

# TRAFFIC FLOW ESTIMATES INFERRED FROM MOBILE PHONE NETWORKS

## ABSTRACT

Traffic volume is the main parameter to quantify transportation demand. Accurate measurements of volumes are essential for the successful implementation of traffic management. The development of models to infer traffic volume data constitutes an important step in the construction of automatic tools for such purposes. Different methodologies and techniques have been developed to estimate traffic volumes. The most common examples are based on time series, historical average, neural networks or non-parametric regression models, using data collected by fixed count stations (e.g. loop detectors), which have cost and coverage limitations.

Recently, mobile phone systems appeared as a complementary solution to fixed stations for increasing coverage areas and accuracy without requiring expensive infrastructure investment. Mobile phone systems include processes which automatically keep databases updated with phone location. In the case of an on-call phone, its location information is updated at the base station (cell) to which the phone is connected. For billing purposes, the system also records parameters related to the call, such as start time/end time, duration or identity of the originating cell. Additionally, when a phone with a call in progress moves from one cell to another, the call must be handed over to the new cell in order to provide uninterrupted service. This transference is controlled by the handover process which inserts a record into the databases to update them with content consisting of phone-ID, cell identity or timestamp, among others. Thus, the analysis of call data records is useful to get information related to the movement of phones between cells. From this we can infer vehicle volume data from this easily-accessible, fast and low-cost method.

This paper explores the use of call data generated by on-call phones moving over a network to estimate vehicle volume at certain points of intercity roads. The main purpose of this research is the definition, analysis and assessment of feasible models to infer traffic volume using anonymous call data collected automatically by mobile phone network operators. This will produce a set of forecasts for the site at which the data was collected.

The development of the models was based on call data records of phones provided by a network operator in a group of cells, which define a set of inter-cell boundaries; and traffic volume statistics collected by counting stations (detectors) at the same points. The collected data was divided into two sets: calibrating data and testing data. The calibrating data set was used to estimate the model's parameters, completely defining the models to predict traffic volumes. The estimated traffic volumes were then compared, in numerical and graphical

format, with the volume measurements collected by the detectors. Different indicators were calculated to evaluate the estimation performance, such as classic error measures (mean absolute error or mean absolute relative error), rank correlation (Spearman) or the linear correlation (Pearson) coefficient. Finally, the results using the testing data set showed that the proposed models produced reasonable estimates and are therefore effective for practical use, without requiring the installation of counting stations.

*Keywords: mobile phones, traffic flow, volume data, traffic monitoring.*

## **1. INTRODUCTION**

Accurate measurements of traffic flow are essential for the successful implementation of traffic management. Traditionally, traffic flow data are automatically collected on a subset of links in a network by fixed detectors embedded in the road, such as inductive loops, infrared, visual cameras, radars, lasers, etc. These measurements can be collected with little effort, although they are limited in terms of providing only localised information. Unfortunately, extending or modifying this type of infrastructure tends to be expensive due to the cost of installation and maintenance of detectors. New alternatives are therefore required to characterise traffic in a fast, accurate and continuous manner, at an acceptable cost.

Recently, tools for collecting accurate volume measurements have been introduced within new technologies. Among these new technologies is mobile phone-tracking, an emerging concept that has been studied during the last few years by academic and private entities in the transportation community. With the ubiquity and ever-increasing capabilities of mobile devices, mobile phones and their locations could potentially become a powerful source to derive traffic data. In modern societies mobile phones have reached high penetration rates, with many countries surpassing the 90 percent penetration rate (CIA World Factbook, 2009). Mobility management in mobile phone networks manages certain location data for proper performance. If a handset is located in a moving vehicle, the information collected from this phone can be used in traffic flow estimation. The main advantages of using mobile phones as probes are the low implementation cost and the need for no additional hardware. Mobile phone network-based solutions do not require any changes to handsets and are therefore able to provide better coverage than fixed detectors. However, the major disadvantage of this type of solutions is the establishment of a relationship between phone flows and vehicle flows since not every vehicle carries a phone (in buses several phones are onboard) and not every phone subscription is contracted from the mobile phone network operator which provides the data.

Thus, the problem of interest consists of finding new algorithms to make use of the data obtained from mobile phones, and, in particular, the vehicles acting as probes to estimate traffic flows on motorways. These algorithms should be developed by means of adequate treatments to relate both measures (phone counts and vehicle counts) through the use of mathematical models. In this paper a set of models to derive volumes of vehicles from (anonymous) mobile phone data is proposed and evaluated in order to be used in transport applications. The procedure used to process the incoming phone data requires a non-

negligible computing time and consequently, the proposed models are intended for applications in which the estimation process does not need to be done in real time (off-line).

## 2. MOBILE PHONE SYSTEMS

Mobile telephony facilitates the exchange of information (voice, data, video, image, and the like) for phone users, independent of time, location, and access arrangement. Mobility management enables a mobile phone system to find the locations of mobile phones so as to deliver incoming calls, messages, or data packets to said mobiles. In order to provide such mobility data, information concerning the location and characteristics of each mobile phone is held in a subscriber database known as a Home Location Register (HLR). To distribute the signalling load on the HLR, one or more subsidiary databases are used by the mobile network to temporarily hold profiles of roaming users (users outside their home area), which is known as a Visitor Location Register (VLR). This information is defined by the Technical Specifications 3GPP TS 23.003 and 3GPP TS 23.008, both published by the Third Generation Partnership Project (3GPP).

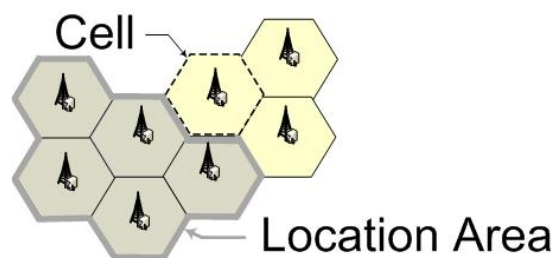


Figure 1 – Definition of a theoretic cell and location area

As Figure 1 shows, a set of base stations (BS) covers the mobile phone service area to connect a mobile phone to one corresponding mobile phone network (for making/receiving calls, sms, etc.). Each BS covers a radio coverage area called “cell”, which is theoretically represented by hexagonal shape. A cell is the basic geographic unit of a mobile phone system. The cell radius depends on many parameters, such as terrain, antenna type, power levels, or even topology and buildings around the antenna. Its size can vary from several hundred meters to several kilometers. The cell radius is smaller in urban areas, where population density is high, and more antennae are needed than in rural areas to provide communication services. A set of adjacent cells are grouped into one location area (LA).

A mobile phone network is characterised by two aspects of mobility. During a call conversation, a radio link connects the mobile phone to a base station (cell). When the user moves from the radio coverage area of one BS (or a cell) to the coverage area of another BS, a radio link in the new base station is required to continue the conversation. The handover is the procedure responsible for transferring that call. Thus, in the case of an on-call phone, the system always knows the base station (cell) to which said phone is connected. However, whenever the phone is turned on and not on-call (idle status) the network does not need to know the precise cell in which a phone is located, only a location area (LA). Therefore, a mobile phone in motion keeps the network informed about changes in the location area. If the phone moves from a cell in one location area to a cell in another

location area, the mobile phone performs a location update (LU) procedure to inform the network about the exact location of the phone. In addition, when a mobile phone is powered on, or if the phone does not produce any kind of activity, a periodic LU procedure will be triggered to identify the current LA of the phone.

Therefore, the mobile network always knows the Location Area Identification (LAI) of an idle phone and the CI (Cell ID) of a phone in communication (Calls, SMS, MMS, etc.). More comprehensive technical details regarding mobile networks can be found in the literature by Wakefield et al. (2007) and Yacoub, (2002).

### **3. INFERENCE OF VOLUME DATA USING MOBILE PHONES**

#### **3.1 State of the art**

The idea of using mobile phones to collect traffic information is a decade old but it has become increasingly widespread. Work reviewing the state of the practice of using mobile phones as a device for traffic monitoring can be found in Caceres et al. (2008), Qiu et al. (2007), and Smith and Fontaine (2006). In terms of volume data, the concept of using mobile phones as probes has been explored by various researchers working on simulated frameworks (Fontaine and Smith, 2005), as well as in field tests (Bekhor et al., 2008; Höpfner et al., 2007; Thiessenhusen et al., 2006). In all cases volume data would be associated to phone transit through a boundary area using processes related to mobility management, detecting boundary crossing rates either at inter-cell boundary level (handover) or at location area boundary level (LU procedure). Most of them focus on the handover event to detect phone transit through boundaries between two cells due to data characteristics (location accuracy, sample size, etc.). The changes from one base station/cell to another while the phone is involved in a call are recorded with timestamps by means of handover parameters, providing a measurement of the number of phones performing handovers in a given period, or in other words, a phone boundary-crossing rate. The main quantitative findings show that phone flow (calls) is closely related to the flow of vehicles, registered by loop detectors installed at those points in the network, and that there are similarities over the course of typical working days (Höpfner et al., 2007; Thiessenhusen et al., 2006; Reades et al., 2007). Typical flow peaks associated with peak times both in the morning and afternoon confirm the influence of phone habits (call times) in typical flow pattern changes.

However, those studies concluded that accurate vehicle flows cannot be obtained directly from mobile phone data due to the characteristics of this source data. The main question is how to correlate the number of crossing phones with the real number of crossing vehicles. Volume data on inter-cell boundaries provided by mobile phones does not yield information on the complete set of vehicles crossing a boundary, but only a statistical sample of all travelling vehicles. Some vehicles may carry either phones of other mobile phone operators or switched-off phones; these phones are not detected as crossing phones. These aspects imply a special treatment to correlate both measures (phone counts and vehicle counts), and this may require a calibration process to be performed using vehicle volume values obtained,

for example, from loop detectors located spatially on the same monitored section. In other pilot projects, different procedures have been used to relate both measures. For example, those based on the flow transfer function (Bekhor et al., 2008) or those based on data assimilation methods, using the so-called nudging method or Newtonian relaxation (Herrera and Bayen, 2008). According to these findings, more accurate estimates for the number of vehicles can be obtained by means of adequate treatments to relate both measures (phone counts and vehicle counts). A simpler scheme for determining and therefore predicting this relationship is through the use of mathematical models.

Mathematical models are widely used to solve specific problems in engineering or in other scientific fields, and they constitute a general framework for the formulation of statistical inference. In this paper models to predict volumes of vehicles from (anonymous) mobile phone call-data records are proposed in order to be used in transport applications. These models also require vehicle volume data (e.g. loop detectors) for the calibration stage.

### **3.2 Estimation of mobile phone count**

In areas close to the road network, the number of moving mobile phones will be directly related to the vehicles which carry them. Our aim is to infer vehicle volume data from information taken from call data of phones which move between cells. This is done using the handover process, which takes place when a phone with an active call enters the service area of another antenna (cell), reporting its new position (cell ID) and the time of entry. Thus, a record that an on-call phone entered a new cell is always made in the data base. An analysis of the moments at which these handovers are produced allows the detection of phones moving from one cell to another, or in other words, when they cross the boundary between cells.

However, using handover data exclusively has certain limitations. In some instances the call may be ended before changing cell, meaning that the handover does not occur, and the phone will not be detected. Other situations must be considered which increase the detected number of phones in motion. We know that when a call is made, the system always stores a record which contains call information such as the start time or the cell from which it originated. Therefore, it is also possible to detect the movement of a handset without a handover taking place. This can occur when two consecutive calls are made from a phone in different cells (A and B) within a short period of time, called T. In this case, we detect that a phone has been in two cells within a short time and so it is considered to have crossed the cell boundary. Consequently, this situation along with the execution of handovers allows a wider set of information to be available for monitoring movement between zones. These types of calls are labelled "in-motion calls" and are classified as:

1. Those in which a phone with an active call changes from cell A to cell B, that is, handover event (Figure 2a);
2. Those which occur when the same phone makes two consecutive calls ( $T < 15$  minutes) in cell A and B, respectively (Figure 2b).

A period of  $T=15$  minutes is used to make the detection of movement between cells A and B more precise, seeing as if  $T$  is a larger number, the phone may move to other cells before making the second call in cell B (Figure 3). Consequently, the analysis of call records stored in data bases, searching for in-motion calls made in the cells of interest, will permit the counting of phones which move from one cell to another as if there were a type of traffic counter station located at the boundary between them.

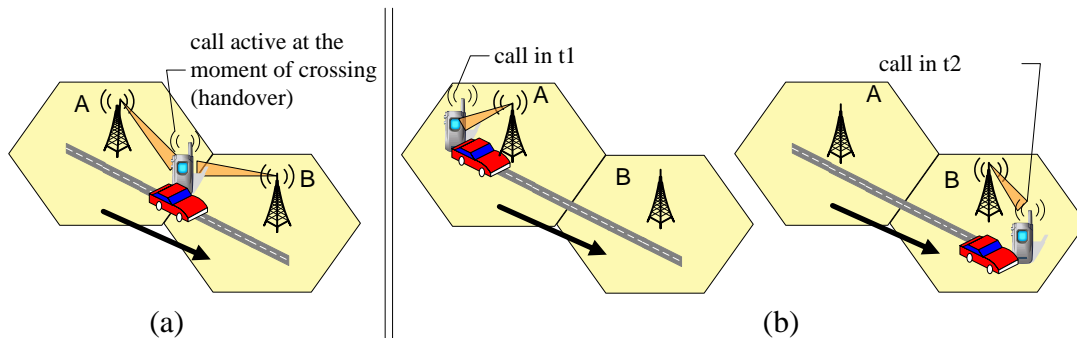


Figure 2 – In-motion calls: (a) handover and (b) two calls in  $t_1$  and  $t_2$  such that  $t_2 - t_1 \leq T$

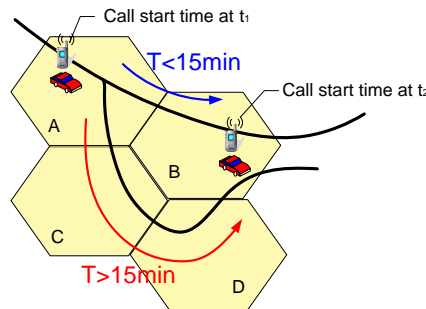


Figure 3 – Definition of  $T$ , time between start times of two consecutive calls

### 3.3 Monitored inter-cell boundaries

A monitored boundary between cells or “virtual traffic counter” (VTC) is composed of routes whose initial node and final node are in different and adjacent cells (Figure 4). In this way, every time one of the situations associated with in-motion calls occurs, it can be said that a telephone has moved through the VTC positioned at said boundary, or in other words, along one of the roads which connect them.

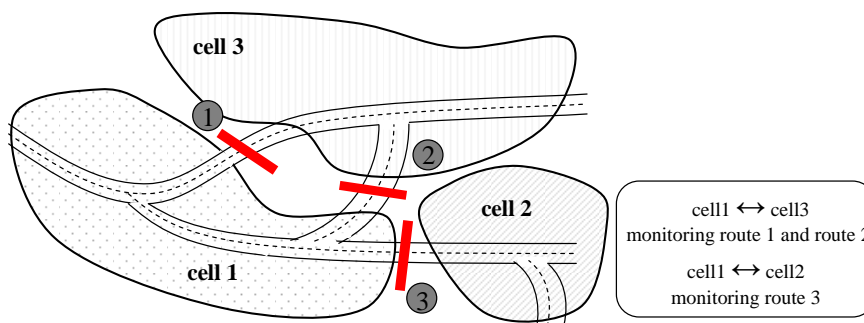


Figure 4 – Inter-cell boundaries monitored by VTCs.

By analysing the topology of a transport network model along with the cell distribution, it is possible to associate the boundaries between cells with sections of road (links). The size of these cells implies the existence of various routes (links) connecting each pair of cells, in such a way that the monitoring of phone movement will, in some cases, provide aggregate volumes of the groups of links which cross the boundary. In Figure 5, the mobility between cells 1 and 2 is associated with the route which passes through point 3 (C1→C2:  $\{\ell_5\}$ ; C1←C2:  $\{\ell_6\}$ ); while for the mobility between cells 1 and 3, given that there are two routes connecting them (points 1 and 2), the number of phones moving from one to the other will be the total which moves through both routes/links (C1→C3:  $\{\ell_1, \ell_3\}$ ; C1←C3:  $\{\ell_2, \ell_4\}$ ). It should be noted that not all links, and therefore boundaries between cells, are valid for monitoring using the detection of in-motion handsets. A valid observed boundary can only be made up of links whose initial node and final node is in distinct and adjacent cells, with no overlapping. Furthermore, there are additional criteria which characterise a boundary as valid: such as adequate coverage of the observed road, absence of uncertainties with regards to alternative routes or the absence of significant pedestrian traffic, among others.

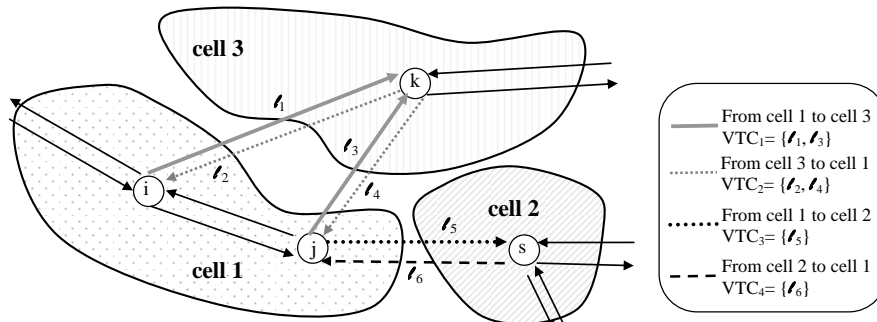


Figure 5 – Links associated to inter-cell boundaries which are valid for monitoring

### Cell coverage area

It is necessary to bear in mind that the cell boundaries are not static, but are to some extent dynamic; they depend on the number of mobile phones, direction of travel and network type (3G, 4G, etc.). Each cell coverage area has an effective radius, defined using signal strength measurements, within which a phone can communicate with a unique base station (Fig. 6a). This radius may vary according to factors as weather, call load, etc. However, this uncertainty in radius size does not affect the proposed approach due to the condition of the “observed” boundary. Only links whose initial and final nodes are located inside the effective cell coverage area are taken into account when creating a valid “monitored boundary” (Fig. 6b). Fluctuations in cell boundaries do not create an effect as long as those nodes remain inside the effective cell coverage area, as can be seen in Fig. 6c.

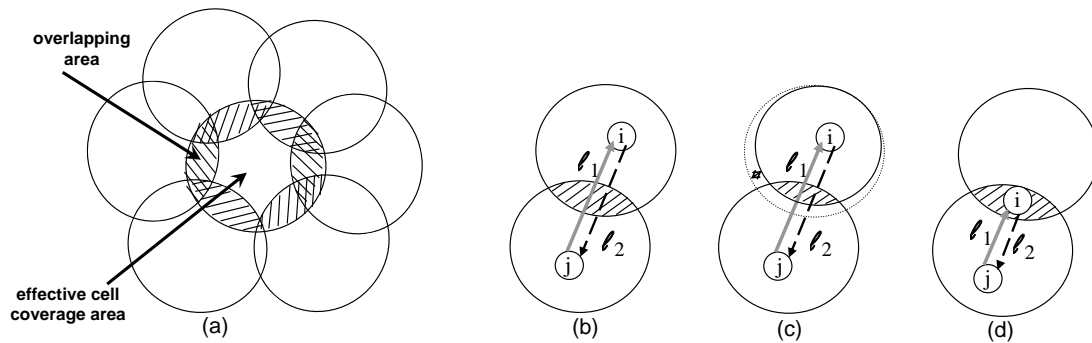


Figure 6 – Cell coverage area

When designing the cell layout for a mobile phone system, the cells overlap at the edges to prevent shadows in coverage. The cell overlap area is defined as the overlap between adjacent cells with regard to the primary coverage. The size of the overlapping area is a design parameter that depends on measurement control parameters. According to standard design criteria, a typical cell overlap area may occupy about 30–40% (EIA/TIA IS-95, 1993), although most network planners agree that the overlapping factor should be 20–30% (Yaipairoj and Harmantzis, 2008; UMTSWorld, 2009). Since the typical radius in non-urban environments is around 1–30 km, the overlapping area between two cells is greater than the area covered by a transport node (e.g. intersection). However, when a node (initial or final) is located in an overlapping area (Fig. 6d), the condition of ‘observed’ boundary is not matched; it belongs to two cells, one of them being the same as the other node (initial or final) in which it is located. Thus, this boundary, and the associated links, cannot be regarded as ‘valid for observation’.

### *Pedestrians*

It should be noted that any mobile phone making in-motion calls can act as a probe. The probes of interest are those located in vehicles moving along a motorway (i.e. passengers with a mobile phone or drivers using a hands-free device). The existence of a railway line running parallel to the motorway, may distort the measurements of moving phones counted along such roadway. A similar problem occurs with pedestrian users who are moving at a lower speed than vehicles in urban environments. Thus, it seems logical to discriminate between phones belonging to a vehicle user, a train passenger or a pedestrian by means of the determination of its speed. Different approaches based on advanced algorithms have been developed, implemented and tested in order to identify the moving phone's means of transport (Ramm and Schwieger, 2007; Schollmeyer and Wiltchko, 2007). According to the current state of the art, and by using exclusively phone data, a method for determining the speed requires the detection of the phone's location at two time reference points, i.e. two consecutive records (call, handover or location update record) for that same phone. Time records and cells allow the calculation of speed and the discrimination between pedestrians and vehicular users. But other problems arise when differentiating users by applying this approach, which is reliable in those cases where a single road crosses the inter-cell boundary and the followed route can be uniquely identified. This usually concerns main motorways/roadways where there are no other roads nearby to confuse identification. However, the identification of the route followed by the phone probe can be problematic



when multiple roadway links exist within the boundary zone between cells, especially in urban areas where there are multiple streets passing through a cell. Although this problem could be solved, the estimated speed does not achieve perfect accuracy in the case of urban environments, in which the speeds of phones belonging to pedestrians and phones belonging to drivers stuck in traffic are similar.

In our approach, the phone records are only used for reporting movement between cells and each record is directly linked with a phone passing through a zone. The identification of the means of transport (pedestrian and vehicular users) has been deliberately excluded by concentrating attention on cells whose boundaries have mainly vehicular traffic and where there are neither pedestrian nor train users.

### 3.4 Field test

The objective of this section is to estimate the volume of vehicles associated with the number of in-motion calls generated at a boundary between cells. To achieve this objective, a model has been proposed, for predictive means, which will provide this data in a similar way to the classic traffic counting stations. A data record was available concerning calls and vehicle flow detected over a series of days at points of the road network at which there is a counting station as well as a cell boundary.

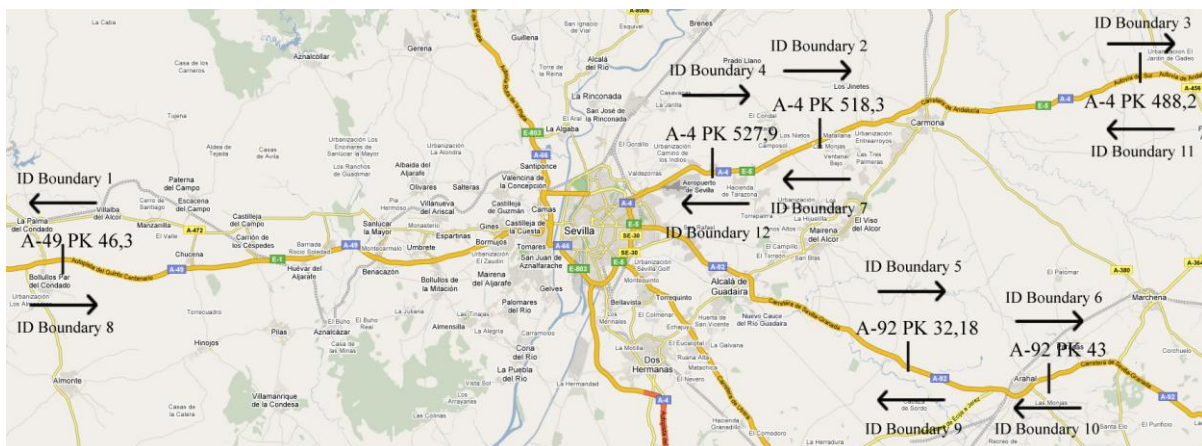


Figure 7 – Kilometre points in the proximity of which there are inter-cell boundaries

To identify the points to be observed, a field study was carried out on the roads within rings around the Seville (Spain) metropolitan area. Those selected were those which were associated with cell boundaries crossed by roads on which the traffic was essentially vehicles, away from urban environments and with a nearby counting station for contrast. Figure 7 shows these six points which make up 12 observed boundaries (one in each direction). For each point, the Traffic Management Centre (DGT) provided the hourly vehicle flow and Vodafone provided the call registers generated in the cells in question.

Given that the aim is to make the extracted information as representative as possible, days which show little variation to a “typical” working day were selected beforehand. For this reason the measurements were taken during days in the middle of the week: Tuesday,

Wednesday and Thursday, in which the preceding or following day was neither a bank holiday or special events (Easter Week, “Feria de Abril”, etc.) Finally, 18 days were selected to take observations on call data, Tuesday, Wednesday and Thursday over six weeks, as can be seen in Table I.

Table I– Observed days

MARCH 2008							APRIL 2008								
M	Tu	W	Th	F	Sa	Su	M	Tu	W	Th	F	Sa	Su	holidays	
					1	2		1	2	3	4	5	6		
3	4	5	6	7	8	9	7	8	9	10	11	12	13		
10	11	12	13	14	15	16	14	15	16	17	18	19	20	observed	
17	18	19	20	21	22	23	21	22	23	24	25	26	27		
24	25	26	27	28	29	30	28	29	30						
31															

### *Time period observed*

On each day, it was noted that the number of calls was significant only between the time period of 08:00-21:00. Outside of this time period the number was too small to make any inference concerning traffic flow. This makes sense given that user behaviour when making calls is conditioned by the hour of the day: either due to the working hours, the call fares or because it can be socially unacceptable to make calls at certain times. This can even be influenced by the cell, depending on the characteristics of the zone to which the service is being provided (if it is a residential area, a shopping centre, etc.). Proof of this is that phone activity during the time period between 01:00 and 06:00 is much reduced, with little probability of calls being made. On the contrary, from 18:00 the frequency of calls increases due to the start of reduced rate calls. Therefore, it is convenient to use dispersion measures to determine the time range in which the number of calls does not show much dispersion with regards to the mean. One of these measures is the coefficient of variation (CV), defined as the standard deviation divided by the mean. It is a dimensionless measure of the dispersion (or variation) of data points around the mean. The higher the CV, the higher the variability; and the lower the CV, the higher is the consistency of the data. Figure 8 shows the CV for the distribution of calls made in observed cells, from data provided by the monitored operator.

From a practical point of view, it is worth focusing the analysis on a time period during which the curves show low CV values. This occurs within the 08:00-21:00 time period, as can be seen in Figure 8 in which the CV value always remains below 0.5. Outside of this interval, the attained CV values are fairly larger than one, indicating that there is a greater variation in the making of calls with regards to the mean. Therefore, the observation time period will be between 08:00 and 21:00 and will be made up of 13 one-hour intervals, during which the making of calls shows, to the greater extent, a uniform behaviour for all cells. This means that the designed estimation models will be more stable for any period of time within the aforementioned range.

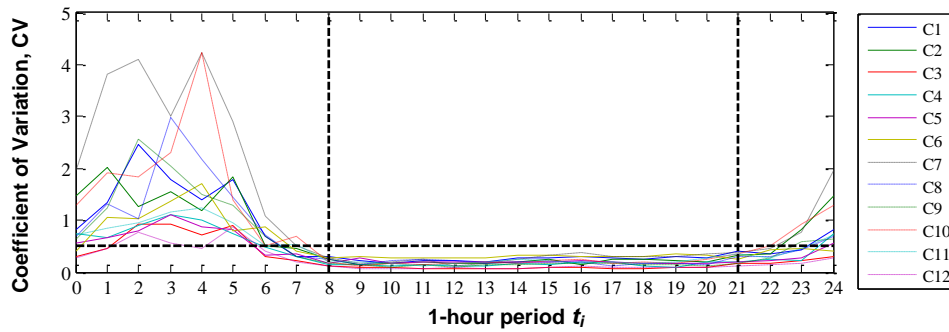


Figure 8 – Coefficient of Variation of call dataset in observed cells

## 4. MODEL DEFINITION

For the definition of models for the estimation of vehicle traffic flow, using information from calls generated by a sample composed of certain mobile phones, a dataset from a bidimensional sample which links both variables  $\{(x_1, y_1), \dots, (x_n, y_n)\}$  is available.  $x$  and  $y$  are, respectively, observed values of in-motion calls and vehicles which have crossed any of the observed boundaries between cells. Therefore, for each of the observed days there is a sample formed of 13 pairs of values: the number of calls and the number of vehicles registered each hour from 08:00 to 21:00. With the objective of having a sufficient sample size to fit a model, the data from all the observed days was aggregated into one single sample. In the end, the size of the resulting sample  $\{(x_1, y_1), \dots, (x_n, y_n)\}$  was  $n=234$ : (18×13) records for each of the twelve km-points studied.

The whole dataset will not be used for the process of parameter estimation, given that part of the dataset, called validation set, will be reserved for a latter use in the evaluation of the estimations and will not participate in the fitting. Generally, a third or quarter of the data is reserved to validate the results. The rest of the data record will be used in the fittings and it is called the calibration set. In the notation of the models which will be formulated below, the following will be used:  $i$  to designate the 1-hour time period,  $d$  for the observed day and  $k$  for the kilometre point on the road or boundary observed.

### 4.1 Formulation of the Physical Model

An estimation model based on the physical phenomenon associated with in-motion calls (Figure 2) will be proposed. To define it, it must be considered that a mobile phone makes an in-motion call when either the user makes a call in each of the two cells which form the boundary within a short period of time (event M, Fig 2b); or the user has an active call and changes from one cell to another (event L, Fig 2a). With regard to the event M, although it requires the making of two calls, only one is counted as in-motion call. The concept of mobility is only associated with the phone from which the calls needed for the detection of movement were made.

In view of these situations, it seems logical that the observations of in-motion calls are affected by affected by the tendency to make calls. A typical user tends to make calls during

a certain period of time, but does not call in other periods, such as late night or early morning. Furthermore, in the case of handovers, the duration of the calls will also affect these observations given that, the handover is more likely to occur as call duration increases when the phone is moving. Therefore, the time dependence in user behaviour when making calls must be taken into account when formulating a model which allows the estimation of the number of vehicles using a given number of in-motion calls.

### *User calling behaviour*

Figure 9 shows the call intensity, represented in terms of number of calls made, in a cell during a typical working day. It can be seen that, essentially, there are three periods of activity during a typical working day, one of very reduced activity and two peaks pronounced during rush hours whose range varies from one cell to another. The first period corresponds to the nighttimes (00:00-07:00), during which the making of calls is improbable. The second period lasts for the most part of a working day (08:00-16:00). The last time period corresponds to the evening and end of the day (16:00-00:00). In turn, it can be seen that call intensity in a working day increases sharply in the morning between 07:00 and 10:00, whereas the transition into nighttimes shows a more gradual decrease. Call intensity can vary from one day to another, week to week or even according to the time of year (summer, Easter Week, etc.). In addition, call activity is also affected by the characteristics of the zone covered by the cell (residential, industrial or commercial areas). However, call intensity patterns normally tend to show a similar behaviour in different zones, varying ratios.

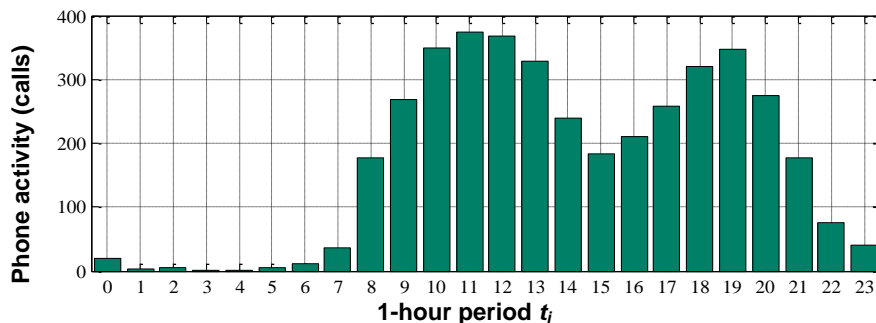


Figure 9 – Number of calls per hour during a day in a typical cell

Similar behaviour is shown by call duration. In Figure 10 the calls made are broken down according to their duration, showing how the trend varies throughout the day. By analysing the mean duration, it can be seen more clearly that users tend to make longer calls after 19:00. The nature of the phone service causes this variation in duration, given that it could be influenced by the call rates or even user availability.

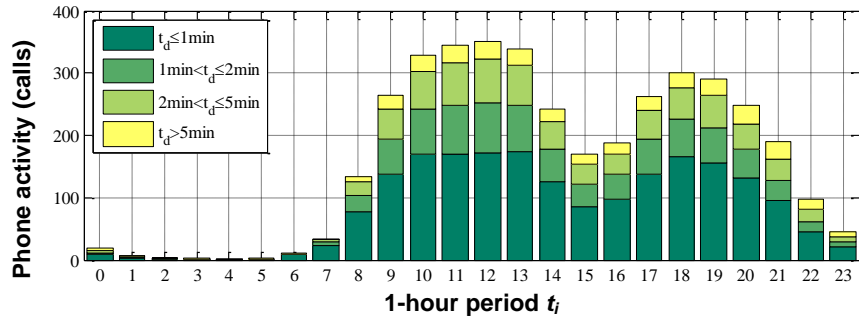


Figure 10 – Number of calls per hour during a day in a typical cell depending on the call length  $t_d$

### Model definition

In this section it is proposed a theoretical model which represents the relation between the number of in-motion calls and the number of vehicles. Its expression considers each of the aforementioned situations of in-motion calls. It makes use of call characteristic such as the probability of making a call or its duration, according to the time period  $t_i$  and the observed boundary  $k$ . The nature of the cells under study, i.e. those away from urban environments, allows the assumption that practically all of the observed in-motion calls during the time period  $t_i$  and at the boundary  $k$ ,  $n_{movcall}(t_i, k)$ , will be made by in-motion users associated with vehicular traffic on the existing roads at that boundary. These users therefore constitute a sample of the vehicle population which has crossed the boundary  $k$  during each interval  $t_i$ ,  $n_{veh}(t_i, k)$ .

With regards to the first type of in-motion calls (event M), it must be introduced a term which reflects the importance of the probability that a phone which uses the monitored operator and which is on board a vehicle, makes two calls during the time period  $t_i$ . The event “making two calls on board a vehicle, using the monitored operator” is a superposition of two independent events  $A \equiv$  ‘making a call on board a vehicle, using the monitored operator’. Whereas said event A has a probability of occurrence  $P_A \equiv P_{vehcall}(t_i)$ , then the probability of the occurrence of event M can be estimated using  $P_{vehcall}(t_i) \cdot P_{vehcall}(t_i)$ . Assuming that the population which generates the in-motion call is the set of vehicles crossing a boundary,  $n_{veh}(t_i, k)$ , the number of in-motion calls generated by event M will be obtained by using:

$$n_M(t_i, k) = n_{veh}(t_i, k) \cdot P_M = n_{veh}(t_i, k) \cdot P_{vehcall}(t_i) \cdot P_{vehcall}(t_i) = n_{veh}(t_i, k) \cdot P_{vehcall}^2(t_i) \quad (1)$$

For the second type of in-motion calls (event L), the proposal is similar but it introduces a term related to the probability that the call performs a handover. Thus, the probability of occurrence of event L requires two terms: on one hand, the aforementioned probability  $A \equiv$  ‘making a call on board a vehicle, using the monitored operator’ and on the other, the probability that said call requires a handover,  $P_h$ . Regarding the latter, the probability that a handover is performed on a call comes from the fact that its duration,  $t_d$ , exceeds the time of remaining in the cell,  $t_p$ . As previously mentioned, assuming that  $n_{veh}(t_i, k)$  is the population which generates the in-motion calls at boundary  $k$  during time period  $t_i$ , the number of in-motion calls generated under conditions associated with event L will be defined as:

$$n_L(t_i, k) = n_{veh}(t_i, k) \cdot P_L = n_{veh}(t_i, k) \cdot P_{vehcall}(t_i) \cdot P_h(t_i, k) \quad (2)$$

The total number of observed in-motion calls in the time period  $t_i$ , corresponds to the sum of both sets of calls,  $\{n_M$  and  $n_L\}$ , and an expression (3) is obtained in which the observed value of the in-motion calls corresponds to the mean number of vehicles making an in-motion call.

$$\begin{aligned} n_{movcall}(t_i, k) &= n_{veh}(t_i, k) \cdot P_{vehcall}^2(t_i) + n_{veh}(t_i, k) \cdot P_{vehcall}(t_i) \cdot P_h(t_i, k) \\ &= n_{veh}(t_i, k) \cdot [P_{vehcall}^2(t_i) + P_{vehcall}(t_i) \cdot P_h(t_i, k)] \end{aligned} \quad (3)$$

According to this expression, the probability that a phone makes an in-motion call, labelled  $P(t_i, k)$ , can be estimated using  $P(t_i, k) = P_{vehcall}^2(t_i) + P_{vehcall}(t_i) \cdot P_h(t_i, k)$ . Therefore, the number of vehicles which cross a boundary  $k$  during time period  $t_i$  will be given by:

$$n_{movcall}(t_i, k) = P(t_i, k) \cdot n_{veh}(t_i, k) \Rightarrow n_{veh}(t_i, k) = \frac{n_{movcall}(t_i, k)}{P(t_i, k)} = \frac{n_{movcall}(t_i, k)}{P_{vehcall}^2(t_i) + P_{vehcall}(t_i) \cdot P_h(t_i, k)} \quad (4)$$

This expression provides a value for the volume of vehicles which cross the boundary between cells in terms of the detected in-motion calls and other variables related to the calls. The number of in-motion calls associated with the crossing of a boundary during the given time period,  $n_{movcall}(t_i, k)$ , is obtained by analysing the records, provided by the network operator, of calls made in the cells involved in the boundary  $k$ . For the other terms, i.e. the probability that a vehicle makes a call and the probability of a handover,  $P_{vehcall}$  and  $P_h$  respectively, an additional statistical treatment is required to take into account the dependence relation of time on the call characteristics. The following points explain in detail the method developed to obtain the aforementioned terms.

#### *Probability of making a call on board of a vehicle, $P_{vehcall}$*

It is usual for a set of data to follow some type of statistical distribution. Adequately defining the probability of occurrence of any event can be tedious if this distribution is not known. The depiction of the values can suggest the use of one of the many known distributions (Poisson, Uniform, Normal, etc.) and thus determine the best one which fits the data. However, there are some ways of approaching such a term of probability without the need for prior knowledge of the distribution: such is the case of an empirical procedure in which the assignation of the probabilities of the events of interest is based on the observed information. In these instances, a large number of tests is needed to infer the probability of occurrence of a event which, in this case, is the making of a call on board a vehicle, using the monitored network operator. Using this empirical focus, the probability is determined based on the proportion of times in which a favourable or successful event occurs with regards to the total number of possible results.

As previously mentioned, the nature of the studied cells, i.e. those away from urban environments and residential areas, allows the assumption that practically all of the observed in-motion calls at each of the boundaries will be made by users in the flow of vehicular traffic. Therefore, the probability of “making a call on board a vehicle, using the monitored operator”,  $P_{vehcall}(t_i)$ , will be equal to the number of calls made on board a vehicle, using said operator, during the time period  $t_i$ , divided by the total number of vehicles which crossed that point during said time period. The total number of vehicles which cross the boundaries is known

due to the counting stations which are near each boundary. On the other hand, the term, “calls made from vehicles, using the monitored operator”, will be directly related to the observed in-motion calls. Bearing in mind that an in-motion call originating from event M implies the making of two calls from the handset which is travelling on board a vehicle, the probability of “making a call on board a vehicle, using the monitored operator” during the time period  $t_i$  is defined as:

$$P_{vehcall}(t_i) = \frac{\sum_{d=1}^D \sum_{k=1}^K [n_L(t_{i,d}, k) + 2 \cdot n_M(t_{i,d}, k)]}{\sum_{d=1}^D \sum_{k=1}^K n_{veh}(t_{i,d}, k)} \quad (5)$$

with  $K=12$  being the total number of studied boundaries and with  $D=18$  being the total number of observed days. As can be seen, a dependence relation has been established in accordance with the considered time period  $t_i$ . The strong dependence relation of time in user behaviour when making calls also influences this probability. Figure 11 represents the variation over time of probability  $P_{vehcall}(t_i)$  empirically obtained. This probability is valid for any of the monitored boundaries seeing as the event “making a call on board a vehicle, using the monitored operator” is independent of the cell in which it occurs, except perhaps in zones where driving is difficult for making a call.

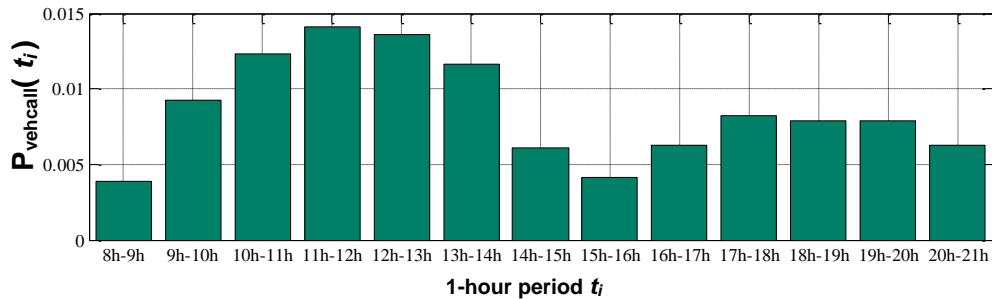


Figure 11 – Probability of making call on board a vehicle

Figure 11 reflects that the probability of an in-motion user (on board a vehicle) making a call follows the trends of daily call activity shown in Figure 9 (for any user) for the most of time periods. There are also two pronounced peaks during rush hours. The morning peak centres around the same time period in both cases; while the evening peak for calls made by in-motion users is less pronounced and wider. This lower tendency may originate from the fact that a in-motion user only makes a call outside of working hours when being necessary to do so. While the tendency to make calls seems to be independent of user mobility during working hours.

### *Probability of handover*

An expression for the probability of handover is proposed by Sallent *et al.* (2003). This work assumes a simplified scenario for handover where a mobile system covering a road network is divided into regular cells. The distance which a mobile phone travels within a cell before crossing the boundary of said cell is modelled as a random variable with uniform distribution throughout the interval  $[0, L]$  metres. The permanence time of a mobile phone in the cell,  $t_p$ , if the phone moves at a constant speed of  $V$  m/s will also be a random variable with uniform

distribution throughout the interval  $t_p \in [0, L/V]$  seconds. Likewise, the call duration,  $t_d$ , can be modelled as an exponential random variable of mean  $T_c$  seconds. The corresponding probability density functions of  $t_p$  and  $t_d$  will be given by:

$$f(t_p) = \frac{1}{L/V} \quad t_p \in \left[0, \frac{L}{V}\right], \quad f(t_d) = \frac{1}{T_c} e^{-\frac{t_d}{T_c}} \quad t_d \in [0, +\infty) \quad (6)$$

A call requires the execution of a handover procedure when its duration exceeds the permanence time of a phone in the cell. Thus, the probability of a handover can be calculated using a conditional probability:

$$P_h = \int_0^{\frac{L}{V}} P(t_d > t_p | t_p) f(t_p) dt_p, \quad \text{where } P(t_d > t_p | t_p) = \int_{t_p}^{\infty} f(t_d) dt_d = e^{-\frac{t_p}{T_c}} \quad (7)$$

By defining a dimensionless factor as  $\alpha = \frac{L}{V \cdot T_c}$  the final result is:

$$P_h = \int_0^{\frac{L}{V}} e^{-\frac{t_p}{T_c}} f(t_p) dt_p = \int_0^{\frac{L}{V}} e^{-\frac{t_p}{T_c}} \frac{1}{L/V} dt_p = \frac{1}{\alpha} [1 - e^{-\alpha}] \quad \text{where } \alpha = \frac{L}{V \cdot T_c} \quad (8)$$

Parameter  $\alpha$  depends on the length  $L$ , the speed  $V$  and the mean call duration  $T_c$ . The parameter  $L$  represents the distance that a mobile phone must travel within a cell until it enters another cell (cross the boundary). This length depends on the roads within the origin cell of the boundary. Something similar occurs with the speed: the value depends on the type of road in question (motorway, main road, etc.). Although speed can vary along the same road according to the time period, due to the level of saturation (traffic) for example, it will be considered uniform within the origin cell of the boundary. Therefore, the value of  $\alpha$  will very much depend on the observed boundary  $k$  in terms of the length and speed associated with the type of road which runs through the origin cell of the boundary. Similarly, bearing in mind the call duration's dependence relation on time,  $\alpha$  will also depend on the analysed time period. Therefore, the probability of handover will be based on the time period and the boundary,  $P_h(t_i, k)$ .

To calculate length  $L$  there are basic functionalities in the network which allow us to know the cell identifier (antenna) to which a phone is connected. In this way, with the help of a GPS, it is possible to establish a correspondence between the position on the road and the cell which provides it with service. Based on this, not only can the length of the road along which a mobile phone travels within a cell be found, but so can the mean travelling speed over said road at any time of day, without the network operator needing to reveal confidential coverage information. Using this simple application, implemented from mobile devices, the roads can be measured various times using the floating car technique, with the aim of being able to characterise the origin cell of each boundary  $k$  in terms of length  $L$  and speed  $V$ .

Finally, it is possible to empirically obtain the duration  $T_c$ , through the analysis of call data provided by the operator, determining an mean value of call duration according to the considered time period,  $T_c(t_i)$ . This duration, which depends on time period  $t_i$ , and the values



of speed and length which correspond to the road which passes through the origin cell of the boundary, allow the definition of handover probability as a function which is variable over time and the boundary,  $P_h(t_i, k)$ . Figure 12 shows the temporal distribution of the probability of handover calculated using (8) for the characteristics of  $L$  and  $V$  at a specified boundary.

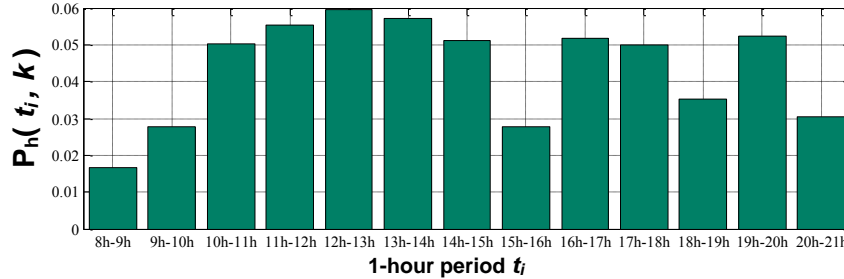


Figure 12 – Probability of handover at one example of boundary  $k$  in each 1-hour period  $t_i$

By substituting the previous expression of  $P_h(t_i, k)$ , the equation (4) would be expressed as:

$$n_{veh}(t_i, k) = \frac{n_{movcall}(t_i, k)}{P_{vehcall}^2(t_i) + P_{vehcall}(t_i) \cdot \frac{1}{\alpha(t_i, k)} [1 - e^{-\alpha(t_i, k)}]}, \quad \text{where } \alpha(t_i, k) = \frac{L_k}{V_k \cdot T_c(t_i)} \quad (9)$$

As explained previously, this expression is obtained by assuming some events as negligible (e.g. the making of a call in each cell while also executing a handover). But the existence of undetected in-motion situations, such as call which ends before the vehicle leaves the cell (no handover) without making a second call in the new cell, produce errors which must be taken into consideration in the definition of the model. Therefore a series of parameters are introduced in the expression (9) to correct the indeterminacies introduced by the hypotheses required to generate the model, which is formulated as:

$$\hat{y}(t_i, k) = \frac{a \cdot x(t_i, k)}{P_{vehcall}^2(t_i) + P_{vehcall}(t_i) \cdot \frac{b_1}{\alpha(t_i, k)} [1 - e^{-b_2 \cdot \alpha(t_i, k)}]} + d \quad (10)$$

In this model the variables  $x$  and  $\hat{y}$  correspond, respectively, to the number of in-motion calls made between the cells of boundary  $k$  during interval  $t_i$ ,  $n_{movcall}(t_i, k)$ , and the estimated number of vehicles.

### *Estimation of model parameters*

To estimate the parameters of the model,  $\Phi = \{a, b_1, b_2, c, d\}$ , there are different adjustment procedures which will use the part of the data record which concerns both variables, the so-called calibration set. In this case, a parameter adjustment technique will be used, in an attempt to satisfy an adjustment criterion. The criterion used takes into account relative error, given that an estimation which is over by 50 is not as important on a real value of 1000 as it is on a real value of 100, with the latter example being more serious. Therefore, given that in this case it is just as important to not make a mistake by going over or falling short, the absolute value must be taken into consideration in the criterion. Consequently, the

adjustment will seek to minimise the absolute value of the relative error between observed and estimated values by resolving the following optimisation problem:

$$\text{Min}_{\phi} \sum_{j=1}^N \left| \frac{y_j - \hat{y}_j}{y_j} \right|, \text{ where } y_j \text{ are observed flows and } \hat{y}_j \text{ estimated flows.}$$

#### 4.4 Evaluation

To evaluate the model, the testing set was utilised, making use of criteria such as the classic direct error measures, evaluation of cumulative distribution function of the absolute relative error using percentiles, rank (or Spearman's) correlation along with the coefficient of linear (or Pearson) correlation. Using these, the estimated values were compared to the observed values. Lastly, some additional models have been proposed to contrast the physical one.

##### *Additional models*

Figure 13 shows the variation over time in the intensity of the number of in-motion calls made and the number of vehicles which cross a determined boundary on any given day. The time periods with high intensity are close in time in both cases, but do not generally coincide. While the peak in vehicle numbers appears at about 07:00 (the start of the working day), at this time the number of calls is low, as it is unusual to make calls. The peak in call numbers is reached at about 10:00. This is due to the fact that the number of in-motion calls is strongly affected by the conduct of the users who make them and this varies with time in accordance with criteria such as reduced call rates or the acceptable times for making calls. Additionally, the flow of vehicles fluctuates with time in accordance with different criteria such as rush hour at the start of the working day, breaks, lunch, etc.

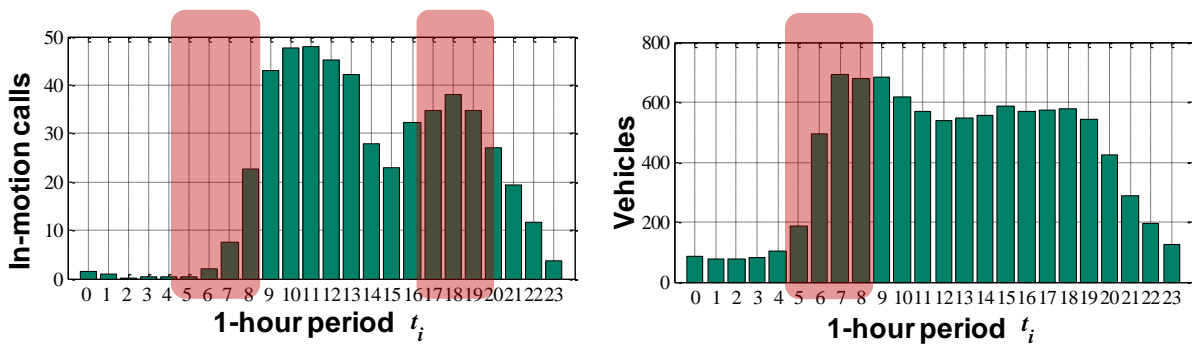


Figure 13 – Specific one-day sample of numbers of in-motion calls and vehicles counted in each 1-hour period  $t_i$

This different time variation in the intensity of calls and vehicles motivates the use of the factors or time coefficients defined in expression (11), which make the model dimensionless. These coefficients weight observations in each time period according to their relative importance within the considered time periods.

$$f_v(t_i) = \frac{\bar{n}_{veh}(t_i)}{\bar{n}_{veh}}, \quad g_{ll}(t_i) = \frac{P_{vehcall}(t_i)}{\bar{P}_{vehcall}} \quad (11)$$

being

$\bar{n}_{veh}(t_i)$  is the average no. of vehicles crossing all the  $K$  boundaries during interval  $t_i$  in each of the  $D$  observed days, so that  $\bar{n}_{veh}(t_i) = \frac{1}{D \cdot K} \cdot \sum_{d=1}^D \sum_{k=1}^K n_{veh}(t_i, d, k)$ ,  $K=12$  and  $D=18$ .

$\bar{n}_{veh}$  is the average no. of crossing vehicles during all the  $H$  hourly intervals from 08:00 to 21:00, so that  $\bar{n}_{veh} = \frac{1}{H} \cdot \sum_{i=1}^H \bar{n}_{veh}(t_i)$ ,  $H=13$ .

$P_{vehcall}(t_i)$  is the probability of making a call on board a vehicle using the monitored operator in the time period  $t_i$  defined in (5).

$\bar{P}_{vehcall}$  is the previously averaged probability in all intervals,  $\bar{P}_{vehcall} = \frac{1}{H} \cdot \sum_{i=1}^H P_{vehcall}(t_i)$ .

The factor  $f_v(t_i)$ , associated with the time variation in vehicle flow, is defined as the relation between the mean number of vehicles counted in a 1-hour period  $t_i$  and the mean number of vehicles counted during the entire observed time period (8:00–21:00) at all points. The calls factor,  $g_{II}(t_i)$ , is defined similarly as the relation between the probability that a vehicle makes a call using the monitored operator in the 1-hour time period  $t_i$ ,  $P_{vehcall}(t_i)$ , and the mean probability for the entire observed time period. These probabilities are calculated empirically from the proportions of occurrences, when the number of in-motion calls and the number of moving vehicles is known. Once the factors of vehicle and call intensity have been obtained, their numerical values will be used as coefficients in the functional form of the models, whose hourly values are shown in Figure 14.

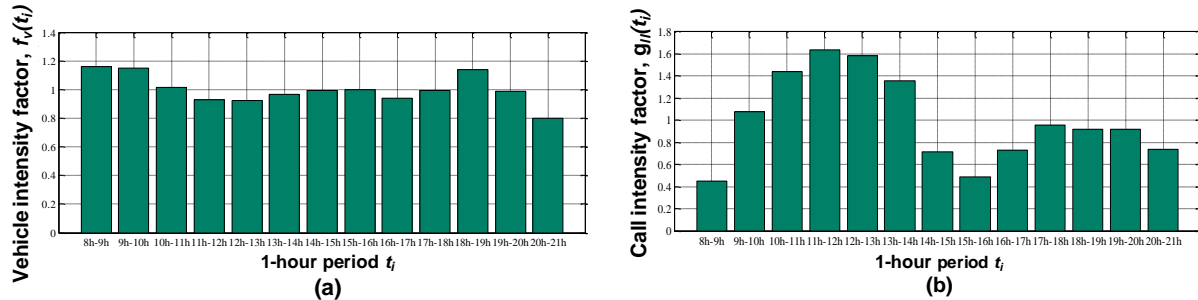


Figure 14 – Fixed values calculated for the factors  $f_v(t_i)$  and  $g_{II}(t_i)$ , (a) and (b) respectively, in each 1-hour period

After an analysis of the most characteristic families, models combining the in-motion calls,  $n_{movcall}(t_i, k)$ , and the hourly coefficients,  $f_v(t_i)$ ,  $g_{II}(t_i)$  and  $P_{vehcall}(t_i)$ , will be proposed using exponential and polynomial functions. In all of these, the dependent variable is  $n_{veh}(t_i, k)$ , which represents the flow of vehicles in time period  $t_i$  at boundary  $k$ . Other models which, along with the physical model, allow comparisons to be made are presented below:

- M1.** Cobb-Douglas Model:  $\hat{y}_1(t_i) = a \cdot f_v^\alpha(t_i) \cdot g_{II}^\beta(t_i)$
- M2.** Modulated Cobb-Douglas Model:  $\hat{y}_2(t_i) = (a + b \cdot x(t_i)) \cdot f_v^\alpha(t_i) \cdot g_{II}^\beta(t_i)$
- M3.** Second order Modulated Cobb-Douglas Model:  $\hat{y}_3(t_i) = (a + b \cdot x(t_i) + c \cdot x^2(t_i)) \cdot f_v^\alpha(t_i) \cdot g_{II}^\beta(t_i)$
- M4.** Linear Model:  $\hat{y}_4(t_i) = a + b \cdot x(t_i)$
- M5.** Quadratic Model:  $\hat{y}_5(t_i) = a + b \cdot x(t_i) + c \cdot x^2(t_i)$
- M6.** Physical Model, defined in (10).

The parameters of these models are estimated from those of the physical model. More details about the formulation of the models can be found in Caceres (2010).

### *Comparative analysis*

Table II shows the error measures of each of the models. The mean error and the mean relative error, both in terms of absolute value, allow the comparison of the estimations with the real values in the way of classical error analysis (MAE, mean absolute error; MARE, mean absolute relative error). Absolute values are used in order to that the prediction error has the same significance either upwards or downwards. Based on these measures of MAE and MARE, the best models are 6, 2 and 3, seeing as the ranges are of the same order in all three. With regards to the percentiles, these are used to indicate that the values of a variable below which a certain percent of observations fall. The best models are those which show the smallest values of the variable for each one of the percentiles. The analysis of the cumulative distribution function of MARE using percentiles for each model shows that models 6, 3 and 2 are better than the others. The 60th percentile shows that 60% of the volume estimates using these models achieve a MARE lower than 0.21.

Table II – Error Measures for each model

	<b>MODEL 1</b>	<b>MODEL 2</b>	<b>MODEL 3</b>	<b>MODEL 4</b>	<b>MODEL 5</b>	<b>MODEL 6</b>
<b>MAE</b>	237.2021	210.6609	210.0152	223.9775	223.2356	203.5965
<b>MARE</b>	0.2318	0.2039	0.2032	0.2169	0.2157	0.2000
<b>25% Perc.</b>	0.0799	0.0769	0.0763	0.0777	0.0799	0.0772
<b>30% Perc.</b>	0.0982	0.0933	0.0933	0.096	0.0967	0.0938
<b>50% Perc.</b>	0.2039	0.1694	0.1694	0.1763	0.1783	0.1646
<b>75% Perc.</b>	0.3596	0.3076	0.3034	0.3351	0.3311	0.2986
<b>85% Perc.</b>	0.4298	0.379	0.3791	0.4117	0.4113	0.3703
<b>100% Perc.</b>	0.9765	0.83	0.8492	1.0723	1.0103	0.9985
<b>Rank corr</b>	0.25102	0.51398	0.51375	0.39904	0.39904	0.56623
<b>Linear corr</b>	0.28436	0.4698	0.48134	0.33953	0.34893	0.53333

Another criterion evaluated is the correlation between the estimates and the real volume values. For this purpose, the linear correlation coefficient, or the Pearson coefficient, and Spearman rank correlation coefficient were studied. The Pearson coefficient is the most widely used measure of linear relationship between two variables. The Spearman rank correlation coefficient is a measure of the strength of the association or concordance between two variables using the relationship between ranks (e.g. a positive Spearman coefficient corresponds to an increasing monotonic trend between variables). In terms of rank correlation, models 6, 2 and 3 clearly stand out from the others; especially model 6 which achieved high values for the size of the considered sample. Similar ranking is achieved in terms of the linear correlation coefficient.

Figure 15 shows the estimations of vehicle flow for each of the six models along with the mean values of observed vehicles, in each 1-hour period and at each boundary. This figure

allows a visual comparison to be made of the precision achieved in each model as well as an appreciation of the order of magnitude of error estimation.

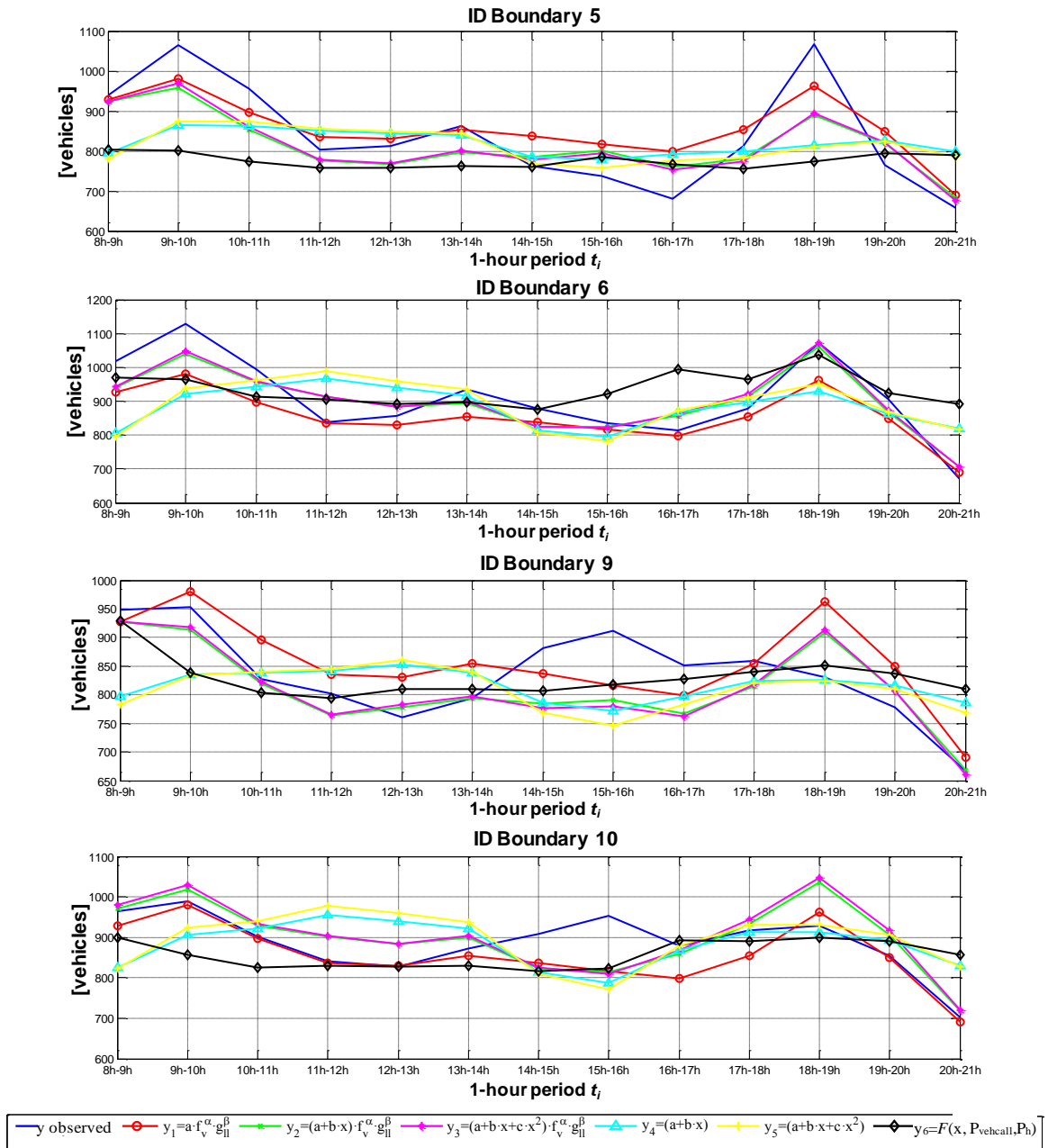


Figure 15 – Vehicle flows crossing boundary 5, boundary 6, boundary 9 and boundary 10 in each 1-hour period observed, and as estimated by the proposed models

In previous figure, it can be seen that model 1 is the one which provides the worst results. This is due to it depending exclusively on the factors related to time variability in the intensity of calls and vehicles. No information about the number of in-motion calls appears in the functional form. The model therefore loses information regarding a proportion of vehicles crossing a boundary. Something similar occurs for the models 4 and 5 but losing significant information on time variability due to the absence of intensity factors. On the contrary, models 2, 3 and 6 show a good fit to observed volumes, within an admissible range of error (<20%). The models 2 and 3 show very similar estimates of vehicle flow for a known number

of in-motion calls. This comes from the functional form of the models, given that both are different only in the quadratic term which multiplies the number of in-motion calls ( $c \cdot x^2$ ). The order of magnitude of this parameter  $c$  does not introduce any significant change in estimated values, whereas the rest of the parameters remain within the same order in both models. Consequently, the differences between the estimates of both models are barely appreciable graphically, but this does not occur using criteria such as MAE, MARE or correlation, among others, as has been previously demonstrated. With regard to these criteria, it can be seen that the estimations of models 2, 3 and 6 are those which best follow the peaks and valleys of the observed values of vehicle flow. Thus, it seems reasonable to think that these three models – models 2, 3 and 6 - are the most suitable for the estimation of vehicle flow using the number of in-motion calls made over a period of time.

In conclusion, the models which achieved the best results among all the criteria were models 6, 2 and 3. By observing their results with regards to the comparison of the studied criteria, it is difficult to establish a clear order of priority given that:

- i. although model 6 presents a slight advantage over models 2 and 3 in terms of MAE, MARE and percentiles of MARE, the order of magnitude is very similar in all three models;
- ii. graphically, the predictions of the three models reflect a similar behaviour when accurately matching the peaks and valleys in the curves of the real flow;
- iii. model 6 shows a certain improvement over models 3 and 2 in terms of rank correlation and linear correlation.

In general, it can be said that physical model (M6) presents a slightly superior predictive capacity when considering spatial dependence relation, based on the geometrical characteristics which should be specifically measured at each monitored point of the road. However said improvement may not be sufficient to justify the model's loss of flexibility with regards to models 2 and 3 which are independent of the boundary. In view of these conclusions, models 2, 3 and 6 were finally considered to be equally viable for the estimation of vehicle flow, and they made clear the following:

1. it is necessary to complement in-motion call data with information about the hourly variability of the behaviour of the users who travel using mechanised means when making calls, or other characteristics associated with the circulation of vehicles and;
2. the applicability of the methodology within admissible error ranges compared to measurements provided by traditional counting stations.

## **5. CONCLUDING REMARKS**

The objective is to present mobile telephony as a real, plausible and competitive alternative for the gathering of traffic flow information. To that end, various models, which are capable of providing said information in a similar way to the classic counting stations, have been

proposed, taking as an input variable the information from calls made by certain users who move through the studied region. The experimental testing of the proposed models has shown evidence that mobile telephony, via in-motion calls, constitutes a valid alternative for the determination of traffic flow in comparison with traditional techniques based on counting stations. The MARE levels are around 20%, an acceptable figure given that the main input variable (calls) has a high level of randomisation. These calls do not represent the whole set of mobile phones crossing an inter-cell boundary but only those pertaining to the monitored operator, and they are also calls made under certain conditions. This affects to a large extent the representative nature of the detected in-motion phones and therefore, the results. Despite these difficulties, the estimation values of the models achieved reasonable results in comparison with the values observed by the counting stations. Therefore it is clear that mobile telephony constitutes a valid alternative for practical use in transport studies and it does not require the installation of additional elements on the roads.

An immediate application of the flows inferred from mobile telephony is within the field of adjustment methodologies of mobility matrices which are traditionally made using the network's counting stations, although a series of modifications, which take the origin of the utilised flows into consideration, would have to be made.

## **ACKNOWLEDGEMENTS**

This project was financed by the Spanish Ministry of Science through R&D National Programmes (ENE2008-05552), and by Vodafone Spain through the Minerva Project (1C-021). The authors acknowledge the collaboration of A.L. Jimenez from the Southwestern Spain Traffic Management Centre (DGT).

## **REFERENCES**

- BEKHOR, S., HIRSH, M., NIMRE, S., FELDMAN, I. (2008) Identifying Spatial and Temporal Congestion Characteristics using Passive Cellular Phone Data. *Transportation Research Board 87th Annual Meeting*, paper no. 1534.
- CACERES, N., WIDEBERG, J., BENITEZ, F.G. (2008). Review of traffic data estimations extracted from cellular networks, *IET Intell. Transp. Syst.*, vol. 2(3), pp. 179-192.
- CACERES, N. (2010). Mobility matrix estimate by using cellular phone data (In Spanish). Ph.D. Thesis Dissertation, University of Seville (Spain). (Accessed on February, 2010). <http://www.esi2.us.es/GT/docs/TesisNCS.pdf>
- CIA World Factbook. (2009). List of countries by number of cellular phones in use. (Accessed October 2009). <https://www.cia.gov/library/publications/the-world-factbook/index.html>
- EIA/TIA IS-95, 1993. Mobile Station-Base Station Compatibility Standard for Dual-Mode Wideband Spread Spectrum Cellular System.
- FONTAINE, M.D., SMITH, B.L. (2005). Probe-based traffic monitoring system using wireless location technology: investigation of the relationship between system design and effectiveness. *Transportation Research Record: Journal of the TRB*, No. 1925.

- HERRERA, J., BAYEN, A. (2008). Traffic flow reconstruction using cellular sensors and loop detector data. *Transportation Research Board 87th Annual Meeting*, paper no. 1868.
- HÖPFNER, M., LEMMER, K., EHRENFORDT, I. (2007). Cellular Data For Traffic Management - First Results of a Field Test. *ITS Europe Conference*, paper no. 2407.
- QIU, Z., JIN, J., CHENG, P., RAN, B. (2007). State of the Art and Practice: Cellular Probe Technology Applied in ATIS. *Transportation Research Board 86th Annual Meeting*, no. 0223.
- RAMM, K., SCHWIEGER, V. (2007) Cellular positioning for traffic state acquisition. *Journal of Location Based Services*, vol 1(2), pp 133-144.
- READES, J., CALABRESE, F., SEVTSUK, A., RATTI, C. (2007). Cellular Census: Explorations in Urban Data Collection, *IEEE Pervasive Computing*, vol. 6, no. 3, pp. 30-38.
- SALLENT, O., VALENZUELA, J.L., AGUSTÍ, R. (2003). *Principios de Comunicaciones Móviles*. Ediciones UPC, Barcelona, Spain.
- SCHOLLMEYER, R., WILTSCHKO, T. (2007). Classification of public transport vehicles using cellular cellular radio data. 6th European Congress & Exhibition on Intelligent Transport Systems and Services, ITS Congress, CD-ROM paper no. 2434.
- SMITH, B.L., FONTAINE, M.D. (2006). Private-Sector Provision of Congestion Data. NCHRP 70-01, University of Virginia, US.
- THIESSENHUSEN, K.U., SCHÄFER, R.P., LANG, T. (2006). Traffic data from cell phones: a comparison with loops and probe vehicle data. *ITS World Congress*, paper no. 1550.
- UMTSWorld. UMTS Network Coverage Planning. (Accessed on November 18, 2009). <http://www.umtsworld.com/technology/coverage.htm>
- WAKEFIELD, T., MCNALLY, D., BOWLER, D., MAYNE, A. (2007). *Introduction to Cellular Communications: Technology, Services, Markets*. Auerbach Publications (Informa Telecoms & Media). Taylor & Francis Group. Boca Raton, New York.
- YACOUB, M.D. (2002). *Wireless Technology: Protocols, Standards and Techniques*. CRC Press.
- YAIPAHOJ, S., HARMANTZIS, F. (2008). Achieving low-cost UMTS networks via pricing, *Netnomics*, vol 9(2), pp. 105-124.