



# Development of origin–destination matrices using mobile phone call data



Md. Shahadat Iqbal<sup>a</sup>, Charisma F. Choudhury<sup>a,b,c,\*</sup>, Pu Wang<sup>b,d</sup>, Marta C. González<sup>b</sup>

<sup>a</sup> Department of Civil Engineering, Bangladesh University of Engineering and Technology, Bangladesh

<sup>b</sup> Department of Civil and Environmental Engineering, Massachusetts Institute of Technology, Cambridge, USA

<sup>c</sup> Institute for Transport Studies, University of Leeds, UK

<sup>d</sup> School of Traffic and Transportation Engineering, Central South University, Hunan, PR China

## ARTICLE INFO

### Article history:

Received 23 August 2013

Received in revised form 12 December 2013

Accepted 2 January 2014

### Keywords:

Mobile phone

Origin–destination

Video count

Traffic microsimulation

## ABSTRACT

In this research, we propose a methodology to develop OD matrices using mobile phone Call Detail Records (CDR) and limited traffic counts. CDR, which consist of time stamped tower locations with caller IDs, are analyzed first and **trips occurring within certain time windows are used to generate tower-to-tower transient OD matrices for different time periods**. These are then associated with corresponding nodes of the traffic network and converted to node-to-node **transient OD matrices**. The actual OD matrices are derived by scaling up these node-to-node **transient OD matrices**. An optimization based approach, in conjunction with a microscopic traffic simulation platform, is used to determine the scaling factors that result best matches with the observed traffic counts. The methodology is demonstrated using CDR from 2.87 million users of Dhaka, Bangladesh over a month and traffic counts from 13 key locations over 3 days of that month. **The applicability of the methodology is supported by a validation study.**

© 2014 Elsevier Ltd. All rights reserved.

## 1. Background

Reliable Origin–destination (OD) matrices are critical inputs for analyzing transportation initiatives. Traditional approaches of developing OD matrices rely on roadside and household surveys, and/or traffic counts. The roadside and household surveys for origin destination involve expensive data collection and thereby have limited sample sizes and lower update frequencies. Moreover, they are prone to sampling biases and reporting errors (e.g. Hajek, 1977; Kuwahara and Sullivan, 1987; Groves, 2006). Estimation of reliable OD matrices from traffic link count data on the other hand is extremely challenging since very often the data is limited in extent and can lead to multiple plausible non-unique OD matrices (e.g. Lo et al., 1996; Van Zuylen and Willumsen, 1980; Chungcheng et al., 2013; Caggiani et al., 2013; Nie et al., 2005). A number of Bayesian methods (e.g. Maher, 1983; Tebaldi and West, 1998; Li, 2005), Generalized Least Squares approaches (e.g. Cascetta, 1984; Bell, 1991; Toledo and Kolehkina, 2013), Maximum Likelihood Approaches (e.g. Spiess, 1987), and Correlation Methods (e.g. Vardi, 1996; Hazelton, 2000, 2003) have been used to tackle the indeterminacy problem. These approaches typically use target matrices based on prior information for generating the plausible route flows and are very sensitive to this prior information as well as to the chosen methodology (e.g. Hazelton, 2001). More recent approaches for OD estimation include automated registration plate scanners (e.g. Castillo et al., 2008) and mobile traffic sensors such as portable GPS devices (e.g. Parry and Hazelton, 2012; Morimura and Kato, 2012; Herrera et al., 2010). The practical successes of these approaches

\* Corresponding author at: Institute for Transport Studies, University of Leeds, UK. Tel.: +44 1133432659.

E-mail addresses: [cfc@alum.mit.edu](mailto:cfc@alum.mit.edu), [c.f.choudhury@leeds.ac.uk](mailto:c.f.choudhury@leeds.ac.uk) (C.F. Choudhury).

have however been limited due to high installation costs of the license plate readers and the low penetration rates of GPS devices (especially in developing countries).

Mobile phone users on the other hand also leave footprints of their approximate locations whenever they make a call or send an SMS. Over the last decade, mobile phone penetration rates have increased manifold both in developed and developing countries: the current penetration rates being 128% and 89% in developed and developing countries respectively (e.g. International Telecommunication Union, 2013). Subsequently, mobile phone data has emerged as a very promising source of data for transportation researchers. In recent years, mobile phone data have been used for human travel pattern visualization (e.g. Phithakkitnukoon et al., 2010; Phithakkitnukoon and Ratti, 2011; Reades et al., 2009; Asakura and Hato, 2004), mobility pattern extraction (e.g. Wang et al., 2012; González et al., 2008; Song et al., 2010; Simini et al., 2012; Candia et al., 2008; Sevtsuk and Ratti, 2010; Asgari et al., 2013; Calabrese et al., 2013), route choice modeling (e.g. Schlaich et al., 2010; Becker et al., 2011), traffic model calibration (e.g. Bolla et al., 2000), traffic flow estimation (e.g. Demissie et al., 2013; Cheng et al., 2006) to name a few. There have been several limited scale researches to explore the feasibility of application of mobile phone data for OD estimation as well. Wang et al. (2011) for instance use a correlation based approach to dynamically update a prior OD matrix using time difference of phone signal receipt times of base stations and Caceres et al. (2007) use a GSM network simulator to simulate the detailed movements of phones that are turned on. But both of these feasibility studies are based on synthetic data in small networks and the practical application is challenging given the need to collect and process detailed location data (which are currently processed by the mobile phone companies for load management purposes but are not stored). The potential to estimate OD matrices using mobile phone Call Detail Records (CDR) (which are stored by operators for billing purposes and hence more readily available) have also been explored (e.g. Mellegard et al., 2011; Calabrese et al., 2011; Wang et al., 2012). Mellegard et al. (2011) have developed an algorithm to assign mobile phone towers extracted from CDR to traffic nodes and Calabrese et al. (2011) have proposed a methodology to reduce the noise in the CDR data but both studies have focused more on computation issues and the relationship between the mobile phone OD and the traffic OD have not been explored in detail. Wang et al. (2012) have used an analytical model to scale up the ODs derived from CDR by using the population, mode choice probabilities and vehicle occupancy and usage ratios and have validated it using probe vehicle data. The methodology however relies heavily on availability of traffic and demographic data in high spatial resolution which may not be always available, particularly in developing countries.

In this research, we propose a methodology to develop OD matrices using mobile phone CDR and limited traffic counts. CDR from 2.87 million users from Dhaka, Bangladesh over a month are used to generate the OD patterns on different time periods and traffic counts from 13 key locations of the city over a limited time are used to scale it up to derive the actual ODs using the microscopic traffic simulator MITSIMLab. The methodology is particularly useful in situations when there is limited availability of detailed travel survey and high resolution traffic data. The ODs are validated by comparing the simulated and observed traffic counts of a different location (which have not been used for calibration).

The rest of the paper is organized as follows. First we describe the data followed by the methodology used for development of the OD matrix. The estimation and validation results are presented next. We conclude with the summary of findings and directions for future research.

## 2. Data

### 2.1. Study area

The central part of the Dhaka city has been selected as the study area and the major roads in the network has been coded. This consists of 67 nodes and 215 links covering an area of about 300 km<sup>2</sup> with a population of about 10.7 million (e.g. DHUTS, 2010). The average trip production rate is 2.74 per person per day with significant portions of walking (19.8%) and non-motorized transport trips (38.3%) (e.g. DHUTS, 2010). The traffic is subjected to severe congestion in most parts of the day, the average speed being only 17 km/h.<sup>1</sup>

The mobile phone penetration rate is approximated to be more than 90% in Dhaka (66.36% being the national average) and Grameenphone Ltd. has the highest market share with 42.7 m mobile phone subscribers nationwide (e.g. Grameenphone Ltd., 2012).

### 2.2. CDR data

The CDR data, collected from Grameenphone Ltd, consists of calls from 6.9 million users (which are more than 65% of the population of the study area) over a month. This comprises of 971.33 million anonymized call records in total made in between June 19, 2012 and July 18, 2012. The majority of the users (63%) have made 100 calls or less over the month. The frequencies of users making certain number of calls over the month and on a randomly selected day (15th July, 2012) are presented in Fig. 1. It may be noted that no demographic data related to the phone users are available.

<sup>1</sup> Excluding the non-motorized vehicles which are restricted from entering the major roads.

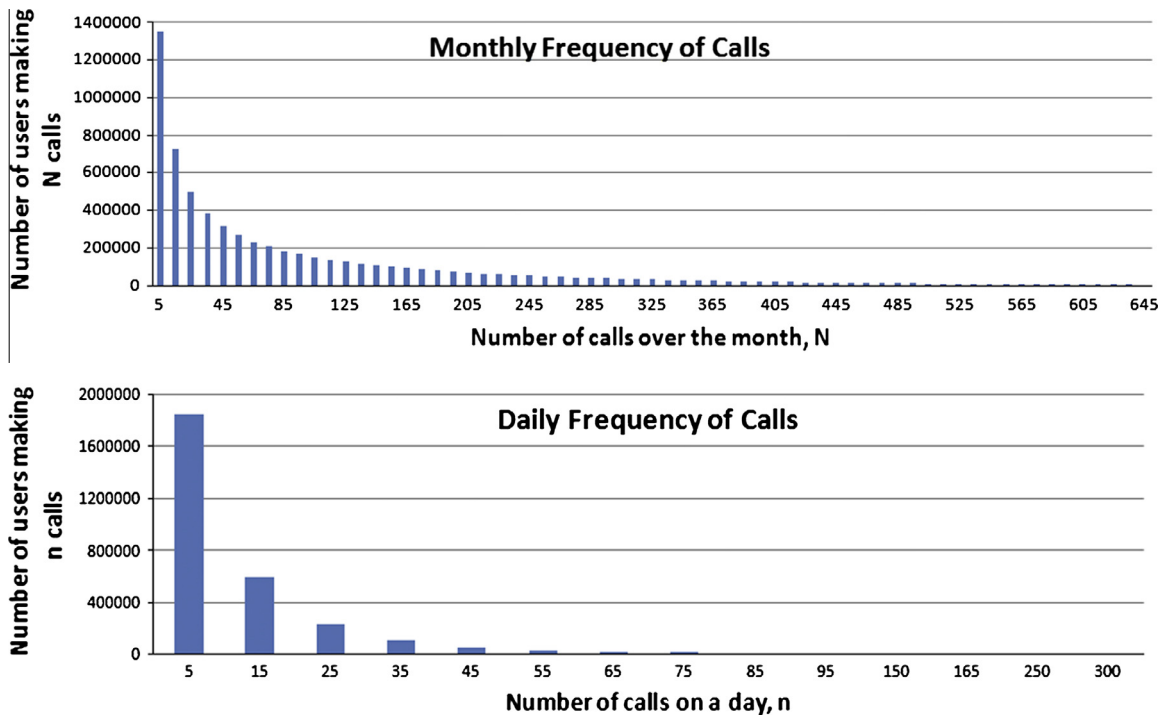


Fig. 1. Frequency of calls per user.

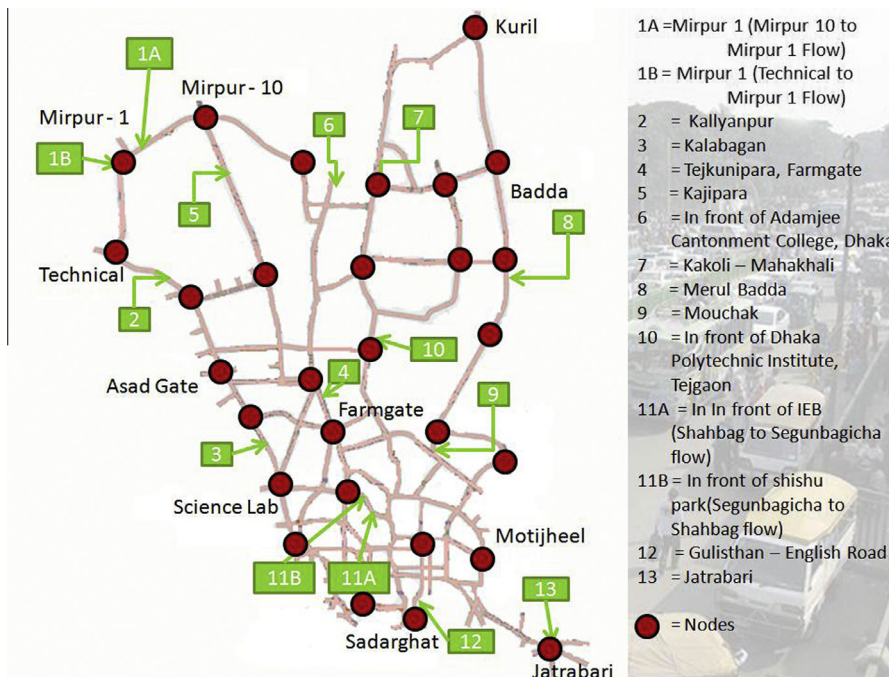


Fig. 2. Locations of video data collection and position of OD generating nodes.

### 2.3. Traffic count data

Video data, collected from 13 key locations of Dhaka city network over 3 days (12th, 15th and 17th July 2012) have been used in this study to extract the traffic counts.<sup>2</sup> The locations (shown in Fig. 2) have been selected such that they cover the

<sup>2</sup> There are no loop detectors or any other automatic traffic counters in Dhaka.

major roads (links) of Dhaka city with flows from major generators and governed by the availability of foot over bridges for mounting video cameras. Care has been taken to avoid roads that have high percentages of non-motorized transport and where lane-discipline is not strictly followed since simulation errors are likely to be higher in these situations due to increased complexity of acceleration and lane-changing behavior of drivers. The data has been collected for 8 h (8.00 am to 12.00 noon and 3.00 pm to 7.00 pm) and analyzed using the software TRAZER (e.g. [Kritikal Solutions Ltd., 2012](#)) to generate classified vehicle counts. Due to inclement weather and poor visibility some portion of the data is non-usable though. Moreover, TRAZER (which is the only commercial software that can deal with mixed traffic streams with ‘weak’ lane discipline) has high misspecification rates in presence of high congestion levels and in those cases, manual counting has been performed instead.

### 3. Methodology

Each entry in the CDR contains unique caller id (anonymized), the date and time of the call, call duration and latitude and longitude of the Base Transceiver Station (BTS). A snapshot of the data is presented in [Fig. 3](#). As seen in the figure, if a person traverses within the city boundary and uses his/her phone from different locations that is captured in the CDR. CDR can thus provide an abstraction of his/her physical displacements over time ([Fig. 3](#)).

However, in the CDR data, a user's location information is lost when he/she does not use his/her phone. As shown in [Fig. 4](#), according to the CDR, a user may be observed to move from zone B to zone E, but his/her initial origin (O) and final destination (D) may actually be located in zone A and zone F. In such cases, segments of the trip information are unobserved in the CDR. However, the mobile phone call records enable us to capture the *transient* origins and destinations which may have the true origin and/or the destination missing for a portion of the person's itinerary, but still retain a large portion of the actual ODs. Thus, we use the concept of transient origin destination (*t*-OD) matrix (as used by [Wang et al. \(2012\)](#)), which uses the mobile phone data to efficiently and economically capture the pattern of travel demand.

The second source of data used in this research is classified traffic counts extracted from video recordings collected from 13 key locations of Dhaka. These counts represent the *ground truth* but are more expensive to collect and limited in extent (only 3 days in this case). This limited point source data therefore cannot be used as a stand-alone source to reliably capture the OD pattern.

In this research, we therefore plan to combine the two data sources. The OD pattern is generated using the CDR data and scaled up to match the traffic counts. The scaling factors are determined using a microscopic traffic simulator platform MIT-SIMLab (e.g. [Yang and Koutsopoulos, 1996](#)) using an optimization based approach which aims to minimize the differences between observed and simulated traffic counts at the points where the traffic counts are available.

The methodology is summarized in [Fig. 5](#) and described in the subsequent sections.

#### 3.1. Generation of tower-to-tower transient OD matrix

The time-stamped BTS tower locations of each user are first extracted from the mobile phone CDR data and used for generating tower-to-tower transient OD matrix. The CDR however contains sparse and irregular records (e.g. [Candia et al., 2008](#)), in which user displacements (consecutive non-identical locations) are often observed with long travel intervals i.e. the first location may be observed at 8:56 and next location may be observed at 18:03 with no information about intermediate locations (if any) or the time when the trip in between these two locations have been made.

| ID                 | Call Date | Call Time | Duration | Latitude | Longitude |
|--------------------|-----------|-----------|----------|----------|-----------|
| AH03JAC8AAAbXtAId  | 20120701  | 09:34:19  | 18       | 23.8153  | 90.4181   |
| AAH03JABiAAJKnPAa5 | 20120707  | 06:15:20  | 109      | 23.8139  | 90.3986   |
| AAH03JABiAAJKnPAa5 | 20120707  | 09:03:06  | 109      | 23.7042  | 90.4297   |
| AAH03JABiAAJKnPAa5 | 20120707  | 10:34:19  | 16       | 23.6989  | 90.4353   |
| AAH03JABiAAJKnPAa5 | 20120707  | 18:44:53  | 154      | 23.6989  | 90.4353   |
| AAH03JABiAAJKnPAa5 | 20120707  | 20:00:08  | 154      | 23.8092  | 90.4089   |
| AAH03JAC5AAAdAYAE  | 20120701  | 09:15:05  | 62       | 23.7428  | 90.4164   |
| AAH03JAC+AAAcVKAC  | 20120707  | 08:56:34  | 242      | 23.7908  | 90.3753   |
| AAH03JAC+AAAcVKAC  | 20120701  | 18:03:06  | 36       | 23.9300  | 90.2794   |
| AAH03JAC5AAAdAYAA  | 20120701  | 11:15:55  | 12       | 23.7428  | 90.4164   |

**Fig. 3.** An excerpt from CDR data (entries of the same user are highlighted) and locations of a random user “AAH03JABiAAJKnPAa5” throughout the day as observed in data.



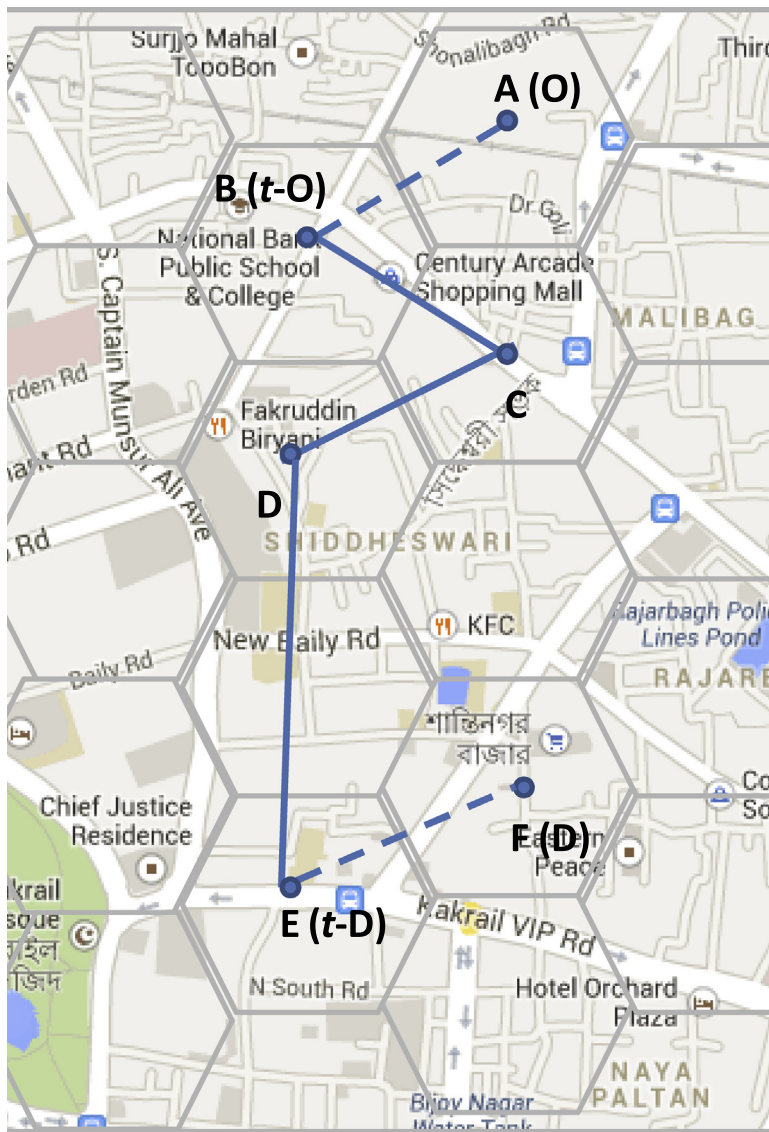


Fig. 4. Actual vs. transient OD.

Another limitation of the CDR data is often there are changes in tower in the data in spite of no actual displacement. This is because the operator often **balances call traffic** among adjacent towers by allocating a new call (or shifting an ongoing call) to the tower that is handling lower call volumes at that moment. In the CDR, such switches in towers are confounded with the actual physical displacements. **To reduce the number of false displacements and better identify timing and origin–destinations of specific trips, we therefore extract displacements that have occurred within a specific time window. A lower bound in the time window (10 min) is imposed to reduce the number of false displacements without affecting the number of physical displacements occurring within short intervals. An upper bound in the time window (1 h) is imposed to ensure that meaningful numbers of trips are retained.** Therefore, a person trip is recorded if in the CDR, subsequent entries of the same user indicate a displacement (change in tower) with a time difference of more than 10 min but less than 1 h.

Further, both call volumes (from CDR data) and traffic volumes (from traffic counts) had significant variations throughout the day. Based on correlation analysis of total mobile call volumes and total traffic counts (Fig. 6), four time periods (7:00–9:00, 9:00–12:00, 15:00–17:00 and 17:00–19:00), have been chosen for analysis.

### 3.2. Conversion of tower-to-tower $t$ -OD to node-to-node $t$ -OD

For application of the  $t$ -ODs in traffic analyses, the origin and destination towers need to be associated with corresponding nodes of the traffic network. The typical tower coverage area can be represented as a combination of three hyperbolas

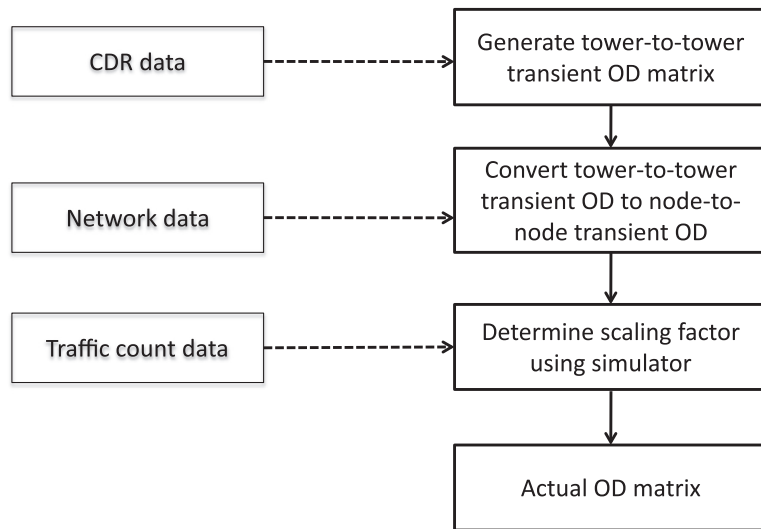


Fig. 5. Framework for developing OD Matrix.

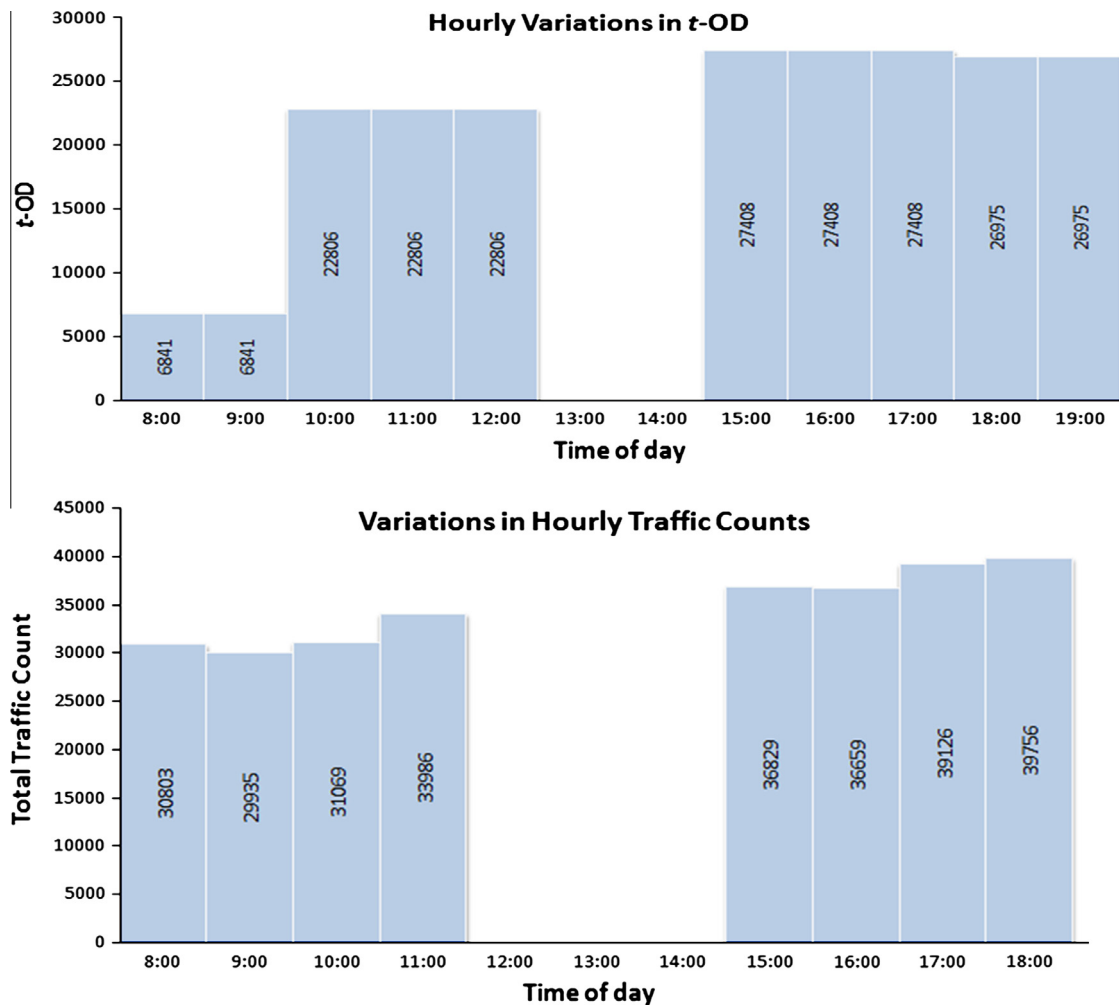


Fig. 6. Hourly variations: (a) traffic count; (b) transient ODs from mobile call records.

(Fig. 7), the size varying depending on tower height, terrain, locations of adjacent towers and number of users active in the proximity (which can vary dynamically).

The population density in the chosen study area is very high (more than 8111 inhabitants/sq. km (e.g. Bangladesh Bureau of Statistics, 2011) and the tower locations are very close to each other (1 km on average). Because of the high user density, it can be assumed that the area between two towers is equally split among the two towers (Fig. 8) that is, each tower  $t$  has a coverage area ( $A_t$ ) approximately defined by a circle of radius  $0.5l$ , where  $l$  is the tower-to-tower distance.

If a unique traffic node  $i$  overlaps with  $A_t$ , the calls handled by  $t$  are associated with node  $i$  (as in the case of Tower 1 in Fig. 8). However, if  $A_t$  has two (or more) candidate nodes for association, then the candidate nodes are ranked based on the proportion of  $A_t$  feeding to each node. That is, the node serving greatest portion of  $A_t$  is ranked 1, the node serving second

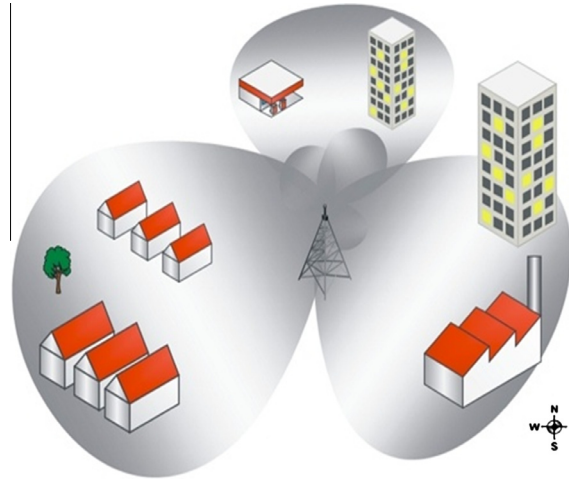
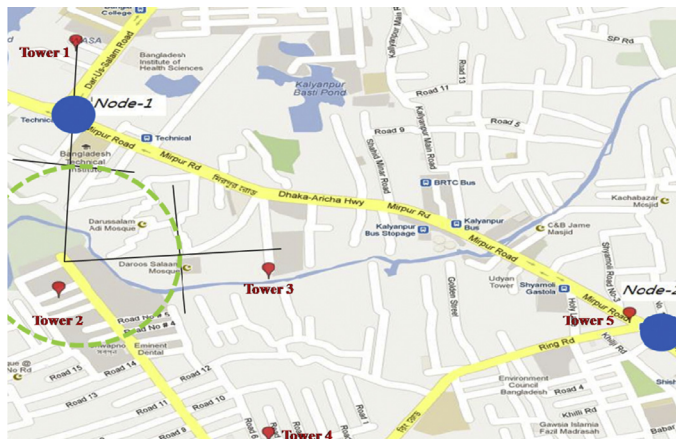


Fig. 7. Typical coverage area of a tower ([http://www.truteq.co.za/tips\\_gsm/](http://www.truteq.co.za/tips_gsm/)).



| ID      | Call Date | Call Time | Origin Tower | Destination Tower |
|---------|-----------|-----------|--------------|-------------------|
| AAH03JA | 20120718  | 15:54     | 6            | 1                 |
| AAH03JA | 20120718  | 16:13     | 1            | 2                 |
| AAH03JA | 20120718  | 16:15     | 2            | 1                 |
| AAH03JA | 20120718  | 18:53     | 1            | 6                 |
| AAH03JA | 20120718  | 20:49     | 6            | 1                 |
| AAH03JA | 20120718  | 23:41     | 1            | 6                 |

(a) Tower-to-tower OD

| ID      | Call Time | Origin |                | Destination |                |
|---------|-----------|--------|----------------|-------------|----------------|
|         |           | Tower  | Candidate Node | Tower       | Candidate Node |
| AAH03JA | 14:54     | 6      | 3              | 1           | 1              |
| AAH03JA | 16:13     | 1      | 1              | 2           | 2 Or 1         |
| AAH03JA | 16:15     | 2      | 2 Or 1         | 1           | 1              |
| AAH03JA | 18:53     | 1      | 1              | 6           | 3              |
| AAH03JA | 20:49     | 6      | 3              | 1           | 1              |
| AAH03JA | 23:41     | 1      | 1              | 6           | 3              |

(b) Intermediate OD with candidate nodes

| ID      | Call Time | Origin Node | Destination Node |
|---------|-----------|-------------|------------------|
| AAH03JA | 14:54     | 3           | 1                |
| AAH03JA | 16:13     | 1           | 1                |
| AAH03JA | 16:15     | 1           | 1                |
| AAH03JA | 18:53     | 1           | 3                |
| AAH03JA | 20:49     | 3           | 1                |
| AAH03JA | 23:41     | 1           | 3                |

(c) Node-to-node OD

Fig. 8. Example of tower to node allocation.

highest portion of  $A_t$  is ranked 2, etc. For example, in Fig. 8, network connectivity (feeder roads) and topography (presence of a canal with no crossing facility in the vicinity) denote that Node 1 and Node 2 are candidate nodes for association with Tower 2. As the major portion of  $A_t$  is connected to Node 2 and the remaining portion is connected to Node 1, they are ranked 1 and 2 respectively for Tower 2. The data format after this step is presented in Fig. 8. As seen in the figure, this typically consists of call records associated with unique nodes and in some cases, a few calls associated with *multiple candidate nodes*. The calls are then sorted and ranked based on the frequency of the unique nodes used by each user (based on analysis of his/her call locations over the month). The frequency of occurrence of the candidate nodes are compared and used as the basis of replacement. For example, frequency analysis of User “AAH03JA” indicates a higher frequency of Node 1. Therefore, in cases where there are ambiguities between Nodes 2 and 1, Node 1 is used (for this particular user). The tower-to-node allocation-sure thus user-specific and uses information derived from each user’s overall travel pattern.

The same process is used for all users and node-to-node  $t$ -OD matrices for each time period of each day are derived.

### 3.3. Finding the scaling factor and determining the actual OD matrix

As discussed, the node-to-node  $t$ -OD matrix ( $t - OD_{ij}$ ) provides the trip patterns for developing the actual OD matrix ( $OD_{ij}$ ). However, in order to determine the actual OD matrix, the  $t$ -OD needs to be scaled to match the real traffic flows. A scaling factor  $\beta_{ij}$  is used in this regard:

$$OD_{ij} = \sum_{ij} (t-OD_{ij}) * \beta_{ij}$$

It may be noted that  $\beta_{ij}$  takes into account the market penetration rates (i.e. not every user has a mobile phone or uses the specific service provider), the mobile phone non-usage issue (i.e. mobile phone calls are not made from every location traversed by the user), the vehicle usage issue (i.e. users may not use cars for every trip), etc. The potential error introduced due to *false displacement* (described in Section 2.1) is also accounted for in the scaling factors.

The scaling factors are determined using the open-sourced microscopic traffic simulator platform MITSIMLab (e.g. Yang and Koutsopoulos, 1996) by applying an optimization based approach. The movements of vehicles in MITSIMLab are dictated by driving behavior models based on decision theories and estimated with detailed trajectory data using econometric approaches. Route choices of drivers are based on a discrete choice based probabilistic model where the utilities of selecting and re-evaluating routes are functions of path attributes, such as path travel times (Bar-Gera, 2007; Zhan et al., 2013) and freeway bias (see Ben-Akiva et al., 2010 for details). The inputs of the simulator include network data, driving behavior parameters and OD matrix. The generated outputs include traffic flow at specified locations in the network. It may be noted that MITSIMLab was calibrated prior to the OD estimation and the desired speed, acceleration and lane-changing model constants have been updated using detailed video trajectories (e.g. Iqbal, 2013; Islam, 2013; Siddique, 2013). The node-to-node OD matrix derived from the mobile phone data are provided as the initial or seed-OD in this case. The simulated traffic flows are compared with the actual traffic flows extracted from video recordings. The objective function seeks to minimize the difference between the actual and simulated traffic flows in each location by changing the scaling factors. The optimization problem can be represented as follows:

$$\text{Minimize, } Z = \sum_{k=1}^K (V_{actual}^k - V_{simulated}^k)^2 \quad (1)$$

$$\text{Such that, } OD_{ij,t} = \sum_{i,j=1}^N t-OD_{ij,t} * \beta_{ij}$$

where,  $V_{simulated}^k$  = traffic flow of link  $k$  of the road network from simulation;  $OD_{ij,t}$  = actual OD between nodes  $i$  and  $j$  in time period  $t$ ;  $t - OD_{ij,t}$  = transient OD between nodes  $i$  and  $j$  in time period  $t$ ;  $\beta_{ij,t}$  = scaling factor associated with the node pair  $i$  and  $j$  and time period  $t$ ;  $K$  = total number of links for which traffic flow data is available;  $N$  = total number of nodes in the network.

However, to make the optimization problem more tractable, group-wise scaling factors are used rather than an individual scaling factor for each OD pair. The grouping is based on the analyses of the CDR data. This simplifies the problem as follows:

$$\text{Minimize, } Z = \sum_{k=1}^K (V_{actual}^k - V_{simulated}^k)^2 \quad (2)$$

$$\text{Such that, } OD_{ij,t} = \sum_{m=1}^M t-OD_{ij,t}^m * \beta_t^m$$

where,  $t - OD_{ij,t}^m$  = transient OD between node pair  $i$  and  $j$  in time period  $t$  where the node pair  $ij$  belong to group  $m$ ;  $\beta_t^m$  = scaling factor for group  $m$  and time period  $t$ ;  $M$  = total number of groups of OD-pairs.



#### 4. Results

The mobile phone network within the study area comprises of 1360 towers which have been assigned to 29 OD generating nodes (812 OD pairs). Out of the one month CDR data, the weekend data have been discarded. For each day, the calls of each user originating from two different towers in each of the time period have been extracted. After application of the transient trip definitions (displacements occurring more than 10 mins but less than 1hr apart) and the tower to node conversion rules (elaborated in Section 3.2), the node-to-node  $t$ -ODs are derived. The total number of node-to-node  $t$ -ODs is presented in Table 1.

Analyses of the node-to-node transientflows indicate that the flows between adjacent nodes are substantially higher than those between non-adjacent nodes (Fig. 9). This is reasonable since given the low travel speed in Dhaka, a traveler may not be able to move very far in the 50 min time window and the  $t$ -ODs mostly capture segments of a longer trip. However, part of it may also be due to the *false displacement* problem discussed in Section 3.1. Therefore, the OD-pairs have been divided into two groups (adjacent and non-adjacent nodes) and the objective function to determine scaling factors has been formulated as follows:

$$\text{Minimize, } Z = \sum_{k=1}^K (V_{actual}^k - V_{simulated}^k)^2 \quad (3)$$

$$\text{Such that, } OD_{ij,t} = \sum_{adj} t-OD_{ij,t}^{adj} * \beta_t^{adj} + \sum_{non-adj} t-OD_{ij,t}^{non-adj} * \beta_t^{non-adj}$$

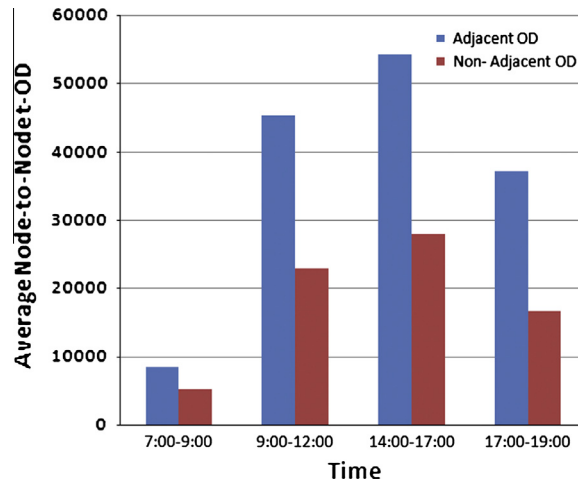
where,  $t-OD_{ij}^{adj}$  = transient OD between node pair  $i$  and  $j$  in time period  $t$  where the node pair  $ij$  are adjacent nodes;  $t-OD_{ij}^{non-adj}$  = transient OD between node pair  $i$  and  $j$  in time period  $t$  where the node pair  $ij$  are non-adjacent nodes;  $\beta_t^{adj}, \beta_t^{non-adj}$  = scaling factors for time period  $t$  and adjacent and non-adjacent nodes respectively.

This yielded eight scaling factors in total that needed to be estimated from the simulation runs of MITSIMLab. Running the optimization process in MATLAB (that invokes MITSIMLab) and using a BOX algorithm (e.g. Box, 1965), the following values of scaling factors have been derived.

It is interesting to note that the scaling factors for adjacent nodes are higher than those of non-adjacent in all time periods other than 15:00–17:00. This does not however indicate that most of the actual trips are to the adjacent nodes since a full trip may consist of several segments each represented by a separate  $t$ -OD.

**Table 1**  
Node-to-node  $t$ -OD.

| Time period | Time        | $t$ -OD                                      |                 |
|-------------|-------------|--|-----------------|
|             |             | Total over the month<br>(including weekends) | Weekday average |
| 1           | 7:00–9:00   | 397355                                       | 13681.86        |
| 2           | 9:00–12:00  | 1915417                                      | 68418.48        |
| 3           | 15:00–17:00 | 2255859                                      | 82226.05        |
| 4           | 17:00–19:00 | 1549109                                      | 53950.57        |



**Fig. 9.** Comparison of  $t$ -ODs between adjacent and non-adjacent nodes.

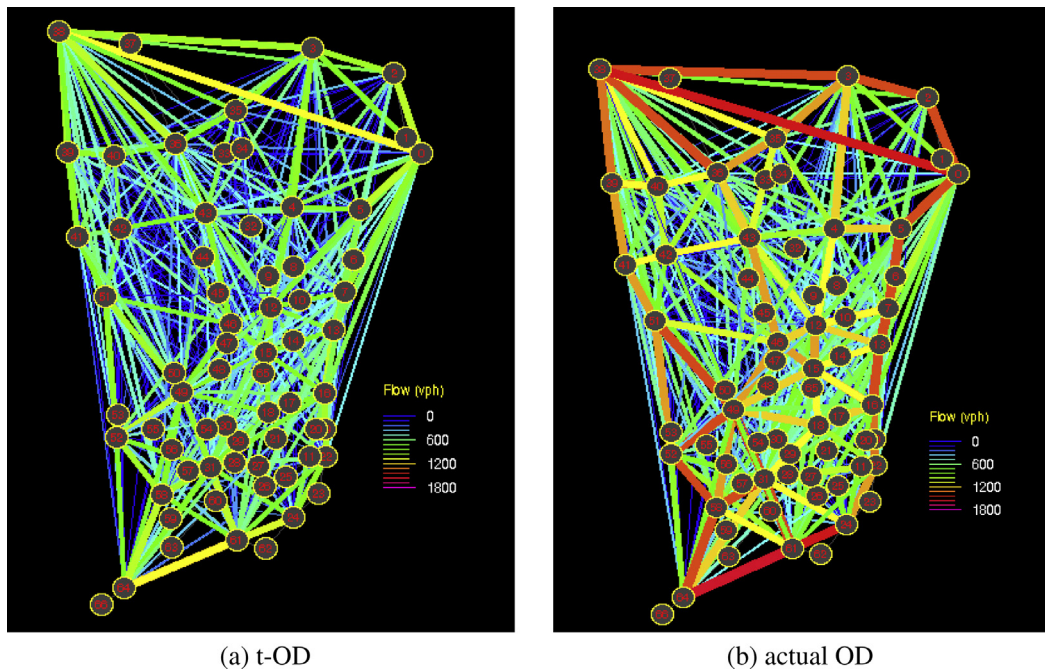


Fig. 10. t-ODs and actual ODs across the network for 7:00–9:00.

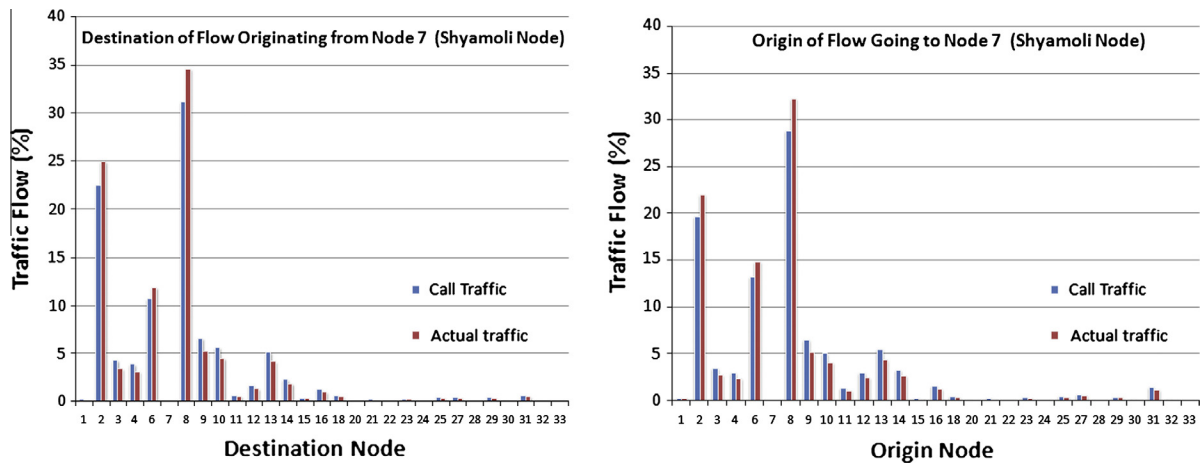


Fig. 11. Example of transient and actual traffic flows to and from a node (Shyamoli) between 7:00 and 9:00.

The graphical representation of the t-ODs and actual ODs across the network for one of the time periods and the variations for an example node are presented in Figs. 10 and 11 respectively.

## 5. Validation

In addition to the aggregate data used for calibration, traffic counts are collected from four additional locations on a different day (not used for calibration). For validation purposes, the scaled up ODs have been applied to simulate the traffic between 9:00–12:00 in MITSIMLab and the simulated traffic counts are compared against the observed counts from these locations. In order to quantify the prediction error, **Root Mean Square Error** and **Root Mean Square Percent Errors** have been calculated and are found to be 335.09 and 13.59% respectively (see Table 2).

**Table 2**  
Scaling factors.

| Time period | OD type      | Scaling factor |
|-------------|--------------|----------------|
| 7:00–9:00   | Adjacent     | 6.787          |
|             | Non-adjacent | 1.712          |
| 9:00–12:00  | Adjacent     | 0.971          |
|             | Non-adjacent | 0.345          |
| 15:00–17:00 | Adjacent     | 1.647          |
|             | Non-adjacent | 3.407          |
| 17:00–19:00 | Adjacent     | 9.404          |
|             | Non-adjacent | 6.779          |

## 6. Conclusion

The main outcome of this research is the methodology for development of the OD matrix using mobile phone CDR and limited traffic count data. The strengths of both data sources are utilized in this approach: the trip patterns are extracted from mobile phones and the ground truth traffic scenario is derived from the counts. The methodology is demonstrated using data collected from Dhaka.

There are several limitations of the current research though. Firstly, in this research a simplified objective function with grouped scaling factors has been used. This overlooks the heterogeneity in call rates from different locations (e.g., more calls may be generated to and from railway stations compared to and from offices with land telephone lines, etc.). A more detailed classification of scaling factor can be used to overcome this bias and may yield better results. Moreover, in this particular context, detailed network data and extensive calibration data were not available which limited the number of traffic nodes used in the study. The transferability of the driving behavioral models in MITSIMLab have also not been tested in detail and only the key behavioral model constants have been updated to better match the traffic patterns in Dhaka (e.g. Iqbal, 2013; Islam, 2013). These factors may have increased the simulation errors and affected the validation results. However, initial validation results indicate promising success in real life application by transport planners and managers. It may be noted that though MITSIMLab has been used in this study to determine the scaling factors, the developed ODs are simulator independent.

Since CDR is already recorded by mobile phone companies for billing purposes, the approach is more economic than the traditional approaches which rely on expensive household surveys and/or extensive traffic counts. It is also convenient for periodic update of the OD matrix and extendable for dynamic OD estimation. This method is particularly effective for generating complex OD matrix where land use pattern is heterogeneous and asymmetry in traveling pattern prevails throughout the day but there is a limitation of traditional data sources.

## Acknowledgments

The data provided for the research has been provided by Grameenphone Ltd., Bangladesh. The funding for this research was provided by Faculty for the Future Program of Schlumberger Foundation and Higher Education Enhancement Project of the University Grants Commission of Bangladesh and the World Bank, the National Natural Science Foundation of China and the New England UTC.

## References

- Asakura, Yasuo, Hato, Eiji, 2004. Tracking survey for individual travel behavior using mobile communication instruments. *Transp. Res. Part C* 12 (3–4), 273–291.
- Asgari, F., Gauthier, V., Becker, M., 2013. A Survey on Human Mobility and Its Applications. arXiv Preprint arXiv:1307.0814.
- BBS, 2011. Population and Housing Census: Preliminary Results, 2011. Bangladesh Bureau of Statistics, Statistics Division, Ministry of Planning, Government of the People's Republic of Bangladesh.
- Bar-Gera, Hillel., 2007. Evaluation of a cellular phone-based system for measurements of traffic speeds and travel times: a case study from Israel. *Transp. Res. Part C* 15 (6), 380–391.
- Becker, R.A., Caceres, R., Hanson, K., Loh, J.M., Urbanek, S., Varshavsky, A., Volinsky, C., Ave, P., Park, F., 2011. Route classification using cellular handoff patterns. In: Proceedings of the 13th International Conference on Ubiquitous Computing. ACM, Beijing, China.
- Bell, M., 1991. The estimation of origin–destination matrices by constrained generalized least squares. *Transp. Res. Part B* 25 (1), 13–22.
- Ben-Akiva, M., Koutsopoulos, H.N., Toledo, T., Yang, Q., Choudhury, C.F., Antoniou, C., Balakrishna, R., 2010. Traffic simulation with MITSIMLab, in fundamentals of traffic simulation. In: Barceló, J. (Ed.), *International Series in Operations Research and Management Science*, vol. 145. Springer, pp. 233–268.
- Bolla, R., Davoli, F., Giordano, A., 2000. Estimating road traffic parameters from mobile communications. In: Proceedings 7th World Congress on ITS, Turin, Italy.
- Box, M.J., 1965. A new method of constrained optimization and a comparison with other methods. *Comput. J.* 8 (1), 42–52.
- Caceres, N., Wideberg, J.P., Benitez, F.G., 2007. Deriving origin destination data from a mobile phone network. *Intell. Transp. Syst., IET* 1 (1), 15–26.
- Caggiani, Leonardo, Ottomanelli, Michele, Sassanelli, Domenico, 2013. A fixed point approach to origin–destination matrices estimation using uncertain data and fuzzy programming on congested networks. *Transp. Res. Part C: Emerg. Technol.* 28 (March), 130–141.

- Calabrese, F., Lorenzo, G.D., Liu, L., Ratti, C., 2011. Estimating origin–destination flows using mobile phone location data. *IEEE Pervasive Comput.* 10 (4), 36–43.
- Calabrese, F., Diao, M., Lorenzo, G.D., Ferreira, J., Ratti, C., 2013. Understanding individual mobility patterns from urban sensing data: a mobile phone trace example. *Transp. Res. Part C* 26 (January), 301–313.
- Candia, J., González, M.C., Wang, P., Schoenharl, T., Madey, G., Barabási, A.L., 2008. Uncovering individual and collective human dynamics from mobile phone records. *J. Phys. A: Math. Theor.* 41 (22), 224015.
- Cascetta, E., 1984. Estimation of trip matrices from traffic counts and survey data: a generalized least squares estimator. *Transp. Res. Part B* 18 (4–5), 289–299.
- Castillo, E., Menéndez, J., Jiménez, P., 2008. Trip matrix and path flow reconstruction and estimation based on plate scanning and link observations. *Transp. Res. Part B* 42 (5), 455–481.
- Cheng, P., Qiu, Z., Ran, B., 2006. Particle filter based traffic state estimation using cell phone network data. In: *Proceedings of the IEEE ITSC 2006*.
- Chungcheng, L., Zhou, X., Zhang, K., 2013. Dynamic origin–destination demand flow estimation under congested traffic conditions. *Transp. Res. Part C: Emerg. Technol.* 34 (September), 16–37.
- Demissie, M.G., de Almeida Correia, G.H., Bento, C., 2013. Intelligent road traffic status detection system through cellular networks handover information: an exploratory study. *Transp. Res. Part C: Emerg. Technol.* 32, 76–88.
- DHUTS, 2010. Dhaka Urban Transport Network Development Study, Draft Final Report. Prepared by Katahira and Engineers International, Oriental Consultants Co. Ltd., and Mitsubishi Research Institute Inc.
- González, M.C., Hidalgo, C.A., Barabási, A.L., 2008. Understanding individual human mobility patterns. *Nature* 453, 779–782.
- Grameenphone Ltd., Bangladesh. <<http://grameenphone.com>> (accessed 15.12.12).
- Groves, R.M., 2006. Nonresponse rates and nonresponse bias in household surveys. *Publ. Opin. Quart.* 70 (5), 646–675.
- Hajek, J. J. (1977). Optimal Sample Size of Roadside-interview Origin–Destination Surveys (No. RR 208).
- Hazelton, M.L., 2000. Estimation of origin–destination matrices from link flows on uncongested networks. *Transp. Res. Part B* 34 (7), 549–566.
- Hazelton, M.L., 2001. Inference for origin–destination matrices: estimation, reconstruction and prediction. *Transp. Res. Part B* 35 (7), 667–676.
- Hazelton, M.L., 2003. Some comments on origin–destination matrix estimation. *Transp. Res. Part A* 37 (10), 811–822.
- Herrera, J., Work, D.B., Herring, R., Ban, X., Jacobson, Q., Bayen, A., 2010. Evaluation of traffic data obtained via GPS-enabled mobile phones: the Mobile Century field experiment. *Transp. Res. Part C: Emerg. Technol.* 18 (4), 568–583.
- International Telecommunication Union. <[http://www.itu.int/en/ITU-D/Statistics/Documents/facts/ICTFactsFigures 2013.pdf](http://www.itu.int/en/ITU-D/Statistics/Documents/facts/ICTFactsFigures%202013.pdf)> (accessed 20.07.13).
- Iqbal, S., 2013. Development of Origin–Destination Trip Matrices Using Mobile Phone Call Data. MSc Thesis, Bangladesh University of Engineering and Technology.
- Islam, M.M., 2013. Acceleration Decision in Heterogeneous Traffic Stream. MSc Thesis, Bangladesh University of Engineering and Technology.
- Kritikal Solutions Ltd., India. <<http://www.kritikalsolutions.com/products/traffic-analyzer.html>> (accessed 15.12.12).
- Kuwahara, M., Sullivan, E.C., 1987. Estimating origin–destination matrices from roadside survey data. *Transp. Res. Part B* 21 (3), 233–248.
- Li, B., 2005. Bayesian inference for origin–destination matrices of transport networks using the EM algorithm. *Technometrics* 47 (4), 399–408.
- Lo, H.P., Zhang, N., Lam, W.H., 1996. Estimation of an origin–destination matrix with random link choice proportions: a statistical approach. *Transp. Res. Part B* 30 (4), 309–324.
- Maher, M., 1983. Inferences on trip matrices from observations on link volumes: a Bayesian statistical approach. *Transp. Res. Part B* 20 (6), 435–447.
- Mellegard, E., Moritz, S., Zahoor, M., (2011). Origin/destination-estimation using cellular network data. In: *Data Mining Workshops (ICDMW), 2011 IEEE 11th International Conference on IEEE*, pp. 891–896.
- Morimura, T., Kato, S., 2012. Statistical origin–destination generation with multiple sources. In: *21st International conference on in pattern recognition (ICPR), November 11–15, Tsukuba, Japan*.
- Nie, Y., Zhang, H.M., Recker, W.W., 2005. Inferring origin–destination trip matrices with adecoupled GLS path flow estimator. *Transp. Res. Part B* 39, 497–518.
- Parry, K., Hazelton, M.L., 2012. Estimation of origin–destination matrices from link counts and sporadic routing data. *Transp. Res. Part B* 46 (1), 175–188.
- Phithakkitnukoon, S., Ratti, C., 2011. Inferring asymmetry of inhabitant flow using call detail records. *J. Adv. Inf. Technol.* 2 (4), 239–249.
- Phithakkitnukoon, S., Horanont, T., Di Lorenzo, G., Shibasaki, R., Ratti, C., 2010. Activity-aware map: identifying human daily activity pattern using mobile phone data. *Hum. Behav. Underst.* 6219 (3), 14–25.
- Reades, J., Calabrese, F., Ratti, C., 2009. Eigenplaces: analyzing cities using the space-time structure of the mobile phone network. *Environ. Plann. B: Plann. Des.* 36 (5), 824–836.
- Schlaich, J., Otterstätter, T., Friedrich, M., 2010. Generating trajectories from mobile phone data. In: *TRB 89th Annual Meeting Compendium of Papers. Transportation Research Board of the National Academies, Washington, DC, USA*.
- Sevtsuk, A., Ratti, C., 2010. Does urban mobility have a daily routine? Learning from aggregate data of mobile networks. *J. Urban Technol.* 17 (1), 41–60.
- Siddique, A.M., 2013. Lateral Movement Models for Heterogeneous Traffic Stream. MSc Thesis, Bangladesh University of Engineering and Technology.
- Simini, F., González, M.C., Marita, A., Barabási, A.L., 2012. A universal model for mobility and migration patterns. *Nature* 484, 96–100.
- Song, C., Koren, T., Wang, P., Barabási, A.L., 2010. Modelling the scaling properties of human mobility. *Nat. Phys.* 6, 818–823.
- Spiess, H., 1987. A maximum likelihood model for estimating origin–destination matrices. *Transp. Res. Part B* 21 (5), 395–412.
- Tebaldi, C., West, M., 1998. Bayesian inference on network traffic using link count data (with discussion). *J. Am. Stat. Assoc.* 93, 557–576.
- Toledo, T., Koleschikina, T., 2013. Estimation of dynamic origin–destination matrices using linear assignment matrix approximations. *IEEE Trans. ITS* 14 (2), 618–626.
- Van Zuylen, H.J., Willumsen, L.G., 1980. The most likely trip matrix estimated from traffic counts. *Transp. Res. Part B* 14 (3), 281–293.
- Vardi, Y., 1996. Network tomography: estimating source–destination traffic intensities from link data. *J. Am. Stat. Assoc.* 91, 365–377.
- Wang, J., Wang, D., Song, X., Di, Sun, 2011. Dynamic OD expansion method based on mobile phone location. In: *Fourth International Conference on Intelligent Computation Technology and Automation, Shenzhen, China*.
- Wang, P., Hunter, T., Bayen, A.M., Schechtner, K., González, M.C., 2012. Understanding Road Usage Patterns in Urban Areas. *Scientific Reports*, 2.
- Yang, Q., Koutsopoulos, H.N., 1996. A microscopic traffic simulator for evaluation of dynamic traffic management systems. *Transp. Res. C* 4 (3), 113–129.
- Zhan, Xianyuan, Hasan, Samiul, Ukkusuri, Satish V., Kamga, Camille, 2013. Urban link travel time estimation using large-scale taxi data with partial information. *Transp. Res. Part C* 33 (August), 37–49.