# Transportation Research Record

## Dynamic Estimation of Urban Zonal Speed from Mobile Sensing Data and Macroscopic Paths
--Manuscript Draft--

| Full Title: | Dynamic Estimation of Urban Zonal Speed from Mobile Sensing Data and Macroscopic Paths |
|---|---|
| Abstract: | In the past decade, call detail records (CDR), a specific type of mobile phone data, has been proven to be an innovative source of human and urban mobility data. When compared to GPS probe data, CDR data often suffer from temporal sparsity and spatial imprecisions. However CDR data are also usually more massive and offer a larger population coverage. Given this advantage, we wonder whether CDR data could be used as a reliable source to estimate the mean traffic speed dynamics.<br><br>In this article, we propose an innovative method, based on the partitioning of the urban area in reservoirs and on the analysis of basic trip features. The mean traffic speed in each reservoir is derived from the solution of a linear system that clusters individual trips per macro-path to make more robust travel time estimations. This method is fast and simple to implement in real-time but requires a prior analysis of the main trip patterns at the city and reservoirs scales.<br><br>We apply this methodology on a large GPS dataset, that we reduce and downsample in order to reproduce the typical CDR data temporal characteristics. As we want to compare the method results on both GPS and CDR, monitoring the data downsampling process allows us to limit discrepancies. Moreover, the original raw GPS data also provides the ground truth reference. This experiment in a very controlled environment is a first step towards studies dealing with real CDR data. |
| Manuscript Classifications: | Data and Information Technology; Traffic Flow Theory and Characteristics AHB45; Planning and Forecasting; Traffic Flow Theory and Characteristics AHB45 |
| Manuscript Number: | |
| Article Type: | Presentation |
| Order of Authors: | Manon Seppecher |
| | Ludovic Leclercq |
| | Angelo Furno |
| | Delphine Lejri |
| | Amélie-May Lupinski |

# Dynamic Estimation of Urban Zonal Speed from Mobile Sensing Data and Macroscopic Paths

**Manon Seppecher**[*]
Univ. Lyon, ENTPE, IFSTTAR, LICIT
UMR_T 9401, F-69518, LYON, France
CITEPA, PARIS, France
Tel: +33 4 72 04 72 95
Email: manon.seppecher@entpe.fr

**Ludovic Leclercq**
Univ. Lyon, ENTPE, IFSTTAR, LICIT
UMR_T 9401, F-69518, LYON, France
Tel: +33 4 72 04 77 16
Email: ludovic.leclercq@entpe.fr

**Angelo Furno**
Univ. Lyon, ENTPE, IFSTTAR, LICIT
UMR_T 9401, F-69518, LYON, France
Tel: +33 4 72 04 77 16
Email: angelo.furno@ifsttar.fr

**Delphine Lejri**
Univ. Lyon, ENTPE, IFSTTAR, LICIT
UMR_T 9401, F-69518, LYON, France
Tel: +33 4 72 04 77 16
Email: delphine.lejri@entpe.fr

**Amelie-May Lupinski**
Univ. Lyon, ENTPE, IFSTTAR, LICIT
UMR_T 9401, F-69518, LYON, France
Tel: +33 4 72 04 77 16
Email: amelie-may.lupinksi@entpe.fr

\* Corresponding author

1  August 1, 2019

**ABSTRACT**

In the past decade, call detail records (CDR), a specific type of mobile phone data, has been proven to be an innovative source of human and urban mobility data. When compared to GPS probe data, CDR data often suffer from temporal sparsity and spatial imprecisions. However CDR data are also usually more massive and offer a larger population coverage. Given this advantage, we wonder whether CDR data could be used as a reliable source to estimate the mean traffic speed dynamics.

In this article, we propose an innovative method, based on the partitioning of the urban area in reservoirs and on the analysis of basic trip features. The mean traffic speed in each reservoir is derived from the solution of a linear system that clusters individual trips per macro-path to make more robust travel time estimations. This method is fast and simple to implement in real-time but requires a prior analysis of the main trip patterns at the city and reservoirs scales.

We apply this methodology on a large GPS dataset, that we reduce and downsample in order to reproduce the typical CDR data temporal characteristics. As we want to compare the method results on both GPS and CDR, monitoring the data downsampling process allows us to limit discrepancies. Moreover, the original raw GPS data also provides the ground truth reference. This experiment in a very controlled environment is a first step towards studies dealing with real CDR data.

## 1. INTRODUCTION

In the past few years, mobile phone data has been proven to be an innovative, accessible and very rich source of information about human mobility. Those data, passively generated by mobile phone users while communicating, moving around the network, or activating and deactivating their devices, are collected by the mobile phone data provider either for billing or network management purposes. The resulting massive databases form an incredibly large source of information about communication activities, but also mobility behaviours of the urban populations. In this paper, we wonder if a specific type of mobile phone data, call detail records, could be used in order to estimate the mean traffic speed dynamics at a zonal scale in an urban area.

Call detail records (CDR) register for each communication event (calls, messages or data browsing) the following information: the user's unique identification key, a timestamp characterising the start of the event, and the location of the base station antenna that processed the event. This basic structure, relating together user id, time and location, confers to those data a strong potential for mobility analysis. Compared to GPS data classically used in mobility and traffic analysis, CDR data present several advantages. The mobile phone penetration rates are generally high. For traffic studies, phone data usually provides a better spatial coverage as they are not restricted to a single users category, unlike taxis GPS probe data. However, CDR data also have drawbacks. First, as the positional information is obtained at the antennas level, the spatial precision depends on the base station network. Second, the data generation depends on the users' communication activities and behaviours, making the temporal acquisition of data uneven and sparse it time; in particular users with little communication activity will generate fewer location data and their mobility will be especially difficult to estimate.

Using CDR data, Gonzalez et al. (1) analysed the behavioural rules intrinsic to human mobility in order to construct realistic individual mobility models. At a more aggregated scale, CDR data were also used to estimate origin destination matrices (2, 3, 4), proving that they are an interesting alternative to the traditional and costly transportation surveys and census data at an urban scale. In Toole et al. (5) the origin-destination matrices estimated thanks to the aforementioned methods were later assigned onto the road network in order to estimate the traffic load and to identify the specific origin-destination flows that caused the highest share of it. Two other types of mobile phone data, handovers and location area updates data, have been specifically exploited for traffic analysis studies. Handovers data record the network-centred events that happen when the communication connection of a user engaged in a call session is transferred from one antenna to another due the user's movement on the network. Location area updates data record for every device (including idle ones) the transfer of the communication connections from one identified group of antennas to another one when the devices owners move over the network. Those data are less dependant to communication activities than CDR data and have a good spatial coverage, however they are not always accessible in practice, mainly for privacy reasons. Derrmann et al. (6) explored the potential of handovers data for the estimation of Macroscopic Fundamental Diagrams. On highway segments, Bar-Gera (7) used handovers data to estimate the traffic speed while Janecek et al. (8) analysed the travel time using both handovers and location area updates.

But, despite this extensive literature in the field of mobility and traffic, the specific question of speed estimation from CDR data appears to have been barely studied. This is probably partly due to the imprecisions of CDR data, which can seem *a priori* insufficient for traffic speed estimation compared to GPS or handovers and location area updates data. Aside from the challenge that speed estimation from CDR data represents, or from the insight it can give on the traffic dynamics

over time at a large urban scale and for a large sampled population, this speed estimation presents a high potential in the field of traffic emission estimations. CDR data are accessible, massive, and representative of large population groups unlike most GPS probe datasets derived from taxi trajectories. While using GPS data for those estimation requires important complementary data and costly scaling up processes (9), an estimation using CDR data can be an efficient and light alternative. In that direction, Li et al. (10) introduced a method to estimate traffic related pollutant emissions from CDR data at a regular zonal scale. However the method used for estimating the emission factors (essential for the emission calculation and which depends on the traffic speed) was not made explicit. Consequently, and to the best of our knowledge, the estimation of the traffic speed dynamics from CDR data still needs to be investigated. This is what we propose to do in this article.

We aim to determine if CDR data can be used in order to estimate dynamically the traffic speed at a zonal scale in an urban area. As a first step into that direction, we propose an innovative method based on the partitioning of the urban area in reservoirs and on the identification of clusters of individual trips sharing similar simple features such as macro-path and arrival time period. Providing that a good estimation of the trip lengths at the city and reservoirs scale is known, the clusters allow to build at each studied time period a simple linear system, that, once solved, returns an estimation of the traffic speed. We chose the city of Lyon, France, as our case study and experiment our method on a GPS dataset that we downsample in order to get the same statistical characteristics than CDR data, following a method initially developed by Chen et al. (11) to estimate the spatial biases introduced by CDR data in mobility studies. As the GPS data use does not come from taxi compagnies, it does not present the specific bias mentioned above. This approach first allows to validate the general method with unbiased data, and then to evaluate the impact of the bias introduced by the downsampling and attempt to limit this impact. Eventually, the objective is to make this estimation method robust enough to be used on real CDR data rather than simulated ones.

This article is organised as follows : Section 2 exposes the principle of our approach and describes the proposed methodology. Section 3 presents our case study, as well as the exploited data. Section 4 focuses on the results we reached. Finally, Section 5 concludes with the achieved results, the limitations of this work and the on-going perspectives.

## 2. METHODOLOGY
**Work scope and definitions**
We propose a method to estimate the temporal dynamics of the mean traffic speed at zonal scale, using a limited and selected set of individual trips characterised by a low level of information. Our method relies on partitioning an urban network into a set of regions, called reservoirs, and on dividing the time dimension into a set of time intervals of equal duration. Our objective is to estimate the mean traffic speed in each reservoir for every time period.

The notions on which the methodology relies are defined below.

**Reservoirs:** We call reservoirs the result of the partitioning of the studied urban area into smaller regions of equivalent sizes and homogeneous characteristics in term of city fabric, demography, road network or traffic dynamics.

**Micro-path:** We call micro-path the ordered sequence of road segments traveled by a sampled vehicle along a trip.

**Macro-path** We call macro-path the ordered sequence of reservoirs crossed by a sampled

1  vehicle along a trip. A macro-path is a scaled representation of a micro-path at the reservoirs level.
2      **Micro-trip:** We call a micro-trip the representation of a sampled individual vehicle trip
3  according to the following features : *(trip Id, vehicle Id, micro-path, exact arrival time, total*
4  *travel time)*
5      **Macro-trip:** We call macro-trip the representation of a sampled individual vehicle trip
6  according to the following features : *(trip Id, vehicle Id, macro-path, arrival interval, total travel*
7  *time)*. A macro-trip is a simplified representation of a micro-trip at the reservoirs scale and with a
8  coarser time resolution.
9      The simplification process that reduces a micro-trip to its macroscopic representation al-
10 lows to group similar (but not identical) micro-trips into clusters based on both the macro-path
11 and the arrival interval. The methodology we propose in the next part relies on those macro-trip
12 categories.

13 **Linear System Construction**
14 *Generic System*
15 Let us consider that the studied area is divided in a given number $R$ of reservoirs. Let us also
16 consider that the time span of the experiment is divided into a set of equal time intervals. Let $I_t$
17 (of size $n_t$) be the set of macro-trips reaching their destinations during period $t$, and let $I_{t,P}$ (of size
18 $n_{P,t} \leq n_t$) be the subset of $I_t$ of macro-trips that match macro-paths $P$.
19     For a macro-trip $i$ of $I_{t,P}$, its traveled time along $P$ can be written as the sum of the traveled
20 times $T^i_{r,P}$ of $i$ in each reservoirs $r$ of $P$:

$$T^i_P = \sum_{r \in P} T^i_{r,P} \tag{1}$$

21     We assume that the traffic speed in reservoir $r$ during interval $t$ $V_{r,t}$ is constant and homo-
22 geneous. This gives:

$$T^i_P = \sum_{r \in P} \frac{L^i_{r,P}}{V_{r,t}} \tag{2}$$

23     with $L^i_{r,P}$ the total distance traveled in reservoir $r$ by instance $i$ along the macro-path $P$.
24     Summing on the $n_{P,t}$ equations 2 characterising to the $n_{P,t}$ macro-trips of $I_{t,P}$, we get:

$$\sum_{i=1}^{n_{P,t}} T^i_P = \sum_{i=1}^{n_{P,t}} \sum_{r \in P} \frac{L^i_{r,P}}{V_{r,t}} \tag{3}$$

$$= \sum_{r \in P} \sum_{i=1}^{n_{P,t}} \frac{L^i_{r,P}}{V_{r,t}} \tag{4}$$

$$n_{P,t} \bar{T}_{P,t} = \sum_{r \in P} n_{P,t} \frac{\bar{L}_{r,P,t}}{V_{r,t}} \tag{5}$$

$$\bar{T}_{P,t} = \sum_{r \in P} \frac{\bar{L}_{r,P,t}}{V_{r,t}} \tag{6}$$

$$\tag{7}$$

1     With $Y_{r,t} = \frac{1}{V_{r,t}}$, and assuming that $\bar{L}_{r,P,t}$ is independent of $t$ , we get:

$$\bar{T}_{P,t} = \sum_{r \in P} \bar{L}_{r,P} Y_{r,t} \tag{8}$$

2     For the time interval $t$, this equation characterises one given macro-path. Applying it to the
3 whole set of macro-paths observed at interval $t$, we can construct the following system $S_t$:

$$S_t = \left\{ \bar{T}_{P,t} = \sum_{r \in P} \bar{L}_{r,P} Y_{r,t} \qquad \forall P \right. \tag{9}$$

4     $S_t$ is linear with $R$ variables made of as many equations as macro-paths observed during
5 interval $t$. This system is usually overdetermined as the number of possible macro-paths is higher
6 than the number of reservoirs. Assuming that $\bar{L}_{r,P}$ can be estimated exogenously, and as $\bar{T}_{P,t}$ can
7 be derived from the macro-trip information, the system can be solved using for exemple a least
8 square optimisation method. The inverse of the solution $Y_{r,t}$ will correspond to the traffic speed $V_{r,t}$
9 for each reservoir $r$.
10     The method we just exposed presents the advantage of considerably reducing the complex-
11 ity of the problem. First, individual trips are gathered for each time period per macro-paths which
12 provides a robust estimation of $\bar{T}_{P,t}$. Second, we can select among all macro-paths the most rep-
13 resentative ones before looking for the system solution, which permits to reduce the system size.
14 Third and foremost, we only need a very limited time information about trips, basically the de-
15 parture and the arrival time. This makes the method perfectly suitable for a CDR data input. The
16 drawback is that it requires a robust and exogenous estimation of mean traveled distance within
17 each reservoir $\bar{L}_{r,P}$. The spatial definition of the reservoirs will also be crucial as it may drive
18 completely different patterns for different macro-paths inside the same reservoir. This topic is still
19 under research.

20 *Parameter reliability and temporal bias*
21 Theoretically, GPS data are precise enough to provide accurate estimations of $\bar{T}_{P,t}$ (corresponding
22 to the mean traffic time along $P$ when reaching destination during $t$) and $\bar{L}_{r,P}$ (the total traveled
23 distance in reservoirs $r$ along path $P$). However, when it comes to using CDR data (or simply
24 sparser GPS data) the direct estimation of those parameters might not be as reliable. Indeed, CDR
25 data include both temporal and spatial biases, due to users potentially large inactivity times and
26 weak spatial resolution of base stations network. Using such data to estimate the system parameters
27 implies reproducing those biases on the parameters and therefore twisting the system and its results.
28     Throughout the study, we make the assumption that $\bar{L}_{r,P}$ can be estimated exogenously
29 thanks to another reliable method (probe vehicles, surveys, automatic network analysis, etc) and
30 therefore do not consider the potential impact of distance bias on the data. We rather focus on the
31 impact of the temporal bias introduced by the low acquisition rate of mobile phone communica-
32 tion events. For a given mobile phone user, the good characterisation of their mobility is highly
33 dependant on their communication activity. The more active they are, the more location data are
34 collected, and statics phases (stays) can be separated from the in-between mobility phases. More
35 precisely, if the user is very active, the departure and arrival time will be estimated with a limited
36 time imprecision. On the contrary, for a barely active user, detecting their movements and there-
37 fore estimating the correct travel time will be much more difficult. This means that the estimated
38 traveled time of a trip observed from CDR data is actually higher than the real one.

Consequently, for a macro-trip $i$ of $I_{P,t}$ obtained with CDR data, we have:

$$T_{P,obs}^i = T_P^i + \varepsilon^i \tag{10}$$

where $T_P^i$ is the real travel time of trip $i$, $T_{P,obs}^i$ is the observed traveled due to data impreci-sions and $\varepsilon^i$ is the temporal bias that exists in-between. Summing over $I_{P,t}$, this gives:

$$\bar{T}_{P,obs} = \bar{T}_P + \bar{\varepsilon}_t \tag{11}$$

Reinjecting $\bar{T}_P$ in (8), we get :

$$S_t' = \left\{ \bar{T}_{P,obs} - \bar{\varepsilon}_t = \sum_{r \in P} \bar{L}_{r,P} Y_{r,t} \qquad \forall P \right. \tag{12}$$

Thus, if we are able to estimate the average temporal bias induced by the human dependent sampling rate of mobile data, then it becomes possible to de-skew the estimation the travel time along each macro-path $P$, and therefore to correct the system.

**Macro-paths and equations selection**

In order to reduce the system size and to limit the incompatibility risks between equations, we propose to filter out of the system the equations that are the least reliable, that is to say that have the least significance among the macro-trips. This significance can indeed be very uneven from one macro-path to another. While some macro-paths are observed for many macro-trips, some others might be representative of only one individual. In this later case, the spatial and temporal mean values will represent a single trip and so will be weakly reliable. Thus, we implement a significance threshold corresponding to the minimum number of times that a macro-path should be observed during period $t$ in order for the corresponding equation to be selected in the system $S_t$. This threshold is set to 5 trips in this study.
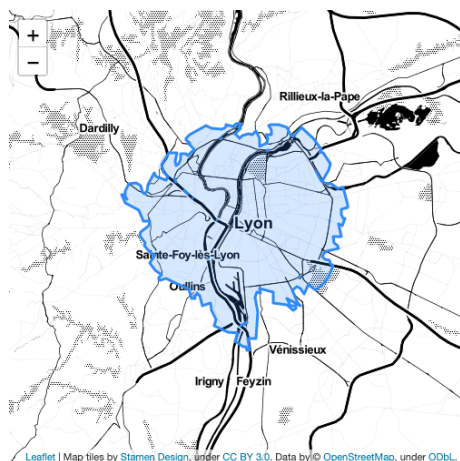
**3. CASE STUDY AND DATA DESCRIPTION**

**Case study**

We select the city of Lyon, France, as our case study. The study area includes both Lyon and the next-by municipality of Villeurbanne, which is located inside Lyon's ring road. It is displayed in figure 1a. We have parted this territory into a number of 11 distinct reservoirs. 10 of them divide the inner city, while the last one corresponds to the ring road. Those reservoirs are displayed in figure 1b. For the inner reservoirs, their geometry are constructed as an aggregation of "IRIS" units, the smallest census unit defined by the French National Institute of Statistics and Economic Studies, which fills both demographic and geographic criteria. This aggregation of reservoirs is made so that it is consistent with Lyon's road network and its traffic characteristics.
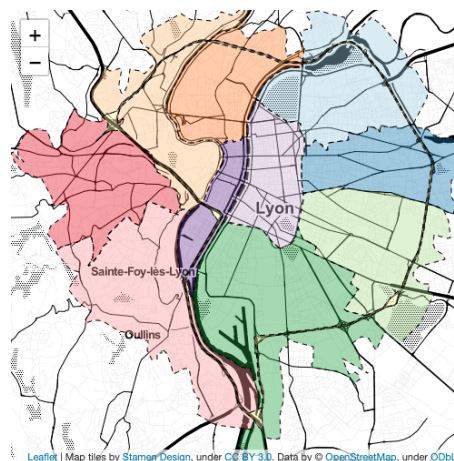
**Working data**

*Input data*

In this study, we intentionally work with a GPS dataset, which we progressively simplified and downsampled in order to recreate the temporal biases of mobile phone data. This approach presents two advantages. First, it allows us to use the raw GPS data as ground truth for the estimation of the traffic speed. Second, it enables us to progressively reduce the data quality, and thus to understand how the jamming process impacts the quality of our results.

(a) Considered area in Lyon, France

(b) Spatial partitioning of the studied area in 11 reservoirs, with the ring road as one of them

**FIGURE 1** : Visualisation of the studied area and reservoirs

1  The data used consist of one year of GPS traces over the Great Lyon area. For this study
2  however, we select the data of one typical weekday, 2018 February 5. The data was provided by
3  a European navigation system provider, and the traces come from vehicles equipped with their
4  navigation system technology. This means that the dataset does not present the bias observed with
5  GPS traces from taxis. Moreover, as each trace corresponds to a vehicle, there is no need to filter
6  out pedestrian or cyclist travellers as usually in CDR-based studies. Though this will make the
7  speed estimation process easier, it also means that the data we simulate do not exactly reproduce
8  CDR data that gather users making no distinction between their means of transport.
9  The data structure is slightly different from typical GPS datasets, in that the GPS sampled
10  locations of the monitored vehicles have already been map-matched by the data provider onto the
11  road network. Thus, each GPS trace is made of an ordered sequence of traveled network links
12  rather than a succession of geolocated points. The original data have the following structure : each
13  monitored *vehicle* is associated to the different *trips* that it performed, and each *trip* corresponds to
14  a sequence of *observations*. An *observation* relates each link of the trip to the considered vehicle,
15  the timestamp of the entrance of the vehicle on this link, the speed of the vehicle on this link as well
16  as the distance coverage that the vehicle does of this link (the first and last link of the trajectory
17  might not be fully traveled, for exemple because of parking spots). This information will especially
18  allow us to reconstruct both the travel distance and the travel time on each link.
19  As the data coverage is actually larger than the study area defined above, we filter out the
20  data tracks that do not enter the inner Lyon area, and split into smaller trips the ones that go beyond
21  those limits. Additionally, we further preprocess the data by filtering out any aberrant trace or static
22  vehicle in order to obtain a clean and reliable data set.

23  *Ground truth data*
24  First, the raw data are used for computing the ground truth reservoir speed values throughout time.
25  For each time interval $t$, and each macro-path $P$, the ground truth speed $V_{r,t}$ is calculated as :

$$V_{r,t} = \frac{TTD_{r,t}}{TTT_{r,t}} \tag{13}$$

1  where $TTD_{r,t}$ and $TTT_{r,t}$ are respectively the total traveled distance and time by the sam-
2  pled users in reservoir $r$ during interval $t$. This information can easily be derived from the *obser-*
3  *vations* dataset.

4  *Distance parameters*
5  To estimate the travel distance $\bar{L}_{r,P}$ we resort here to the full resolution of GPS data. In practice,
6  this estimation is done thanks to a few days of data and supposed known thereafter. It will be used
7  as a parameter in the systems constructions. On the other hand, the speed estimation is applied
8  over another selection of days where only CDR-like data are supposed to be available.

9  *Spatiotemporal aggregation : from trips to macro-trips*
10 Reshaping the original preprocessed data, we build a $T_0$ dataset of micro-trips as defined in sec-
11 tion 2, i.e. using the five following features: *trip id*, *vehicle id*, *micro-path*, *exact arrival time*, and
12 *total traveled time*. From this initial dataset $T_0$, we are going to derive the two degraded datasets
13 used for our experiments.
14       The first of those datasets corresponds to the macro-trips data set. First, we apply the spatial
15 aggregation which relies on the city partitioning proposed above. This partition of the studied area
16 allows us to relate each of the road network links to the reservoir it is included in. Subsequently,
17 each vehicle *micro-path*, i.e. the succession of network links making up the trip itinerary, can be
18 summarised as the ordered sequence of reservoirs that the vehicle crossed along its way.We call
19 this ordered sequence of crossed reservoirs *macro-path*, in opposition to the micro-paths made of
20 ordered sequence of links. Then, for each trip, we reduce the exact arrival time to the time interval
21 it belongs it. In this study, we select time intervals of 15 minutes. Thus, an arrival time at *11:22:36*
22 will be transformed into *11:15:00*.
23       This aggregating process allows us to transform trips into macro-trips of the shape *(trip
24 Id, vehicle Id, macro-path, arrival interval, total travel time)*. We will call this new dataset $T_1$.
25 For each vehicle trip, the exact arrival time is replaced by the closest anterior quarter of hour. In
26 that way we part the day in a predefined grid of 96 time intervals. We can now identify similar
27 vehicle trips, i.e vehicle trips that travel along the same macro-path and end in the same time
28 interval. So far, we have erased the knowledge about the exact travelled path and exact arrival
29 time. However, it is important to stress here that this data degradation has not impacted yet the
30 information about total traveled time. It simply allows to categorise the different trajectories and
31 to relate them together, in anticipation of the construction of the linear systems.

32 *Temporal downsampling*
33 In order to evaluate the potentials of our speed estimation method with mobile phone data, we
34 apply on the dataset $T_1$ a downsampling process which aims to simulate the temporal imprecisions
35 of mobile phone data compared to GPS data. To do so, we adapt the method developed in Chen
36 et al. (11) for spatial bias analysis to our data and problematic. We introduce in the data temporal
37 gaps that reproduce the characteristic inter-event time of CDR data. We specifically focus on the
38 impact of our downsampling on the detected beginning and end of the trip. This downsample
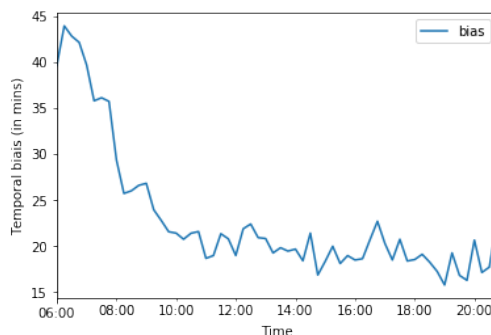
**FIGURE 2** : Temporal bias evolution throughout the day

1 introduce a bias both on the trip start time and on the trip arrival time, and thus implies an increase
2 of the global traveled time. The detected arrival interval can also change depending on how big
3 was the resulting time gap at the end of the trip. This new resulting dataset will be named $T_2$.
4     In reality, using CDR data not only introduces temporal bias but also distance imprecisions
5 due to the sampling done at the antenna location and the lack of intermediary trajectory points. For
6 the moment, we neglect this limit and consider that the macro-paths remain observable.

## 4. RESULTS AND DISCUSSION
8 In this section, we present the results of the application of our method to a specific day, 2018
9 February 5.

**Bias characterisation**
11 Using both $T_1$ and $T_2$, we can compare for each macro-trip the observed travel time (biased) to the
12 real travel time. The difference characterises the temporal bias at an individual level. Measuring
13 this bias for every macro-trip and averaging their values for every time period, we can estimate the
14 mean bias $\varepsilon_t$ throughout the day. The result of this estimation with our data is displayed in figure 2.
15 In the early morning the mean temporal bias is quite high, up to 45 minutes of difference between
16 the real travel time and the observed one. This average bias strongly drops between 6 am and 10
17 am before stabilising around 20 minutes. This difference is due to the changes around that time in
18 the communication activity rates of people. This is understandable as this time bias depends itself
19 on the human communication rhythms, much sparser at night. Thanks to this bias characterisation,
20 it becomes possible to understand the impact of the CDR data temporal characteristics on the
21 estimation of the travel time.

**Method application to macro-trip data**
23 As the macro-trip data contains only a few trips during the night, the study time span is restrained
24 to the day hours in-between 6 am and 9 pm. Taking this restriction into account, our macro-
25 trip dataset $T_1$ is made of 69 520 individual macro-trips, with a total of 1 862 distincts macro-
26 paths observed. Out of those macro-paths, only 82 of them are observed more than 5 times (our
27 significance threshold) during at least one time period, and therefore will be represented by system
28 equations.
29     The results of the method application on dataset $T_1$ are displayed in figure 3. The first ten
30 reservoirs of each figures correspond to the urban inner reservoirs, while the last one characterises

**FIGURE 3** : Results of methodology applied to dataset $T_1$ of macro-trips



**FIGURE 4** : Results of methodology applied to dataset $T_2$ of macro-trips without bias correction

**FIGURE 5** : Results of methodology applied to dataset $T_2$ of macro-trips with bias correction

1   the speed dynamics on Lyon's ring road. On each subplot, the ground truth traffic speed is repre-
2   sented in orange. The blue curves characterises the results obtained from the system resolution.
3   As it happens that the resolution does not converge, the returned speed signal can return aberrant
4   values. We apply a physical filter to remove the few diverging speed values. A speed of 60 km/hour
5   is selected as the threshold for detecting aberrant values in every reservoir except on the ring road
6   where the limit is raised up to 90km/hour. Values above those thresholds are arbitrarily replaced by
7   the previous consistent one. This explains the step that we can particularly observe for the reservoir
8   10 during the evening. This physical filtering removed the aberrant speed estimations. However
9   the obtained signals still present high volatility. In order to smooth the results, we apply a second
10   filter represented in green on the plots. Assuming that the signal noise corresponds to higher speed
11   peaks, we apply a low-pass frequency filter. The Butterworth filter is selected and applied with a
12   forward-backward filtering method (once in each direction).
13         From figure 3 we observe that the system generally renders correctly the speed dynamics.
14   This is especially true for reservoirs 5, 7, 8, 9 and 10. We also notice an overestimation of the
15   traffic speed in particular for the reservoirs 0 and 2, which we do not explain yet.
16         In figure 4 and figure 5 we display the results of the method applied on dataset $T_2$ respec-
17   tively with and without the correction of the temporal bias. Figure 5 highlights the very negative
18   effect of the temporal bias introduced by the downsampling process on the results. This evidences
19   that the temporal bias induced by human uneven and large communication times makes CDR data
20   unusable as they are to estimate the traffic speed, even at a regional level. After correcting the
21   system however, we can observe in figure 5 that the estimation tend to better fit the reference curve
22   again.
23         Table 1 summarises the achieved results, by comparing the root mean square error over the
24   full day in each reservoir, in-between the filtered results (green curves) and the reference values.

1  The results obtained for $T_1$ and de-skewed $T_2$ are sensibly similar. For some reservoirs, the results
2  obtained with $T_2$ after bias correction are slightly better than with $T_1$. Those results should be
3  further explored in order to explain this slight improvement, which can maybe come from the
4  biased arrival time used in $T_2$ for the system resolution.

| Reservoirs | Original macro-trips ($T_1$) | $T_2$ without bias correction | $T_2$ with bias correction |
|:---:|:---:|:---:|:---:|
| 0 | 16.58 | 18.56 | 15.16 |
| 1 | 3.38 | 13.53 | 3.58 |
| 2 | 14.25 | 14.99 | 14.46 |
| 3 | 6.81 | 16.30 | 5.53 |
| 4 | 3.66 | 19.66 | 2.91 |
| 5 | 3.11 | 18.70 | 3.58 |
| 6 | 7.58 | 22.17 | 6.57 |
| 7 | 2.89 | 18.58 | 2.68 |
| 8 | 2.16 | 12.40 | 1.97 |
| 9 | 4.89 | 14.35 | 3.66 |
| 10 | 9.04 | 48.55 | 9.06 |

6  **TABLE 1** : Synthesis of the RMSE for each reservoir according to the dataset used

7  **5. CONCLUSION**
8  In this article we have proposed a new methodology to estimate the dynamics of regional traffic
9  speeds from mobile sensing data. Our method is based on the partitioning of the urban area in
10  reservoirs, and on the identification of groups of sampled trips sharing common macro-paths and
11  arrival time period. This clustering of macro-trips provides a robust estimation of the travel times
12  along each path. Cross-referencing the macro-path travel time estimations within a linear system
13  allows us to estimate the traffic speed dynamics, providing that exogenous travel distance data are
14  available. The structure of this method is particularly fitted to a CDR data input, as it requires very
15  little temporal or itinerary information at the individual level and to takes into account the inherent
16  temporal bias that characterises those data.
17      Thanks to the application of our method to the set of GPS trips reduced to minimal tem-
18  poral and path information, we could validate the global methodology. We then downsampled the
19  trips temporal dimension in order to simulate the human uneven communication rhythms and to
20  reproduce the temporal limits of CDR data. The direct application of the method the downsampled
21  data showed that correcting the estimation of the travel time was a necessary condition to properly
22  estimate the speed. Comparing the original and biased trip data specifically enabled us to estimate
23  the temporal bias that CDR data can present compared to the ground truth GPS data. Removing the
24  average bias permitted to correct the travel time estimation and to obtain satisfactory speed results
25  compared to the ground truth data.
26      In our future works, we plan to explore the sensibility of our method the different param-
27  eters such as the size of the reservoirs, the time period duration, or the significance threshold of
28  the system equation. About the specific question of the reservoir size, it is clear that the larger
29  the reservoirs are, the more likely they are to encompass very different micro-paths under a same
30  macro-path. In such cases, it might be interesting to divide a macro-path into several sub-macro-

1 paths, analysing the potential modes in travel distance and time distributions. Last but not least,
2 our objective is eventually to apply the developed method on real CDR data.

## 1 ACKNOWLEDGEMENTS

## 4 References

5 [1] Gonzalez, M. C., C. A. Hidalgo, and A.-L. Barabási, Understanding individual human mo-
6     bility patterns. *Nature*, Vol. 453, 2008, pp. 779 EP –.

7 [2] Iqbal, M. S., C. F. Choudhury, P. Wang, and M. C. González, Development of origin–
8     destination matrices using mobile phone call data. *Transportation Research Part C: Emerging
9     Technologies*, Vol. 40, 2014, pp. 63 – 74.

10 [3] Çolak, S., L. P. Alexander, B. G. Alvim, S. R. Mehndiratta, and M. C. González, Analyzing
11     cell phone location data for urban travel: current methods, limitations, and opportunities.
12     *Transportation research record: Journal of the transportation research board*, , No. 2526,
13     2015, pp. 126–135.

14 [4] Alexander, L., S. Jiang, M. Murga, and M. C. González, Origin–destination trips by purpose
15     and time of day inferred from mobile phone data. *Transportation Research Part C: Emerging
16     Technologies*, Vol. 58, 2015, pp. 240 – 250, big Data in Transportation and Traffic Engineer-
17     ing.

18 [5] Toole, J. L., S. Çolak, B. Sturt, L. P. Alexander, A. Evsukoff, and M. C. González, The path
19     most traveled: Travel demand estimation using big data resources. *Transportation Research
20     Part C: Emerging Technologies*, Vol. 58, 2015, pp. 162 – 177, big Data in Transportation and
21     Traffic Engineering.

22 [6] Derrmann, T., R. Frank, F. Viti, and T. Engel, Estimating urban road traffic states using
23     mobile network signaling data. In *2017 IEEE 20th International Conference on Intelligent
24     Transportation Systems (ITSC)*, 2017, pp. 1–7.

25 [7] Bar-Gera, H., Evaluation of a cellular phone-based system for measurements of traffic speeds
26     and travel times: A case study from Israel. *Transportation Research Part C: Emerging Tech-
27     nologies*, Vol. 15, No. 6, 2007, pp. 380–391.

28 [8] Janecek, A., D. Valerio, K. A. Hummel, F. Ricciato, and H. Hlavacs, The Cellular Network
29     as a Sensor: From Mobile Phone Data to Real-Time Road Traffic Monitoring. *IEEE TRANS-
30     ACTIONS ON INTELLIGENT TRANSPORTATION SYSTEMS*, 2015.

31 [9] Shang, J., Y. Zheng, W. Tong, E. Chang, and Y. Yu, Inferring Gas Consumption and Pollution
32     Emissions of Vehicles throughout a City, 2014.

33 [10] Li, Q., Y. Cheng, F. Ding, X. Wan, and B. Ran, Citywide Hourly Traffic Emissions Estimation
34     Using Cellular Activity Data. In *TRB 95th Annual Meeting Compendium of Papers*, 2016.

35 [11] Chen, G., S. Hoteit, A. Carneiro Viana, M. Fiore, and C. Sarraute, Enriching sparse mobility
36     information in Call Detail Records. *Computer Communications*, 2018.