ is a set of states.
- 5 • $C = \{c_0, \dots, c_{M-1}\}$ is a finite state of possible emissions.
- 6 • π is the probability distribution of the initial state, given that π is a probability distribution:
7
$$\sum_{i=1}^N \pi(i) = 1.$$
- 8 • A is a set of transition probability . The probability to transit from hidden state v_i to hidden
9 state v_j is denoted as $\{a(v_i, v_j)\}$. Besides, $\forall v_i \in V, \sum_{v_j \in V} a(v_i, v_j) = 1.$
- 10 • B is a set of emission probability . The probability to emit c_k from hidden state v_i is denoted
11 as $\{b(v_i, c_k)\}$. Besides, $\forall v_i \in V, \sum_{c_k \in C} b(v_i, c_k) = 1.$

12 Our map-matching problem can be modeled with a Hidden Markov Model: hidden states are
13 modeled as the set of vertices (nodes) V from the Multi-Layer Transportation Graph. Emissions
14 are modeled as the set of antennas C from Cellular Network. Hidden Markov Model allows to
15 solve the following problem: given a sequence of observations (sequence of antennas on a cellular
16 trajectory), the model finds the most likely sequence of hidden states (sequence of nodes on the
17 transportation network).

18 5.2. HMM parameters

19 5.2.1. Initial Probability

20 As the definition of the initial probability, all the nodes in the transportation network are equally
21 assigned with a probability of $1/N$ with N representing the total number of nodes in the trans-
22 portation network:

$$23 \quad (1) \quad \pi(i) = \frac{1}{N}$$

24 5.2.2. Transition Probability

25 The transition probability corresponds to the probability that a mobile phone user moves, on
26 the underlying transportation network from hidden state v_i at time $t - 1$ to hidden state v_j at
27 time t . Various transition probabilities have been proposed in the literature. For instance, as in the
28 definition by Luo *et al.* [16], transition probability only depends on the network connectivity. The
29 one used by Thiagarajan *et al.* [13] depends instead on the distance between transportation nodes.
30 However, all of these approaches use road transportation network to define transition probability.
31 Thus, these definitions require to be adapted to the case of a multilayer network, in order to take into
32 account the attributes of each layer. Hence, we choose the definition proposed by Asgari *et al.* [15],
33 i.e., the transition probability depends on the average speed over an edge and the edge length.
34 Average speed values are used to define edge weights as follows (see Tab. 1):

$$35 \quad (2) \quad W_{ij} = \begin{cases} w_{ij} & \text{if } v_i \text{ and } v_j \text{ are adjacent in } G \\ 0 & \text{otherwise} \end{cases}$$

value of w_{ij}	Condition
1/80	$\Psi(v_i) = \Psi(v_j) = \textit{subway}$
1/25	$\Psi(v_i) = \Psi(v_j) = \textit{tramway}$
1/15	$\Psi(v_i) = \Psi(v_j) = \textit{bus}$
1/50	$\Psi(v_i) = \Psi(v_j) = \textit{road}$
1/10	$\Psi(v_i) \neq \Psi(v_j)$

Table 1: Edge classification and weights for multilayer transportation network G

1 Finally, the transition probability is defined as the inverse of the shortest path cost between
 2 two nodes v_i and v_j :

$$3 \quad (3) \quad a(v_i, v_j) = \left(\sum_{\forall (v_m, v_n) \in SP_{v_i v_j}} w_{mn} \cdot d(v_m, v_n) \right)^{-1}$$

4 where (v_m, v_n) is any edge on the shortest path $SP_{v_i v_j}$ connecting the two nodes v_i and v_j in
 5 graph G . The shortest path cost of $SP_{v_i v_j}$ is the sum of the distances over each edge (v_m, v_n)
 6 belonging to $SP_{v_i v_j}$ weighted by w_{mn} . $d(v_m, v_n)$ denotes the geodesic distance between each two
 7 nodes v_m and v_n .

8 5.2.3. Emission Probability

9 The emission probability corresponds to the probability that an individual user is in the hidden
 10 state v_i at time t given that an emission (e.g communication event at an antenna) c_k is observed on
 11 the cellular network at time t . In the literature on mobile phone data, various emission probability
 12 has been proposed. Luo *et al.* [16] define a score inversely proportional to the distance between
 13 the hidden state and the observation. Jagadeesh *et al.* [24] prefer to use a Gaussian distribution
 14 with zero mean and an empirically estimated standard deviation of the measurement error between
 15 hidden states and observations. Similarly to Asgari *et al.* [15], since detailed information regarding
 16 the underlying cellular network is unavailable, we use Voronoi tessellation to model the area covered
 17 by each antenna. Finally, the emission probability is defined as a decreasing function of the distance
 18 between the antenna location and the hidden state:

$$19 \quad (4) \quad b(v_i, c_k) = \begin{cases} 1 & \text{if } d_{ik} < r_{max} \\ \left(\frac{r_{max}}{d_{ik}} \right)^\beta & \text{if } r_{max} < d_{ik} < \tau \cdot r_{max} \\ 0 & \text{otherwise} \end{cases}$$

20 where d_{ik} is the geodesic distance between c_k and intersection v_i , and $\beta = \frac{\ln(10)}{\ln(\tau)}$ is the decreasing
 21 factor which has been defined to obtain an emission score ten times lower for $d_{ik} = \tau \cdot r_{max}$.
 22 Finally, $\tau \cdot r_{max}$ is a threshold corresponding to the maximum distance at which a cell phone can
 23 be covered by a given cellular antenna. Considering the fact that the communication power is
 24 generally proportional to the inverse square of the distances [25], coefficient $\beta = 2$, which leads to
 25 $\tau = 3$, is chosen.

26 5.3. Preprocessing

27 This preprocessing step aims at reducing the noise in cellular phone trajectory. This is a key step
 28 to improve map-matching accuracy process. The cleaning algorithm of our approach follows these
 29 three sequential steps:

- 1 1. apply a recursive look ahead filter [29] . This filter is based on the mobile phone travel speed
2 on the cellular network. If the speed is higher than a given parameter, the outlier record is
3 removed. In the algorithm, this speed is set at 500 km/h.
- 4 2. through investigation into the data, we have decided to aggregate records with a given thresh-
5 old of two minutes to reduce the oscillation effect (also called ping-pong effect) on the cellular
6 trajectory. Moreover, this value of two minutes is low enough to avoid losing information on
7 the cellular trajectory. The antennas detected within the threshold are replaced by a single
8 antenna, *i.e.*, the one closest to the coordinates of the barycenter of the diverse antennas.
- 9 3. remove consecutive records detected at the same antennas. We consider in this case that the
10 user is static, thus no information is lost by simply removing the record.

11 5.4. Map-Matching algorithm

12 After applying the cleaning algorithm described above, map-matching can be used on cleaned
13 cellular trajectory. Our approach is a two-steps map-matching algorithm, reported in Pseudo-
14 code 1.

15 The first phase consists in an optimized Viterbi algorithm [30]. The inputs of the Viterbi
16 process are the following: the transportation network modeled as a multiplex network G , the
17 possibles states (set of the nodes of G), the emissions (set of antennas from the cellular network),
18 the previously defined HMM parameters (Sec. 5.2) and the cellular trajectory. By calculating all
19 possible paths given the cellular trajectory, the Viterbi process output is the likely sequence of
20 graph nodes, one for each time instant in the input. For real time application, due to a large
21 number of states and emissions, the execution time of the Viterbi algorithm is critical [25]. The
22 standard Viterbi algorithm applied in a scenario with 6,110 states (less complex network than
23 the one used in the study), 3,706 antennas and a set of 2,300 observation sequences completed
24 in around two hours. We executed the algorithm on a server machine equipped with an Intel
25 Xeon E5-2640 2.4 GHz multi-core machine, with 56 virtual cores and 128 GB of DDR4 RAM. To
26 improve performance, we implemented an optimized version of the Viterbi algorithm by leveraging
27 the sparseness of cellular trajectory, which represents the core optimization for real time application
28 used by Algizawy *et al.* [25]. The optimization process consists in eliminating all multiplications
29 by zero thus reducing the search space by keeping only with emittable states from each state
30 observable. The execution time of the optimized Viterbi algorithm drops to 3 seconds instead of 2
31 hours to reconstruct a set of 2,300 observation sequences. By exploiting the optimized version of
32 the Viterbi algorithm, we run the map-matching algorithm on top of a much larger transportation
33 network, composed of 29,012 nodes. In order to further reduce time execution, multiprocessing
34 Python libraries such as Joblib have been used.

35 Concerning the second step of the map-matching algorithm, after inferring the most likely
36 states sequence using the optimized Viterbi implementation presented above, the final trajectory
37 is inferred by applying a traditional shortest path (Dijkstra) detection algorithm on the underlying
38 transportation graph between two consecutive nodes.

Procedure 1 Map-Matching algorithm

Input:Multi-modal Network, $G(V, E)$ States (Network Nodes), $V = \{v_0, \dots, v_{N-1}\}$ Emissions (Mobile Network Antennas), $C = \{c_0, \dots, c_{M-1}\}$ Cell phone trajectory, $T = (o_0, o_1, o_2, \dots, o_{l-1})$ where $o_k \in C$ and l is the length of the sequenceInitial probabilities, π_i such that $i \in V$ Transition probabilities, a_{ij} such that $v_i, v_j \in V$ Emission probabilities, b_{ik} such that $i \in V$ and $k \in C$ **Output:**Maximum probability, $OutputProb$ Most-likely expanded node sequence, $FinalPath = \langle V_{o_0}, \dots, V_{o_{l-1}} \rangle$ *First step: Optimized Viterbi Algorithm*

```

1:  $StateProb \leftarrow \{\}$ 
2:  $Path \leftarrow \{\}$ 
3: for all  $y$  in  $V$  do
4:    $StateProb[0][y] = \pi_y \cdot b_{y,o_0}$ 
5:    $Path[y] \leftarrow y$ 
6: end for
7: for  $t \leftarrow 0$  to  $l - 1$  do
8:   for all  $y$  in  $V | b_{y,o_t} \neq 0$  do
9:      $(Prob, Pred) \leftarrow (\max_{z \in V | a_{z,y} \neq 0} (StateProb[t-1][z] \cdot a_{z,y} \cdot b_{y,o_t}), z)$ 
10:     $StateProb[t][y] \leftarrow Prob$ 
11:     $NewPath[y] \leftarrow Path[Pred] + y$ 
12:   end for
13:    $Path \leftarrow NewPath$ 
14: end for
15:  $(Prob, Pred) \leftarrow (\max_{y \in V} (StateProb[l-1][y]), y)$ 
16:  $OutputProb \leftarrow Prob$ 
17:  $OutputPath \leftarrow Path[Pred]$ 

```

Second step: Final itinerary reconstruction

```

18:  $FinalPath \leftarrow OutputPath[0]$ 
19: for all  $k$  in  $OutputPath \setminus \{OutputPath[0]\}$  do
20:    $FromNode \leftarrow OutputPath[k-1]$ 
21:    $ToNode \leftarrow OutputPath[k]$ 
22:    $IntPath \leftarrow ShortestPath(FromNode, ToNode, G)$ 
23:    $FinalPath \leftarrow FinalPath + IntPath$ 
24: end for
25: return( $FinalPath, OutputProb$ )

```

1 6. EVALUATION

2 6.1. Datasets

3 In order to evaluate our approach, we consider two datasets related to the metropolitan area of
4 Lyon, France:

- 5 • We collected GPS traces for a group of users in Lyon metropolitan area, who agreed to be
6 jointly tracked via a GPS tracking smartphone app and by the mobile phone operator as
7 well. For such specific users, Orange provided mobile phone data traces, having the same
8 characteristics as described above. This dataset is used for the microscopic validation of our
9 approach, GPS traces being used as ground truth, as reported in Sec. 6.2.
- 10 • Anonymized individual mobile phone (passive) data provided by Orange, the major French
11 mobile phone operator covering the week from 9th Sept. 2015 to 15th Sept. 2015. This
12 dataset include records of all users who visit, on the same day, at least one base station in
13 two areas of Lyon, i.e., Part Dieu (PD) and Sainte-Foy (SF) zones. It is worth mentioning
14 that the identifiers of such users are not the same across different days for privacy reasons.
15 Only timestamps, anonymized user id, antennas id information are provided. This dataset
16 is used for the macroscopic validation of our approach, as reported in Sec. 6.3.

17 6.2. Microscopic validation

18 In order to validate our model for microscopic user mobility, we applied our HMM-based map-
19 matching on the cellular trajectories and compared the inferred trajectories with the corresponding
20 GPS traces, considered as ground truth. In Fig. 3, six results from the proposed approach are
21 shown.

22 In Fig. 3a and Fig. 3b, using a fine-grained cellular network and the multi-modal transport
23 network, our algorithm shows a good accuracy, despite the sparsity of the cellular trajectory, in
24 an urban context. Moreover, in Fig. 3a our algorithm can properly infer a trip on the public
25 transportation network (the transportation mode used is the tramway).

26 The algorithm is also particularly effective in accurately map-matching cellular trajectories
27 along major roads in inter-urban contexts, as clearly shown in Fig. 3c. Indeed, the complexity of
28 the map-matching problem is reduced when a user is moving in a non-urban environment, compared
29 to an urban context. In such situations, our model is highly accurate and can fairly reconstruct
30 such kind of multi-modal trajectories.

31 Finally, in Fig. 3f, an example of reduced accuracy of our map-matching solution is shown. In
32 case of multiple events in a short spatial range, our approach considers the user as mobile and
33 attempts to infer a path whereas she/he is static. This explains why some loops appear on the
34 inferred cellular trajectory.

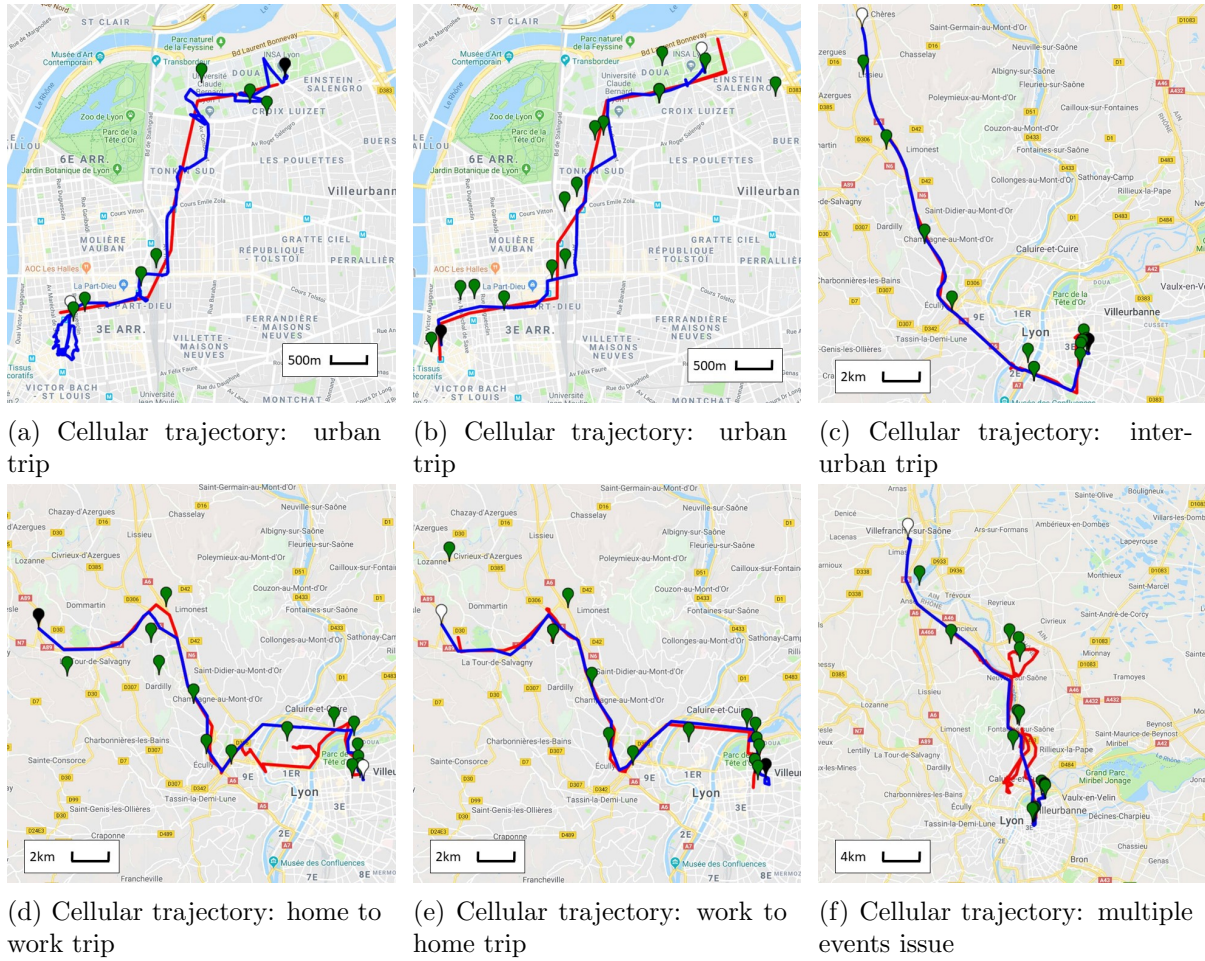


Figure 3: Set of cellular trajectory. The blue line represents the GPS trace (ground truth), the red line is the trajectory inferred by our approach (output of the map-matching algorithm), the green markers correspond to the cellular trajectory (input of the map-matching algorithm), the white marker represents the beginning of the trip and the black one, the end.

1 The approach has also been evaluated in a quantitative way. In order to estimate the per-
 2 formance of a map-matching approach, a largely used metric is the percentage of reconstructed
 3 paths (i.e., the reconstructed mobile phone itineraries) that correspond to ground-truth ones (i.e.,
 4 the GPS traces map-matched to the network). However, given that the map-matching problem is
 5 considered in our work according to a multi-modal setting and that our GPS traces do not include
 6 additional information, such as accelerometer data, traditionally leveraged in state-of-the-art ap-
 7 proaches (e.g., [31]) to more precisely compute the ground truth data (i.e., map-matching the GPS
 8 traces on the multi-modal transportation network via additional data collected via smartphone
 9 sensors), we cannot easily determine precise ground truth paths on the multi-modal transportation
 10 network (since the transportation mode is not known for our GPS raw data). As an alternative
 11 approach, we have thus computed the spatial accuracy between the sequence of nodes resulting
 12 from our map-matched mobile phone itineraries and the raw GPS traces, directly considered as
 13 our ground truth. As the considered trajectories do not match perfectly in terms of timestamps
 14 with the recorded events (i.e., GPS and mobile phone logs can be collected in different moments),
 15 we have considered two different kinds of matching. The first one is a spatial matching, in which

1 we assign the closest (in term of geodesic distance) GPS record to the mobile phone map-matched
 2 node. This approach tends to overestimate the performance of our approach. The second matching
 3 is a temporal one, where we assign the temporally-closest GPS record to the map-matched node,
 4 by looking at the timestamps of the two events. As the two timestamps do not perfectly match
 5 (dozens of seconds of difference) this approach tends to underestimate the performance of our so-
 6 lution. Based on such assumptions, if the geodesic distance between a map-matched node and the
 7 corresponding GPS point is under a certain value of spatial tolerance, we considered correct the
 8 map-matching for that node. Based on such twofold approach, the resulting performance of our
 9 approach has been computed, and can be safely said to lay between the two curves (i.e., cumulative
 10 distribution functions) reported in Fig. 4 (temporal and spatial matching). Therefore, by looking
 11 at the curves, we can say that for 70% correctly map-matched nodes, the spatial inaccuracy is
 12 bounded between 260m and 390m.

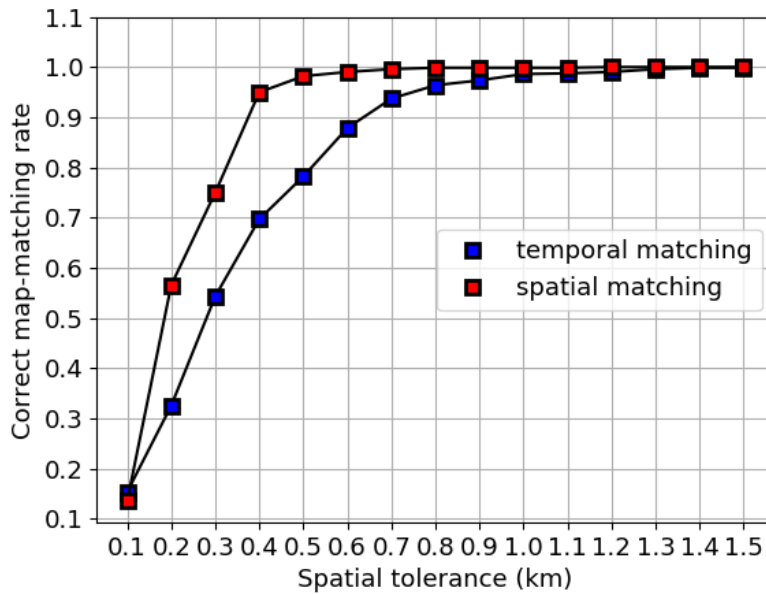


Figure 4: Spatial accuracy of the approach

13 6.3. Macroscopic validation

14 To validate our approach according to a more aggregate and larger-scale perspective, we propose
 15 in the following a macroscopic evaluation. Such an evaluation is organized in two different parts,
 16 depending on the way the aggregate ground truth is generated, as discussed in the following two
 17 sub-sections.

18 6.3.1. Validation via on OD matrix

19 The first validation only focuses on assessing the potential of mobile phone data to reconstruct
 20 aggregate mobility flows at larger spatial scales. To that purpose, we consider as our ground truth
 21 an OD matrix available at the sub-regional scale (i.e., smaller administrative areas). Such OD
 22 matrix has been independently reconstructed by our research laboratory by combining survey-
 23 based data from the administration of the city of Lyon and traffic counts data collected via loop
 24 detectors largely deployed in the city. This evaluation aims to compare the aggregate flows that

1 can be observed via our mobile phone data with respect to an external source of aggregate mobility
 2 flows (the OD matrix).

3 As previously discussed, our mobile phone traces are available for two different large zones of
 4 the city of Lyon, i.e., Part Dieu and Sainte Foy, each composed of multiple sub-zones called *IRIS*
 5 *sectors*. The analysis of the OD matrix is therefore limited to the IRIS sectors available in these
 6 two large areas of the city. Specifically, the areas of Part Dieu and Sainte Foy are divided into
 7 8 and 9 IRIS sectors, respectively. For the sake of clarity, we separately consider two OD matrices
 8 describing respectively: *i*) the daily number of trips generated from each IRIS sector inside the
 9 Part Dieu area to each IRIS sector in Sainte Foy (reported in Fig. 5a); *ii*) the daily number of trips
 10 generated from the IRIS zones inside the Sainte Foy area to each IRIS zone in Part Dieu (reported
 11 in Fig. 6a).

12 By using mobile phone data, we reconstructed similar matrices, by simply counting the trips
 13 (whose duration is lower than 1 hour) starting in one of the IRIS zone of Part Dieu and ending in
 14 one of the sectors of Sainte Foy (for the direction PD to SF). The same approach has been used in
 15 the other direction (SF to PD). The mobile-phone-based matrices are reported in Fig. 5b and in
 16 Fig 6b for the direction PD to SF and SF to PD, respectively.

17 To compare the two pairs of matrices (obtained in the two directions), we performed a linear
 18 regression and compute the R^2 linear regression coefficient (see Fig. 5c and Fig. 6c, respectively).

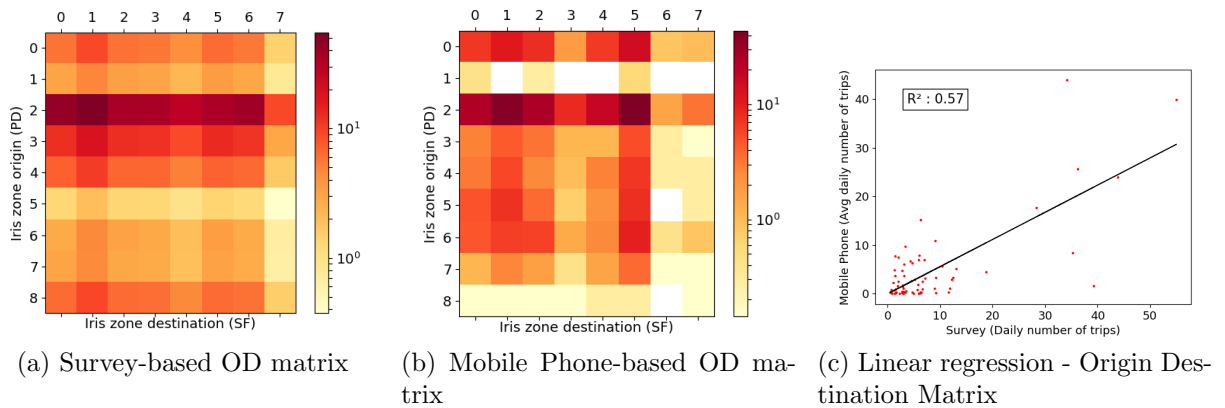


Figure 5: Spatial analysis for direction Part Dieu (PD) to Sainte Foy (SF)

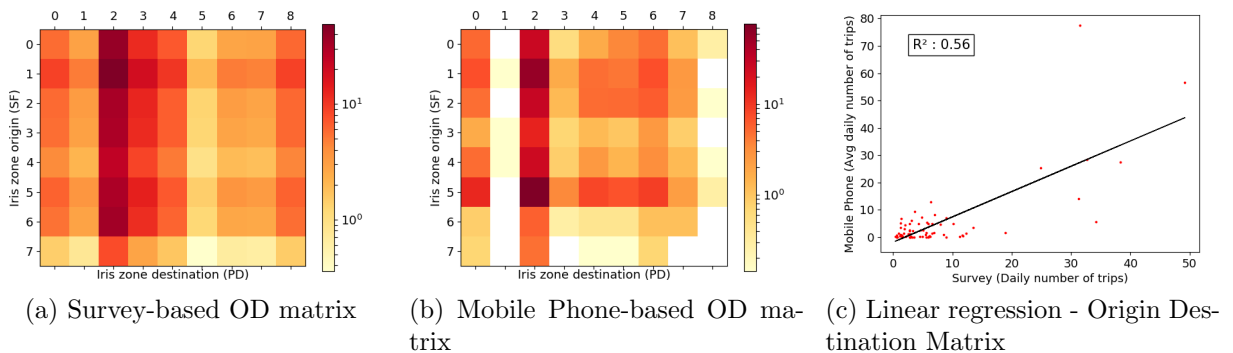


Figure 6: Spatial analysis for direction Sainte Foy (SF) to Part Dieu (PD)

1 The non-negligible levels of the observed R^2 linear regression coefficient (i.e., higher than 0.5 in
 2 the two scenarios) shows that mobile phone data possess the potential to capture the main trends
 3 of the survey-based OD Matrix. Nevertheless, we state that such levels of correlation could be
 4 significantly increased by considering a larger set of mobile phone data (only one week is analyzed
 5 leading up with only about 2000 trips per direction). In addition, we remark that we have considered
 6 a very simplistic definition of mobile phone trips, i.e., mobile phone traces whose first event takes
 7 place in Part Dieu (resp. Sainte-Foy) and whose last event takes place in Sainte-Foy (resp. Part
 8 Dieu), with a duration of the trip under a given threshold (1 hour). This definition does not include
 9 a formal identification of the starting/ending point of the trip (which would require the user to be
 10 stationary for a longer time), thus including in the count many trips that just pass through the
 11 two zones of interest. Such trips are indeed not considered in the survey-based OD matrix, thus
 12 inducing very likely a bias in the mobile phone reconstructed matrix.

13 Finally, we have compared the daily temporal demand profile from our OD matrix and the daily
 14 profile of trips from mobile phone data between Part Dieu and Sainte-Foy. A strong similitude can
 15 be observed in both directions, indicating that the daily temporal profile observed via mobile phone
 16 data are comparable to those observed based on survey data. In addition, given that Part Dieu is
 17 an area with a strong presence of offices and commercial activities, while Sainte-Foy is mainly a
 18 residential area, the mobile phone data temporal profiles can be interpreted as follow: most of the
 19 trips between Part Dieu and Sainte Foy come from people who live in Sainte Foy, go to work to
 20 Part Dieu in the morning (peak on Fig.8b between 8 and 9 am, in the direction Sainte-Foy to Part
 21 Dieu) and come back home in the afternoon (peak on Fig.7b between 4 and 7 pm, in the direction
 22 Part Dieu to Sainte-Foy). Thus, despite the few cellular trajectories analyzed (only one week of
 23 mobile phone data), our mobile phone data captures well-known activity-based mobility patterns,
 24 such as home-to-work travel.

25 With both spatial and temporal analysis, the relevance of the mobile phone data to study
 26 macroscopic mobility has been proven.

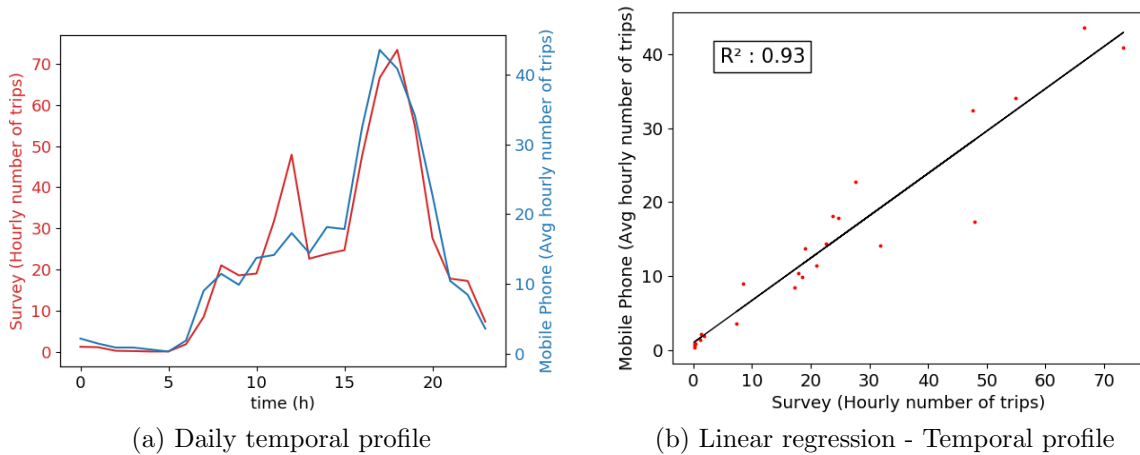


Figure 7: Temporal analysis for direction Part Dieu (PD) to Sainte Foy (SF)

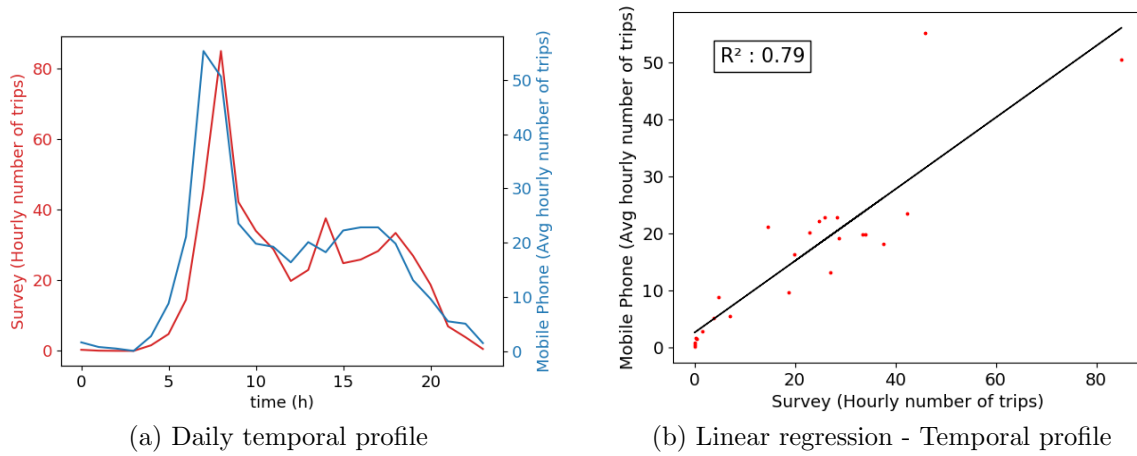


Figure 8: Temporal analysis for direction Sainte Foy (SF) to Part Dieu (PD)

6.3.2. Validation via simulation

Our second approach for macroscopic validation focuses instead on the proposed mobile phone map-matching approach, leveraging a simulation-based solution to identify multiple multi-modal popular paths (as well as the traffic shares over them) connecting the two observed areas. The evaluation aims to answer the following question: is our algorithm able to properly infer the distribution of flows over the most-traversed paths between the two considered areas of Part-Dieu and Sainte-Foy? In order to determine, in a fairly realistic way, a sort of ground truth describing the common paths between Sainte-Foy (SF) and Part-Dieu (PD), we used a combination of several tools. First, we used an A* shortest-path (SP) algorithm, based on the work of [32]. This SP algorithm uses a heuristics-directed search and it includes link penalties for multi-path search. It also incorporates a link penalty depending on a hierarchical description of the network. This method provided us with a set of routes efficient for cars only. We completed these results with the Google Map itineraries calculation, in order to confirm the results produced by the SP algorithm and to add supplementary routes for public transportation and bicycles. For public transportation, we also relied on the website of the SYTRAL, the public transportation authority in Lyon. At last, our choices were confirmed by our knowledge of the city. Especially, for the SF to PD direction, we added a supplementary route which seemed to be reliable, even if it was proposed neither by our A* nor by Google Map algorithm.

Assigning users to the different alternative paths is not straightforward. In our case, we derived the assignment coefficients results from a length-based C-logit approach following the work of [33]. The C-logit model solves a Stochastic User Equilibrium problem by considering both the cost of each alternatives and the commonality factor between alternatives. The cost is the mean travel time, provided as a static data. A numerical parameter β , presented in the above-mentioned article, was set to 70. The θ parameter on which the logit formula relies was set to $\theta = 0.009$ (see [33]). To determine θ , we ran a static traffic simulation on the city of Lyon-Villeurbanne in which we tried to minimize the difference between observed and modeled flow, while calibrating θ . Observations were taken from loop detectors and furnished by the city authorities.

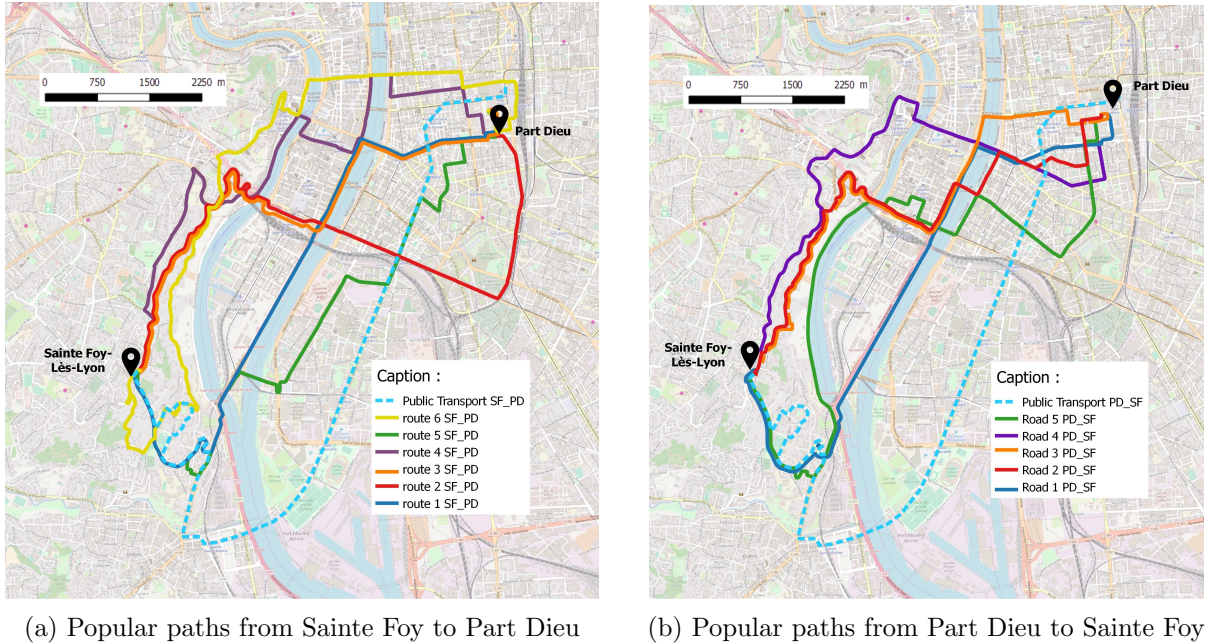


Figure 9: Popular paths between Part-Dieu to Sainte-Foy

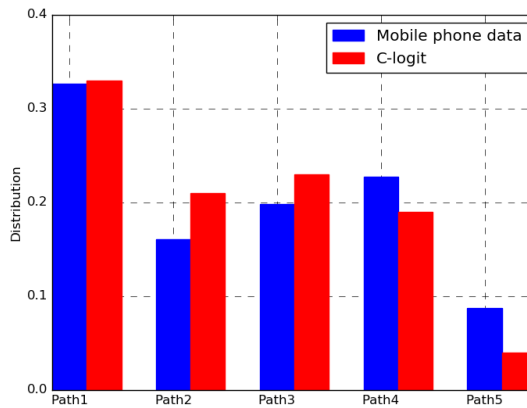


Figure 10: Macroscopic flow between Part-Dieu and Sainte-Foy

1 In Fig. 10, we report on the comparison between macroscopic flows, as inferred from static traffic
 2 simulations (in red) and the aggregated results retrieved by using our cellular-trajectory-based map-
 3 matching approach (in blue). It is worth highlighting that, given the biases and incertitude present
 4 in both approaches to compute path distribution, expecting a perfect match between the two
 5 approaches is rather unrealistic. However, we believe this comparison can provide qualitative and
 6 global insights on the capacity of our solution to properly match trajectories on the multi-modal
 7 transport network. Also, it complements the previously described microscopic validation, which
 8 has already proven a good hit-rate at an individual scale.

9 As a first interesting result, our approach does not lead to completely unrealistic and unexpected
 10 traffic flows, like for instance all cellular trajectories matching with only one or two expected
 11 popular paths from simulations. Besides, the two approaches lead to consistent results in terms of

1 popular/unpopular paths: *Path 1* is the most used in both cases, while *Path 5* has the lowest score
2 in the two approaches as well.

3 **7. CONCLUSION AND DISCUSSION**

4 Cellular-network signaling data have the great potential to provide fine-grained spatio-temporal
5 information to reconstruct users' mobility at both microscopic and macroscopic scales. In this
6 paper, we performed an empirical study, based on real cellular traces collected in the city of Lyon,
7 France, by a major telecommunication operator, aimed at investigating such potential. We devel-
8 oped a HMM-based map-matching algorithm for mapping sparse and noisy cellular trajectories to
9 the underlying transportation network.

10 As a practical basis for our approach, we developed automatic tools to build a large multi-modal
11 transportation network.

12 Taking into account the azimuth of antennas allows to increase cellular network segmentation.
13 Network modeling at a fine-level of granularity allows to properly apply map-matching in urban
14 complex environment.

15 By providing a formal definition of the HMM parameters, our methodology follows three main
16 steps: the cleaning process, an optimized implementation of the Viterbi algorithm, and the deter-
17 mination of the shortest path on the sequence of nodes returned by Viterbi algorithm.

18 To validate our approach, we have analyzed an original case study, related to French city of
19 Lyon, by leveraging both real cellular traces collected by a major network operator and GPS data
20 collected via a mobile phone application. This data has been leveraged to perform a microscopic
21 validation proving the accurate map-matching capability of our approach, even in a complex urban
22 context. Moreover, we have demonstrated the possibility to retrieve popular paths between two areas
23 by comparing the spatial distribution of flows as computed by both our approach and simulations.

24 Future directions should consider improvement with dynamic HMM parameters, in order to
25 build a transition matrix depending on actual traffic conditions. In addition, some limitations of
26 our algorithm have been shown in relation to oscillations (or ping-pong effect) in the user's com-
27 munication activity. Therefore, a better understanding of this recurrent phenomenon is required.
28 The latter should allow to create an advanced filtering approach to remove this oscillation effect
29 from cellular trajectories and further improve the map-matching accuracy. Finally, we aim to de-
30 velop a comprehensive comparative evaluation of our solution with respect to the most relevant
31 state-of-the-art approaches related to mobile phone data map-matching.

32 **ACKNOWLEDGMENTS**

33 The authors would like to thank Razvan Stanica and Marco Fiore for their valuable contribution
34 in the modelling choices and in the analysis of mobile phone data. The authors would like to thank
35 Orange for providing the mobile phone data used in this study. This work has been supported by
36 the French ANR project, PROMENADE, grant number ANR-18-CE22-0008.

37 **AUTHOR CONTRIBUTION STATEMENT**

38 The authors confirm contribution to the paper as follows: study, conception and design: LB, AF,
39 JK, NEEF; analysis and interpretation of results: LB, AF, JK, NEEF; LB was the lead writer of
40 the manuscript. All authors reviewed the results and approved the final version of the manuscript.

41 **References**

42 [1] Diala Naboulsi, Marco Fiore, Stephane Ribot, and Razvan Stanica. Large-Scale Mobile Traffic
43 Analysis: A Survey. *IEEE Communications Surveys & Tutorials*, 18(1):124–161, 21 2016.

- 1 [2] Marta C. González, César A. Hidalgo, and Albert-László Barabási. Understanding individual
2 human mobility patterns. *Nature*, 453(7196):779–782, 6 2008.
- 3 [3] Md. Shahadat Iqbal, Charisma F. Choudhury, Pu Wang, and Marta C. González. Development
4 of origin–destination matrices using mobile phone call data. *Transportation Research Part C:
5 Emerging Technologies*, 40:63–74, 3 2014.
- 6 [4] Angelo Furno, Marco Fiore, Razvan Stanica, Cezary Ziemlicki, and Zbigniew Smoreda. A Tale
7 of Ten Cities: Characterizing Signatures of Mobile Traffic in Urban Areas. *IEEE Transactions
8 on Mobile Computing*, 16(10):2682–2696, 10 2017.
- 9 [5] Angelo Furno, Marco Fiore, and Razvan Stanica. Joint spatial and temporal classification of
10 mobile traffic demands. In *IEEE INFOCOM 2017 - IEEE Conference on Computer Commu-
11 nications*, pages 1–9. IEEE, 5 2017.
- 12 [6] Rex W Douglass, David A Meyer, Megha Ram, David Rideout, and Dongjin Song. High
13 resolution population estimates from telecommunications data. *EPJ Data Science*, 4(1):4, 12
14 2015.
- 15 [7] Hugo Barbosa, Marc Barthelemy, Gourab Ghoshal, Charlotte R. James, Maxime Lenormand,
16 Thomas Louail, Ronaldo Menezes, José J. Ramasco, Filippo Simini, and Marcello Tomasini.
17 Human mobility: Models and applications. *Physics Reports*, 734:1–74, 3 2018.
- 18 [8] John R. B. Palmer, Thomas J. Espenshade, Frederic Bartumeus, Chang Y. Chung, Necati Er-
19 can Ozgencil, and Kathleen Li. New Approaches to Human Mobility: Using Mobile Phones
20 for Demographic Research. *Demography*, 50(3):1105–1128, 6 2013.
- 21 [9] Injong Rhee, Minsu Shin, Seongik Hong, Kyunghan Lee, Seong Joon Kim, and Song Chong.
22 On the Levy-Walk Nature of Human Mobility. *IEEE/ACM Transactions on Networking*,
23 19(3):630–643, 6 2011.
- 24 [10] Yu Zheng, Xing Xie, and Wei-Ying Ma. Understanding Mobility Based on GPS Data, 9 2008.
- 25 [11] Y U Zheng. Trajectory Data Mining: An Overview. *ACM Trans. On Intelligent Systems and
26 Technology*, 6(3), 2015.
- 27 [12] Paul Newson and John Krumm. Hidden Markov map matching through noise and sparseness.
28 In *Proceedings of the 17th ACM SIGSPATIAL International Conference on Advances in Ge-
29 ographic Information Systems - GIS '09*, page 336, New York, New York, USA, 2009. ACM
30 Press.
- 31 [13] Arvind Thiagarajan, Lenin Ravindranath, Hari Balakrishnan, Samuel Madden, and Lewis
32 Girod. Accurate, low-energy trajectory mapping for mobile devices, 2011.
- 33 [14] Wei Wu, Yue Wang, Joao Bartolo Gomes, Dang The Anh, Spiros Antonatos, Mingqiang
34 Xue, Peng Yang, Ghim Eng Yap, Xiaoli Li, Shonali Krishnaswamy, James Decraene, and
35 Amy Shi Nash. Oscillation Resolution for Mobile Phone Cellular Tower Data to Enable Mo-
36 bility Modelling. In *2014 IEEE 15th International Conference on Mobile Data Management*,
37 pages 321–328. IEEE, 7 2014.
- 38 [15] Fereshteh Asgari, Alexis Sultan, Haoyi Xiong, Vincent Gauthier, and Mounîm A. El-Yacoubi.
39 CT-Mapper: Mapping sparse multimodal cellular trajectories using a multilayer transportation
40 network. *Computer Communications*, 2016.

- 1 [16] An Luo, Shenghua Chen, and Bin Xv. Enhanced Map-Matching Algorithm with a Hid-
2 den Markov Model for Mobile Phone Positioning. *ISPRS International Journal of Geo-*
3 *Information*, 6(11):327, 2017.
- 4 [17] Mohammed A. Quddus, Washington Y. Ochieng, and Robert B. Noland. Current map-
5 matching algorithms for transport applications: State-of-the art and future research directions.
6 *Transportation Research Part C: Emerging Technologies*, 15(5):312–328, 2007.
- 7 [18] Christopher E White, David Bernstein, and Alain L Kornhauser. Some map matching algo-
8 rithms for personal navigation assistants. *Transportation Research Part C: Emerging Tech-*
9 *nologies*, 8(1-6):91–108, 2 2000.
- 10 [19] Meng Yu. Improved positioning of land vehicle in its using digital map and other accessory
11 information. *The Hong Kong Polytechnic University*, 2006.
- 12 [20] Dragan Obradovic, Henning Lenz, and Markus Schupfner. Fusion of Map and Sensor Data in
13 a Modern Car Navigation System. *The Journal of VLSI Signal Processing Systems for Signal,*
14 *Image, and Video Technology*, 45(1-2):111–122, 11 2006.
- 15 [21] Mohammed A. Quddus, Robert B. Noland, and Washington Y. Ochieng. A High Accuracy
16 Fuzzy Logic Based Map Matching Algorithm for Road Transport. *Journal of Intelligent Trans-*
17 *portation Systems*, 10(3):103–115, 9 2006.
- 18 [22] Jong-Sun Pyo, Dong-Ho Shin, and Tae-Kyung Sung. Development of a map matching method
19 using the multiple hypothesis technique. In *ITSC 2001. 2001 IEEE Intelligent Transportation*
20 *Systems. Proceedings (Cat. No.01TH8585)*, pages 23–27. IEEE, 2001.
- 21 [23] Gunnar Schulze, Christopher Horn, and Roman Kern. Map-Matching Cell Phone Trajecto-
22 ries of Low Spatial and Temporal Accuracy. *IEEE Conference on Intelligent Transportation*
23 *Systems, Proceedings, ITSC*, 2015-Octob:2707–2714, 2015.
- 24 [24] George R. Jagadeesh and Thambipillai Srikanthan. Online Map-Matching of Noisy and Sparse
25 Location Data with Hidden Markov and Route Choice Models. *IEEE Transactions on Intel-*
26 *ligent Transportation Systems*, 18(9):2423–2434, 2017.
- 27 [25] Essam Algizawy, Tetsuji Ogawa, and Ahmed El-Mahdy. Real-Time Large-Scale Map Matching
28 Using Mobile Phone Data. *ACM Transactions on Knowledge Discovery from Data*, 11(4):1–38,
29 7 2017.
- 30 [26] Aric A Hagberg hagberg, lanlgov Los, Daniel A Schult, and Pieter J Swart swart. Exploring
31 Network Structure, Dynamics, and Function using NetworkX. *PROCEEDINGS OF THE 7TH*
32 *PYTHON IN SCIENCE CONFERENCE SCIPY*, 2008.
- 33 [27] Geoff Boeing. OSMNX: New Methods for Acquiring, Constructing, Analyzing, and Visualizing
34 Complex Street Networks. *SSRN Electronic Journal*, 5 2016.
- 35 [28] Données métropolitaines du Grand Lyon, url : <https://data.grandlyon.com/>.
- 36 [29] A. Sultan. *Méthodes et outils d’analyse de données de signalisation mobile pour l’étude de*
37 *la mobilité humaine*. PhD thesis, Ecole doctorale : Informatique, Télécommunications et
38 Electronique de Paris, 2016.

- 1 [30] A. Viterbi. Error bounds for convolutional codes and an asymptotically optimum decoding
2 algorithm. *IEEE Transactions on Information Theory*, 13(2):260–269, 4 1967.
- 3 [31] Jingmin Chen and Michel Bierlaire. Probabilistic multimodal map matching with rich smart-
4 phone data. *Journal of Intelligent Transportation Systems: Technology, Planning, and Oper-*
5 *ations*, 2015.
- 6 [32] Peter Hart, Nils Nilsson, and Bertram Raphael. A Formal Basis for the Heuristic Determination
7 of Minimum Cost Paths. *IEEE Transactions on Systems Science and Cybernetics*, 4(2):100–
8 107, 1968.
- 9 [33] Zhong Zhou, Anthony Chen, and Shlomo Bekhor. C-logit stochastic user equilibrium model:
10 formulations and solution algorithm. *Transportmetrica*, 8(1):17–41, 1 2012.