

# **Getting the best of both worlds - a framework for combining disaggregate travel survey data and aggregate mobile phone data for trip generation modelling**

## **Andrew Bwambale**

Choice Modelling Centre  
Institute for Transport Studies  
University of Leeds  
34-40 University Road, LS2 9JT, Leeds, United Kingdom  
Email: [ts13ab@leeds.ac.uk](mailto:ts13ab@leeds.ac.uk)

## **Charisma F. Choudhury**

Choice Modelling Centre  
Institute for Transport Studies  
University of Leeds  
34-40 University Road, LS2 9JT, Leeds, United Kingdom  
Email: [C.F.Choudhury@leeds.ac.uk](mailto:C.F.Choudhury@leeds.ac.uk)

## **Stephane Hess**

Choice Modelling Centre  
Institute for Transport Studies  
University of Leeds  
34-40 University Road, LS2 9JT, Leeds, United Kingdom  
Email: [S.Hess@its.leeds.ac.uk](mailto:S.Hess@its.leeds.ac.uk)

## **Md. Shahadat Iqbal**

Lehman Centre for Transportation Research  
Department of Civil and Environmental Engineering  
Florida International University  
10555 W. Flagler Street, EC 3729, Miami, FL 33174  
Email: [miqba005@fiu.edu](mailto:miqba005@fiu.edu)

Submission Date 15 May 2020

## 1 **Abstract**

2 Traditional approaches to travel behaviour modelling primarily rely on household travel survey  
3 data, which is expensive to collect, resulting in small sample sizes and infrequent updates.  
4 Furthermore, such data is prone to reporting errors which can lead to biased parameter  
5 estimates and subsequently incorrect predictions. On the other hand, mobile phone call detail  
6 records (CDRs), which report the timestamped locations of mobile communication events,  
7 have been successfully used in the context of generating travel patterns. However, due to their  
8 anonymous nature, such records have not been widely used in developing mathematical models  
9 establishing the relationship between the observed travel behaviour and influencing factors  
10 such as the attributes of the alternatives and the decision makers. In this paper, we propose a  
11 joint modelling framework that utilises the advantages offered by both travel survey data and  
12 low-cost CDR data to optimise the prediction capacity of traditional trip generation models. In  
13 this regard, we develop a model that jointly explains the reported trips for each individual in  
14 the household survey data and ensures that the aggregated zonal trip productions are close to  
15 those derived from CDR data. This framework is tested using data from Dhaka, Bangladesh  
16 consisting of household survey data (65419 persons in 16750 households), mobile phone CDR  
17 data (over 600 million records generated by 6.9 million users), and aggregate census data. The  
18 model results show that the proposed framework improves the spatial and temporal  
19 transferability of the joint models over the base model which relies on household travel survey  
20 data alone. This serves as a proof-of-concept that augmenting travel survey data with mobile  
21 phone data holds significant promise for the travel behaviour modelling community, not only  
22 by saving the cost of data collection, but also improving the prediction capability of the models.

23

24

25 *Keywords:* Trip generation, CDR data, mobile phone data, household travel survey data, census  
26 data, population synthesis, transferability, Bangladesh, developing country

27

28

29

30

31

32

33

## 34 **Acknowledgements**

35 The research in this paper used mobile phone data made available by Grameenphone Ltd,  
36 Bangladesh, household travel survey data provided by the Japan International Cooperation  
37 Agency (JICA), and aggregate census data obtained from the Bangladesh Bureau of Statistics  
38 (BBS). We would like to thank the Economic and Social Research Council (ESRC) of the UK,  
39 the Institute for Transport Studies, University of Leeds and FP7 Marie Curie Career Integration  
40 Grant of the European Union (PCIG14-GA-2013-631782) for funding this research. Stephane  
41 Hess was supported by the European Research Council through the consolidator grant 615596-  
42 DECISIONS.

## 43 **1 Introduction**

44 Traditional approaches to developing travel behaviour models rely on household travel surveys  
45 to establish the mathematical relationship between the choices made by the travellers, the  
46 attributes of the network and socio-demographic characteristics of the travellers. However,  
47 household surveys are often affected by low response rates and reporting errors (e.g. Rolstad  
48 et al., 2011, Groves, 2006). Further, the surveys are expensive to conduct which leads to small  
49 sample sizes and lower update frequencies. Consequently, transport models designed to fit  
50 household travel survey data alone can result in biased parameters capturing the noise in the  
51 data rather than the actual relationships in the population.

52 On the other hand, there has been growing interest in the use of mobile phone data for mobility  
53 modelling over the last few decades. Among the various transport-related applications, such  
54 data has been widely used to estimate origin-destination matrices (e.g. Çolak et al., 2015, Iqbal  
55 et al., 2014, Pan et al., 2006, White and Wells, 2002) and trip generation (e.g. Çolak et al.,  
56 2015). Since mobile phone data generally covers significant proportions of the population  
57 (GSM Association, 2017), the data is able to reliably capture the aggregate travel patterns.  
58 However, due to its anonymous nature, mobile phone data is not traditionally used in  
59 developing mathematical models of travel behaviour that establish the relationship between  
60 observed travel behaviour and causal factors such as the attributes of the alternatives and the  
61 decision makers. The existing mobility models based on mobile phone data alone cannot be  
62 used to reliably test alternative or future travel demand scenarios, and yet this is one of the core  
63 roles of transport models.

64 We are thus in a situation where traditional survey data is small in size, potentially  
65 unrepresentative and inaccurate, but contains information on key causal variables. On the other  
66 hand, mobile phone data is larger in size, more representative and accurate but missing  
67 information on key causal variables. This situation motivates the present research where we  
68 propose a framework that brings in a third type of data, namely census information, which is  
69 representative and contains detailed socio-demographic variables but does not have travel  
70 behaviour information. We thus combine household travel survey data, aggregate census data,  
71 and mobile phone data using a combination of population synthesis techniques (to generate  
72 realistic disaggregate artificial populations to assist with forecasting) and mathematical  
73 modelling to jointly optimise the aggregate and the disaggregate fit of travel behaviour models.  
74 In terms of the aggregate fit, we seek to minimise the error between the modelled and the zonal  
75 trip productions derived from call detail record (CDR) data, while in terms of the disaggregate  
76 fit, we seek to ensure that the model parameters represent the genuine sensitivities of  
77 individuals in the population. The framework is calibrated and tested in the context of trip  
78 generation models.

79 In the context of trip generation, the traditional models based on household survey data  
80 establish the mathematical relationship between the number of trips made by an individual or  
81 household with the socio-demographics (see Bwambale et al., 2015 and the cited references).  
82 But the household survey data is prone to under-reporting of the number of trips (e.g. Zhao et  
83 al. 2015, Stopher et al. 2007, Itsubo and Hato 2006). Aggregating models based only on  
84 household survey data for estimating the zonal travel patterns can lead to errors, with serious  
85 consequences for the different steps of the four-stage model. This prompts us to investigate  
86 various ways of adjusting the parameter scales of the traditional trip generation model by using  
87 a joint optimisation process to combine it with the trip patterns derived from the mobile phone  
88 data. We adopt a joint optimisation approach because CDR data too is inherently noisy, and  
89 thus not error-free. Given the lack of knowledge about which datasource really represents the  
90 ground truth, it would also be unrealistic to benchmark one dataset over the other.

91 In the proposed joint modelling framework, the base trip generation model is first estimated  
92 using household travel survey data alone to obtain the parameter priors (i.e. the sensitivities).  
93 The parameter scales are then adjusted in three different approached (without changing the  
94 prior parameter signs). The joint models hence explain the reported trips for each individual in  
95 the household survey data and ensure that the aggregated zonal trip productions are close to  
96 those derived from CDR data. This ensures that the joint models do not lose the travel  
97 behaviour sensitivities reflected in the household survey data and is computationally tractable.

98 The rest of the paper is organised as follows, section 2 presents a brief review of the literature,  
99 section 3 presents the data used in this study, section 4 presents the modelling framework,  
100 section 5 presents the model results, and section 6 presents the summary and conclusions of  
101 the study.

## 102 **2 Literature review**

103 This section presents a brief review of the literature on related work in applying mobile phone  
104 data to trip generation and other mobility studies, as well as an overview of different population  
105 synthesis techniques.

### 106 **2.1 Previous applications of mobile phone data to trip generation**

107 The estimation of trip generation from CDR data remains a challenging area of research, with  
108 only one study so far covering this subject to the best of our knowledge (Çolak et al., 2015).  
109 This is mainly due to the spatio-temporal discontinuities in the data as it only reports mobile  
110 phone positions associated with calls (voice, message, data), thereby making it difficult to  
111 capture movements when the phone is not in use. Çolak et al. (2015) attempt to address the  
112 issue of missed movements to and from the home location by introducing a home-based trip  
113 where the first or the last reported position of the day in the CDR data is at a non-home location.  
114 Although this partly addresses the problem, the challenge still remains as several other home-  
115 based trips made during the day can be missed if the mobile phone is not in use. Nonetheless,  
116 it is important to note that CDR data is likely to become more reliable in the near future with  
117 the increasing use of apps by means of mobile internet data services (Gerpott and Thomas,  
118 2014), which will increase the frequency of recorded mobile phone positions, thereby reducing  
119 the spatio-temporal discontinuities in the data. Besides CDR data, trip generation has also been  
120 previously estimated from GSM data, which is more continuous compared to CDR data (e.g.  
121 Bwambale et al., 2015). However, GSM data remains rare as it is typically not stored by mobile  
122 network operators due to storage space constraints.

### 123 **2.2 Related studies on mobile phone data and population synthesis**

124 The availability of large-scale mobile phone data over the last few decades has motivated a lot  
125 of research in quantifying human mobility and activity patterns using synthetic data generation  
126 methods (e.g. Chen et al., 2014).

127 From an epidemiology perspective, Vogel et al. (2015) combined CDR data with synthetic  
128 populations to model the spread of Ebola in West African countries and obtained promising  
129 results with respect to the Ebola predictions of the Centre for Disease Control and Prevention  
130 (CDC). Still in West Africa, Cárcamo et al. (2017) developed an intelligent epidemiology  
131 simulation software based on synthetic populations made up of agents with realistic travel  
132 behaviour derived from CDR data. In France, Panigutti et al. (2017) compared the spread of a  
133 simulated epidemic using CDR and census survey travel patterns, finding greater similarity in  
134 areas with high population and connectivity, potentially due to the higher calling rates.

135 In the field of transport, Zilske and Nagel (2014) generated artificial CDR data from synthetic  
136 passengers in a simulated traffic scenario and re-used the data to approximate the amount of

137 missed traffic at different calling rates to quantify the error introduced by CDR location  
138 discontinuities. The study found that the errors were inversely proportional to the calling rates  
139 and proposed scaling procedures based on observed data such as traffic counts. This led to a  
140 subsequent study where simulated CDR data and a synthetic population were combined with  
141 link traffic counts to generate all-day trip chains (Zilske and Nagel, 2015). This study found  
142 that even highly biased CDR data could reasonably reproduce the traffic state across different  
143 time periods. This approach of using observed traffic counts to scale CDR data has also been  
144 tested in Dhaka in the context of transient origin-destination (OD) matrix estimation (Iqbal et  
145 al., 2014).

146 Calabrese et al. (2011) developed a methodology to determine the origin-destination flows  
147 utilising 829 million mobile phone locations data for 1 million devices. Those mobile phone  
148 locations data were generated using the cell tower triangulation algorithm and have a lower  
149 resolution and higher uncertainty compared to GPS data. Data of this type was the primary  
150 source of location data for Location Based Services (LBS) before smartphones began to acquire  
151 a significant share of the mobile phone market. In the case of a smartphone, location data can  
152 also be collected through different smartphone applications that use the phone's GPS  
153 technology, WAP data, and user-provided information (Rao and Minakakis, 2003; Huang et al.,  
154 2018). Therefore, smartphone LBS data provide more details (with higher resolution, and  
155 higher frequency) footprints of the user's activities. However, the penetration rate of such  
156 application data is very low compared to CDR data. Several studies have used LBS data from  
157 different sources to implement it in transportation engineering applications. Some of the  
158 applications include travel data collection (Greaves et al., 2015; Safi et al., 2015, 2016;  
159 Patterson and Fitzsimmons, 2016; Xiao, Juan, and Zhang, 2016), activity analysis (Xiao et al.,  
160 2012; Zhou et al., 2016 ), travel behaviour analysis (Vlassenroot et al., 2015; Ferrer and Ruiz,  
161 2014; Deutsch et al., 2012 ), and travel mode detection (Zhou et al., 2016; Wu et al., 2016;  
162 Shin et al., 2015).

163 Still in the field of transport, population synthesis has been applied on real-world mobile phone  
164 datasets. Ros and Albertos (2016) updated MATSim (an agent-based multi-simulation  
165 software) by fusing census and CDR data from Spain to generate synthetic populations with  
166 mobility patterns observed in the CDR data. It may be noted that in this particular case, the  
167 mobile operator also provided the age and the gender of the users, which ensured a reliable  
168 dependence structure between the travel patterns and socio-demographics in the final synthetic  
169 population. However, mobile phone data is usually anonymous, which makes direct socio-  
170 demographic linkage impossible. In our earlier work (Bwambale et al., 2017), we developed a  
171 demographic group prediction model based on mobile phone usage behaviour extracted from  
172 CDR data (as part of a latent class model for trip generation), and can potentially be used for  
173 generating synthetic populations, however, this also requires a sub-sample of CDR data with  
174 known demographics, which is rarely available.

175 Kressner (2017) combined consumer and anonymous mobile phone data (wireless signalling  
176 and GPS data) from the United States to generate synthetic individual-level trip diaries. The  
177 socio-demographics in the disaggregate consumer data were benchmarked against the marginal  
178 census totals, while the synthetic travel was benchmarked against the mobility patterns  
179 extracted from the aggregate mobile phone data of several operators. Although this approach  
180 performed quite well in terms of aggregate-level validation, the disaggregate dependency  
181 structure between the individual's socio-demographics and trips could be seen as arbitrary.  
182 Zhang D. (2018) proposed an integrated model using Exponential Random Graph and Bayesian  
183 approaches to combine HHS and CDR data to generate a synthetic 'connected' population. The  
184 proposed model aims to reproduce the marginal and joint distributions of individuals and

185 household level socio-economic characteristics, a geographical pattern of the observed  
186 community structure, and the statistics of the observed social network.

187 To maintain the underlying dependence structure between the individual's socio-demographics  
188 and trips, Janzen et al. (2017) combined household travel survey data, register data (national  
189 statistics) and CDR data from France to correct the under-reporting of long-distance trips in  
190 travel surveys using population synthesis techniques. The socio-demographics in the travel  
191 survey data were matched against those in the register data, while the reported long-distance  
192 trips in the travel survey data were matched against those derived from the CDR data. However,  
193 a potential issue with this approach is that it assumes uniform under-reporting for all the  
194 respondents in the travel survey data, and yet this might vary, at least across different  
195 demographic groups, with some cases of over-reporting. Furthermore, the assumed higher  
196 reliability of CDR data versus travel survey data is contentious and needs to be approached  
197 impartially. This is why we propose an optimisation approach between the two datasets.

### 198 **2.3 Existing methods of population synthesis**

199 Population synthesis is widely applied in activity-based models, and various techniques have  
200 been proposed to do this. This section presents a brief review of these methods.

201 The most widely applied technique is iterative proportional fitting (IPF), which works by fitting  
202 a contingency table based on disaggregate survey data to the marginal totals in aggregate census  
203 data, constrained by a set of control variables (Beckman et al., 1996). Since its development,  
204 various improvements based on the original concept have been proposed to enhance its  
205 applicability to new challenges. These improvements have mainly focussed on addressing the  
206 zero-cell problem (Guo and Bhat, 2007), simultaneous control of household and individual-  
207 level attribute distributions (Casati et al., 2015, Zhu and Ferreira Jr, 2014, Ye et al., 2009, Guo  
208 and Bhat, 2007), improving the computational speeds (Pritchard and Miller, 2012), and non-  
209 integer conversion to integers (Choupani and Mamdoohi, 2015) etc.

210 Another popular technique is combinatorial optimisation, which focusses on selecting a subset  
211 of households in the disaggregate sample data that closely fit the marginal distributions in the  
212 census data for the same area (Voas and Williamson, 2000). This is done by randomly selecting  
213 an initial subset of households from the sample data, and iteratively replacing these with those  
214 remaining in the sample data, if and only when this leads to improvements in the fit of the  
215 subset. Although this approach has been reported to be superior (Ryan et al., 2009), the IPF  
216 method remains the most popular due to its low data requirements, reliability, and faster  
217 optimisation (Choupani and Mamdoohi, 2015, Sun and Erath, 2015).

218 Besides the two methods above, other techniques have been proposed including, the sample-  
219 free method (Barthelemy and Toint, 2013), Markov chain Monte Carlo simulation (Farooq et  
220 al., 2013), and the Bayesian network framework (Sun and Erath, 2015), among others.

## 221 **3 Data**

222 This section describes the study area, the data used, and the data processing conducted prior to  
223 model estimation. The study combines different data types (i.e. household travel survey data,  
224 census data, and CDR data) collected at different times between 2009 and 2012. Despite this  
225 limitation, these periods are considered close enough to facilitate cross-comparison.

### 226 **3.1 Data description**

#### 227 *3.1.1 Study area*

228 The study location is Dhaka Metropolitan Area (DMA) in Bangladesh. The area covers  
229 approximately 303 square kilometres and is one of the world's most crowded places with a

230 population density of 30551 persons per square kilometre (BBS, 2013). Due to the high  
 231 population density, the cell tower density is also very high. The area is served by 1361 towers,  
 232 with most these located in the central business district. The average tower-to-tower distance is  
 233 approximately 1 kilometre (Iqbal et al., 2014). The total daily trip production from DMA  
 234 residents was approximately 20.8 million in 2010, with 85.46% of these being home-based  
 235 (JICA, 2010).

### 236 3.1.2 CDR data

237 The CDR data used in this study was provided by Grameenphone Ltd and covers the working  
 238 days (i.e. Mondays to Thursdays) between 24 June 2012 and 07 July 2012 (2 weeks). The  
 239 dataset contains information from 6.9 million anonymous users representing about 57% of the  
 240 population (BBS, 2012), who together generated over 600 million records during this period  
 241 An excerpt of the randomised CDR data is presented in Table 1, where the location information  
 242 refers to tower positions as opposed to triangulated positions.

243 **Table 1: Excerpt of the CDR data (anonymised and randomised)**

Unique ID	Date	Time	Duration	Tower Longitude	Tower latitude
AAH03JACKAAAgfBALW	20120624	13:41:49	15	23.9339	90.2931
AAH03JAC8AAAAbZfAHB	20120624	13:41:25	73	23.7931	90.2603
AAH03JAC4AAAacvbABC	20120624	13:27:39	8	23.7761	90.4261
AAH03JAC9AAAAbWFAVM	20120624	13:27:27	41	23.7097	90.4036
AAH03JABkAAHvEkaQE	20120624	13:32:38	530	23.7386	90.4494

### 244 3.1.3 Household travel survey data

245 The household travel survey data used was collected between March 2009 and March 2010 as  
 246 part of the Dhaka Urban Transport Network Development Study (JICA, 2010). The sampling  
 247 of households in each zone was based on the population shares at a rate of approximately 1%.  
 248 The total sample covers 67461 individuals and 17270 households, representing an average  
 249 household size of approximately four persons. The collected information includes each  
 250 individual's socio-demographic details (e.g. gender, age, working status, income, household  
 251 size and housing type) and a single day trip diary. Table 2 presents the summary statistics of  
 252 the data.

253 **Table 2: Summary statistics of the household survey data**

Gender		Age		Working status		Trip rate shares	
Male	53%	0-9 years	15%	Employed	35%	0 trips	43%
Female	47%	10-14 years	9%	Unemployed	38%	1-2 trips	41%
		15-19 years	8%	Student	27%	3-4 trips	14%
		20-29 years	22%			5+ trips	2%
		30-49 years	32%				
		50-59 years	8%				
		60+ years	5%				

### 254 3.1.4 Census data

255 The 2011 Bangladesh Population and Housing Census data was used (BBS, 2012). The Census  
 256 was conducted from 15 to 19 March 2011. The available data reports the aggregate totals of  
 257 the selected person and household level attributes at different geographical scales (e.g. village,  
 258 ward, and zone (Thana)).

259 Since we could not access the detailed census data due to privacy reasons, we used population  
 260 synthesis techniques (Ye et al., 2009) to generate realistic artificial populations for the different

261 study area zones by combining the aggregate census data with the household survey data as  
 262 explained later in Section 3.2.2.

263 It may be noted that the fusion of household survey data and census data could only be done at  
 264 the zone (Thana) level due to differences in the study area delimitations at smaller geographical  
 265 scales. The variables available in both datasets are summarised in Table 3.

266 **Table 3: Variables in both the census and the household survey data**

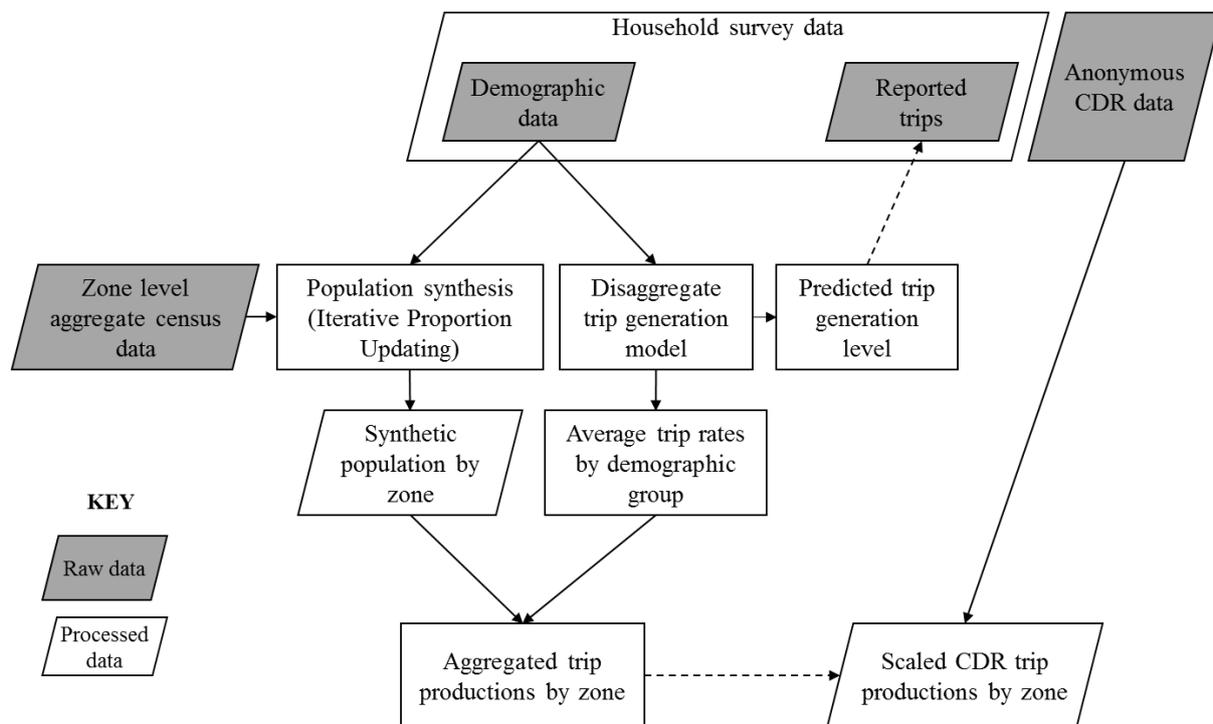
Data	Household survey data	Census data
Individual attributes	Gender	Population by gender
	Age-group	Population by age-group
	Working status (employed, unemployed, student)	Population by working status
	Occupation <sup>1</sup> (agriculture, industry, services)	Population by occupation
Household attributes	Household size	Number of households by household size
	Household type (permanent, semi-permanent, thatched etc.)	Number of households by household type

267 **3.2 Data processing and combination**

268 *3.2.1 General concept*

269 Figure 1 presents a summary of the data processing framework. The subsequent sections  
 270 discuss the key aspects of this framework.

271



272

273

**Figure 1: Data processing framework**

<sup>1</sup> Due to the differences in the definition of the Occupation categories, this data was however not usable for the synthesis.

274 The overarching idea is to minimise the difference between the zonal trip productions derived  
 275 from CDR data and those obtained by aggregating the disaggregate trip generation model,  
 276 without compromising the behavioural sensitivities reflected in the household survey data.  
 277 Model aggregation is based on a synthetic population generated using the Iterative Proportional  
 278 Updating technique (Ye et al., 2009).

### 279 3.2.2 Population synthesis

280 Among the various software applications for population synthesis, we used PopGen (Ye et al.,  
 281 2009), which is capable of conducting Iterative Proportional Updating (IPU). This algorithm  
 282 simultaneously controls for both the person and the household-level attribute distributions  
 283 during the fitting procedure, and has been proven to perform better than the simpler synthesis  
 284 methods.

285 As seen in Figure 1 (top left), the algorithm relies on two raw datasets, the household survey  
 286 data and the zone level aggregate census data to generate the zone-specific synthetic  
 287 populations by means of IPU. The household and individual level control variables used in the  
 288 IPU process are presented in Tables 4 and 5 respectively. It may be noted that we did not use  
 289 the individual's occupation as there are differences in the definitions of the categories used in  
 290 the household survey and the census data.

291 **Table 4: Household-level control variables used in PopGen**

<b>HSETYP</b>	<b>Housing type</b>	<b>HHLDSIZE</b>	<b>Household size</b>
HSETYP1	Pucka (Permanent house)	HHLDSIZE1	1
HSETYP2	Semi-pucka (Semi-permanent house)	HHLDSIZE2	2
HSETYP3	Kutchra (Thatched house)	HHLDSIZE3	3
HSETYP4	Jhupri (Slum house)	HHLDSIZE4	4
		HHLDSIZE5	5
		HHLDSIZE6	6
		HHLDSIZE7	7
		HHLDSIZE8	8+

292

293

**Table 5: Individual-level control variables used in PopGen**

<b>GEND</b>	<b>Gender</b>	<b>AGEP</b>	<b>Age-group</b>
GEND1	Male	AGEP1	0-9 years
GEND2	Female	AGEP2	10-14 years
		AGEP3	15-19 years
		AGEP4	20-29 years
		AGEP5	30-49 years
		AGEP6	50-59 years
		AGEP7	60+ years

294

295 Figure 2 presents the distribution of the Average Absolute Relative Differences (AARD)<sup>2</sup>  
 296 across the zones. This metric gives the mean deviation of the person weighted sums with

2

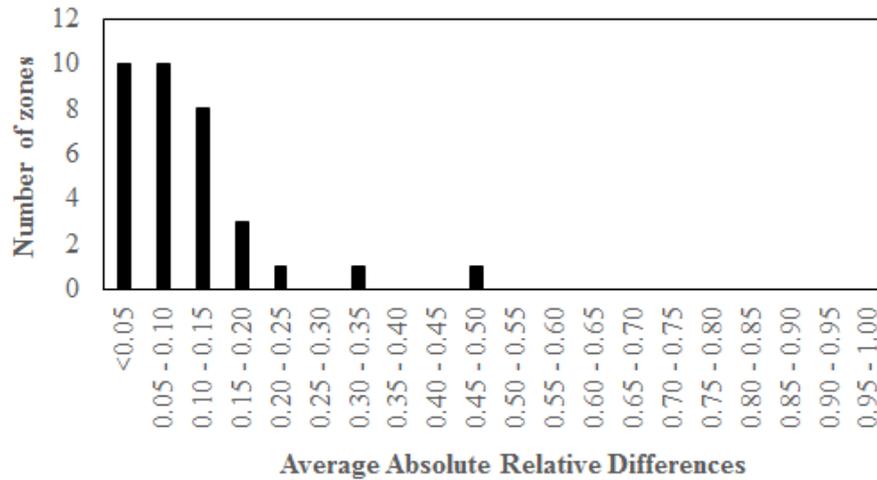
$$AARD = \frac{1}{N} \sum_{i=1}^N \frac{|w_i - c_i|}{c_i}$$

Where,  $c_i$  is the  $i^{th}$  household or person-level constraint obtained from the census data (e.g. the number of men, women, and households by household size etc.),  $w_i$  is the weighted frequency of persons with the  $i^{th}$  attribute in the generated synthetic population, and  $N$  is the total number of constraints.

297 respect to the household and person aggregate census totals (the constraints). As observed, the  
 298 AARD values for most zones are concentrated in the lower ranges of the axis, an indication  
 299 that the population synthesis was successful.

300 Furthermore, comparisons of the synthetic versus the actual estimates for each attribute at the  
 301 person and the household levels are presented in Figures 3 and 4 respectively, where the  
 302 distributions are observed to have a close match.

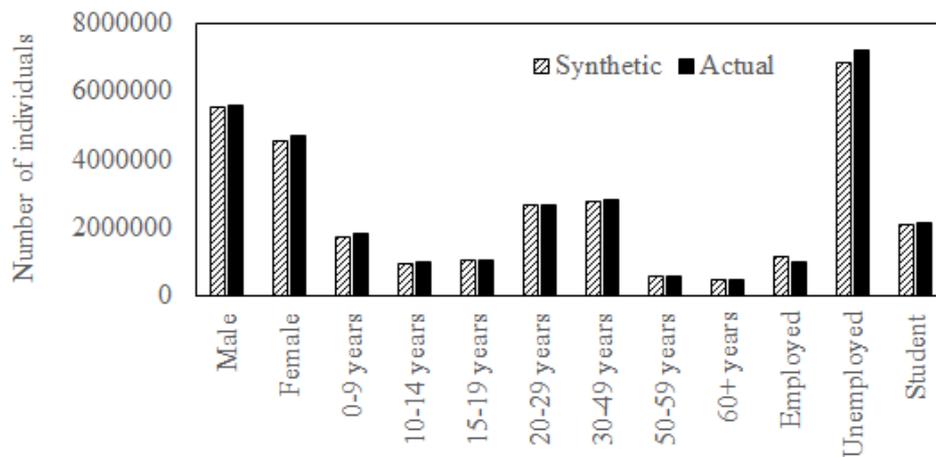
303



304

305 **Figure 2: Distribution of the AARD values**

306

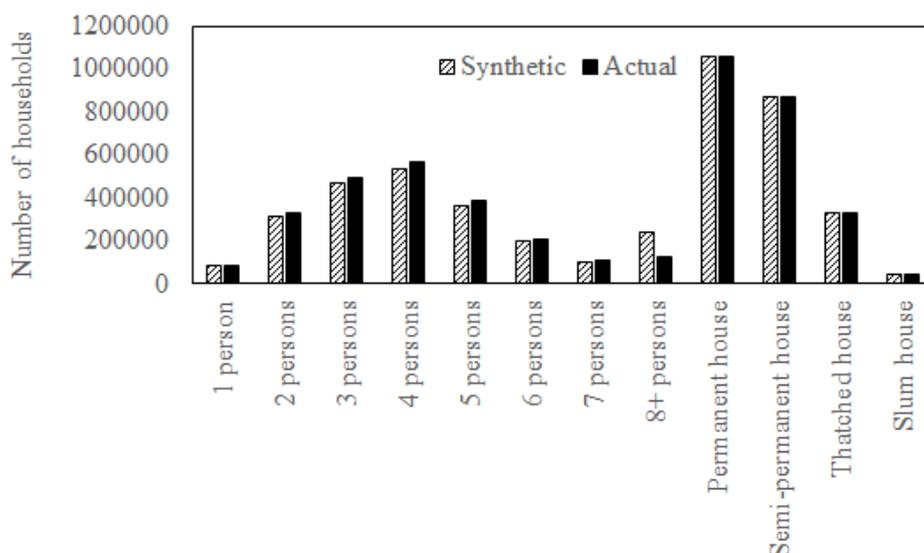


307

308

309 **Figure 3: Distribution of the individual-level estimates**

309



310  
311

**Figure 4: Distribution of the household-level estimates**

### 312 3.2.3 Extraction of unscaled zonal trip productions from CDR data

313 The CDR data for the entire observation period was first analysed to identify each user's home  
314 location, which was defined as the most frequently observed cell tower at night (i.e. between 8  
315 pm and 6 am). The labelled cell towers (i.e. home/others) for each user were then arranged  
316 according to the date and observation timestamp.

317 Home-based trips were extracted by considering any two consecutive CDR events from  
318 different cell towers, with one of those being the home cell tower. From the CDR data, we can  
319 note the distance between adjacent towers varies between 0.02 and 7.00 kilometres. Most areas  
320 of Dhaka are densely populated and about 75% of the towers have an adjacent distance of less  
321 than 0.5 kilometres (90% have an adjacent distance of less than 1 kilometre). Furthermore, a  
322 previous study in Dhaka found that the mean walking trip distance is about 0.45 kilometres  
323 (JICA, 2010). Therefore, a lower distance threshold of 0.5 kilometres between subsequent  
324 towers was considered as the optimum for minimising the number of very short trips within  
325 the neighbourhood and false trips due to tower jumps<sup>3</sup>.

326 An upper threshold of 24 hours or midnight (whichever came first) was specified based on the  
327 assumption that a user typically travels from and back to home within the same effective day.  
328 Consequently, if the first and the last CDR events for the day were not at the home cell tower,  
329 corresponding raw trips were added (Çolak et al., 2015). This led to the unscaled zonal trip  
330 productions shown in Figure 1.

### 331 3.2.4 Scaling the CDR trip productions

332 The home cell towers derived from the CDR data were mapped to the zones with the aid of a  
333 GIS software (QGIS Development Team, 2018). The total trips for each zone were then scaled  
334 using the ratio of the zonal population (from the census) to the number of users classified as  
335 residents of the zone from the CDR data (see Çolak et al., 2015 for details). We however  
336 acknowledge that this straight scaling procedure may bias the results if the CDR data sample  
337 is biased, for example in terms of the socio-economic status of the mobile phone owners.

<sup>3</sup> A false trip occurs when the user is not making a trip but there is a change in the tower as the operator reassigns the call to a different tower (due to load management purposes).

## 338 **4 Modelling framework**

339 We propose an approach that combines two modelling strategies, that is, discrete choice  
340 modelling at the individual level and ordinary least squares at the aggregate level (shown in  
341 patterned text boxes in Figure 1).

### 342 **4.1 Disaggregate trip generation model (Base model)**

343 Trip generation have been found to be affected by household characteristics (e.g. household  
344 size, income, car-ownership, etc.) and composition (e.g. numbers of children, employed people,  
345 etc.) (see Bwamable et al. 2015 and Bwamable et al. 2018 for details). Discrete choice models  
346 have been the most preferred approach for modelling trip generation over the last few decades  
347 (e.g. Bwambale et al., 2015, Pettersson and Schmöcker, 2010, Agyemang-Duah and Hall,  
348 1997). Although the ordered response choice mechanism has been the most preferred approach  
349 for modelling trip generation, the method was intractable in this particular study where model  
350 performance is being optimised at both the aggregate and disaggregate levels through scaling  
351 as discussed later in this paper. While less appealing from a theoretical point of view, the  
352 unordered response choice mechanism was found to be a more feasible approach and was  
353 adopted. It is important to note that the unordered response choice mechanism has been found  
354 to give intuitive results even in contexts with ordered choices such as car ownership (Bhat and  
355 Pulugurta, 1998).

356 To implement the unordered response choice mechanism, we rely on the random utility theory  
357 (Marschak, 1960). Let  $U_{nt}$  be the utility of individual  $n$  making  $t$  trips. This can be expressed  
358 as;

$$U_{nt} = \beta'_t X_n + \varepsilon_{nt} \quad (1)$$

359  
360 Where  $X_n$  is a vector of the socio-demographic attributes of individual  $n$ ,  $\beta_t$  is a vector of the  
361 model parameters to be estimated, and  $\varepsilon_{nt}$  is the random component of utility. Since the  
362 individual socio-demographics are constant across the alternatives, we specify a different set  
363 of parameters for each trip generation level to reflect the fact that each attribute has a  
364 differential impact on the utility for each trip generation level.

365 Under the assumption that the error terms ( $\varepsilon_{nt}$ ) are distributed independently and identically  
366 across alternatives and individuals using a type I extreme value distribution, the trip generation  
367 choice probabilities can be calculated using the multinomial logit (MNL) model (McFadden,  
368 1974) as expressed below;

$$P_{nt} = \frac{\exp(\beta'_t X_n)}{\sum_{t^*} \exp(\beta'_{t^*} X_n)} \quad (2)$$

369  
370 Where  $P_{nt}$  is the probability of individual  $n$  making  $t$  trips.

371 Despite the requirements of the MNL model, it may be noted that the error terms are not likely  
372 to be independent in the real world.

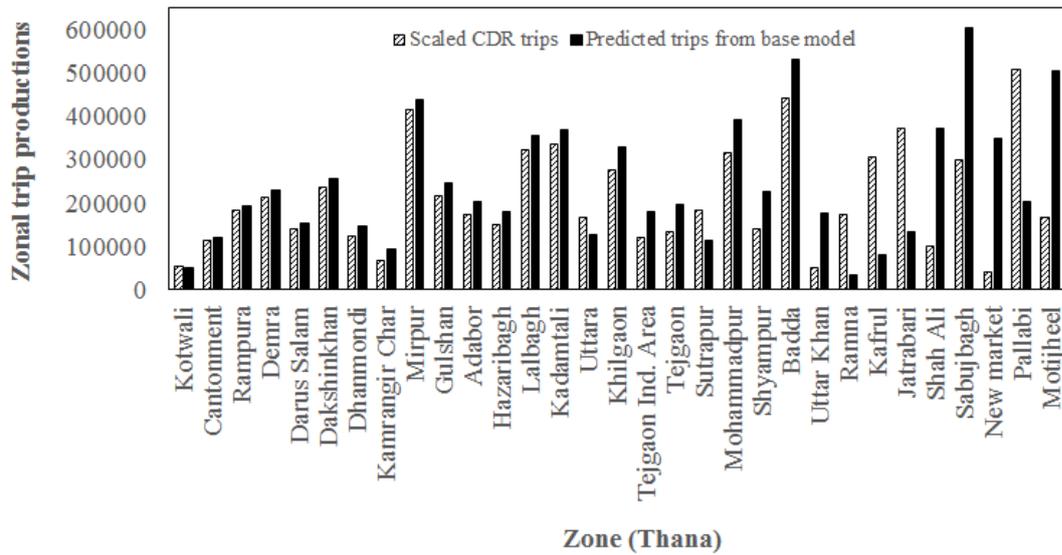
373 If we were to rely on the household travel survey data alone, the model parameters would be  
374 estimated by maximising the log-likelihood function below.

$$LL(\beta_t) = \sum_n \sum_t K_{nt} \ln(P_{nt}) \quad (3)$$

375

376 Where the dummy variable  $K_{nt} = 1$  if and only if individual  $n$  makes  $t$  trips, otherwise  $K_{nt} =$   
 377 0.

378 However as mentioned earlier, fitting the model to match the trips reported in the household  
 379 travel survey data alone can lead to biased parameter estimates due to reporting errors, thereby  
 380 resulting in misrepresentation of the aggregate travel demand as reflected in Figure 5, where  
 381 the predicted aggregate zonal trips from the base model are different from those derived from  
 382 the CDR data, especially towards the right hand side of the figure.



383

384 **Figure 5: Distribution of the CDR trip productions**

385 The relative absolute errors derived from Figure 5 were plotted on a map to check whether  
 386 there is a spatial correlation to the errors as shown in Figure 6.

387 From Figure 6, it is observed that there is no obvious spatial correlation to the errors. The  
 388 magnitude of the error is largest in a single central zone. But apart from that, larger magnitudes  
 389 are observed both in the centre of the metropolitan area, as well as, in some outskirts areas. For  
 390 the centre, the errors are most likely caused by the relatively high number of either false trips  
 391 in the CDR data (due to the high tower density) or unreported short walking trips in the  
 392 household survey data, while for the outskirts, the errors are most likely caused by the missed  
 393 short trips that could not be captured by the CDR data due to the low tower density in those  
 394 areas.

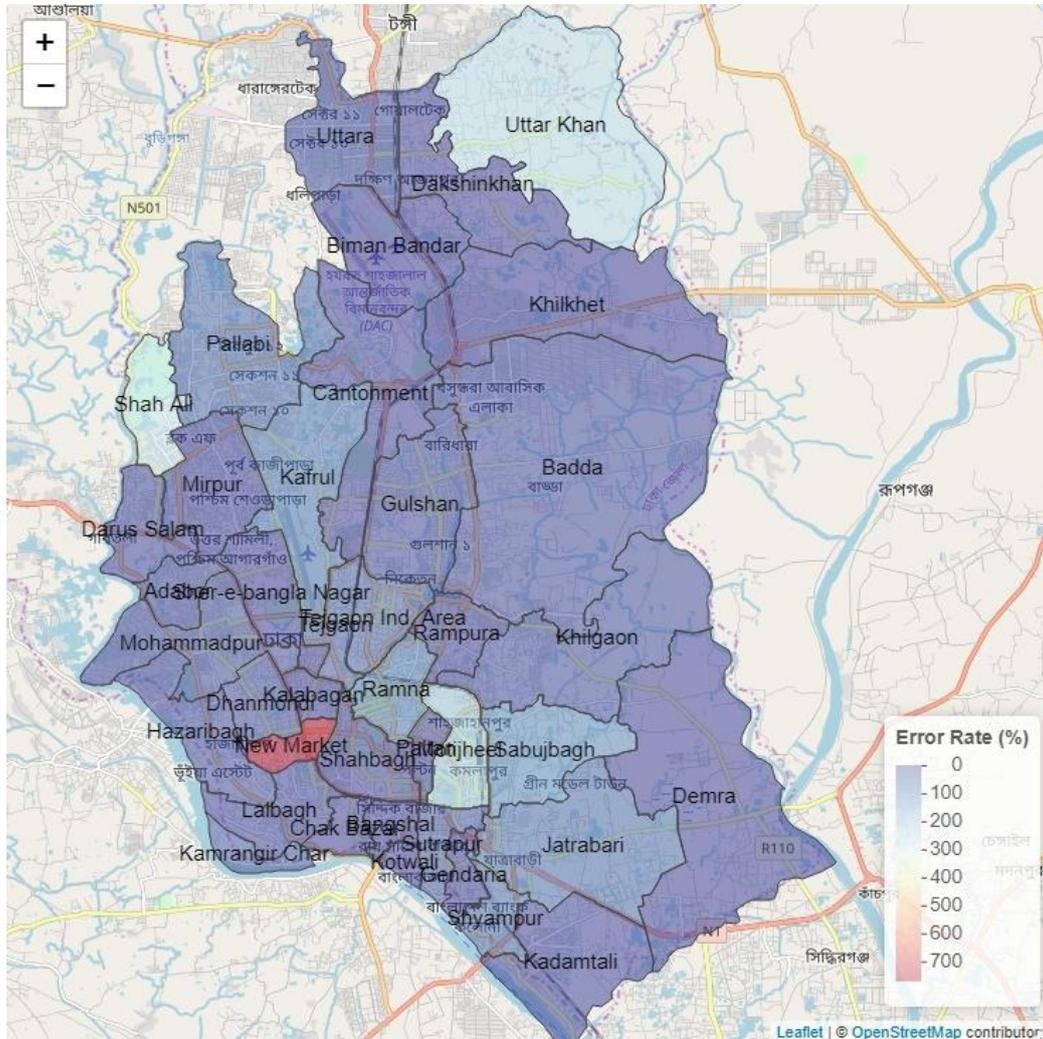
#### 395 **4.2 Joint trip generation model**

396 The priors of the parameter signs and relative magnitudes are obtained from the pre-estimated  
 397 base model. The parameter scales are then adjusted (without changing the prior parameter  
 398 signs). The joint model thus simultaneously optimises performance at both the aggregate and  
 399 disaggregate levels with respect to the CDR and the household travel survey data, respectively.

400 As mentioned earlier, this combined approach ensures that the resulting model does not lose  
 401 the travel behaviour sensitivities reflected in the household travel survey data, by maintaining  
 402 the sensitivities from the base model. Adjusting the parameter scales has an impact on the  
 403 choice probabilities for each trip generation outcome, which influences the expected trip rates  
 404 of the individuals. The framework of the joint trip generation model is described below. Let  
 405  $\hat{U}_{nt}$  be the updated utility of individual  $n$  making  $t$  trips. This can be expressed as;

$$\hat{U}_{nt} = \alpha\beta'_t X_n + \varepsilon_{nt} \quad (4)$$

406  
407



408

409

Figure 6: Spatial distribution of errors in trip productions (CDR data versus base model)

410

411 Where  $\alpha$  is a vector of the scaling factors to be estimated. The  $\beta$  parameters are priors derived  
412 from the base model, and are not re-estimated in the joint framework. The specification of the  
413 scaling factors is discussed later on.

414 The updated trip generation choice probability can be expressed as follows;

$$\hat{P}_{nt} = \frac{\exp(\alpha\beta'_t X_n)}{\sum_{t^*} \exp(\alpha\beta'_{t^*} X_n)} \quad (5)$$

415

416 Where  $\hat{P}_{nt}$  is the updated probability of making  $t$  trips by individual  $n$ .

417 However, to estimate the scaling factors, we need to fulfil two objectives. The first objective is  
418 to explain the reported trips for each individual in the household survey data. The second  
419 objective is to ensure that the aggregated zonal trip productions are close to those derived from

420 CDR data. Both outcomes have a probability attached to them and the simultaneous estimation  
421 maximises the joint probability of the two outcomes.

422 To estimate the aggregate zonal trip productions, we rely on the synthetic population generated  
423 in section 3.2.2. As mentioned earlier, the synthetic population was designed to match both the  
424 person and the household-level attribute distributions during the fitting procedure, thus making  
425 it more reliable. We have a synthetic population of  $M$  simulated individuals identified as  $m$   
426 with  $m = 1, \dots, M$ , and a study area made up of  $Z$  zones identified as  $z$  with  $z = 1, \dots, Z$ .  
427 Let  $\hat{P}_{mt}$  denote the updated probability of making  $t$  trips by simulated individual  $m$ . It may be  
428 noted that  $\hat{P}_{mt}$  is equivalent to  $\hat{P}_{nt}$  if both the simulated individual and the actual respondent  
429 in the household survey data have the same demographics (i.e. the values of  $\hat{P}_{mt}$  depend on the  
430 calculations of  $\hat{P}_{nt}$ ). Now, let  $\hat{T}_z$  denote the aggregate zonal trip production for zone  $z$ . This  
431 can be calculated by taking the weighted average trips for each simulated individual, in which  
432 the updated MNL probabilities are the weights, and summing across the zonal synthetic  
433 population as follows;

$$\hat{T}_z = \sum_{m=1}^M \left[ Y_{mz} \left( \sum_{t=1}^T (t * \hat{P}_{mt}) \right) \right] \quad (6)$$

434  
435 Where the dummy variable  $Y_{mz} = 1$  if and only if simulated individual  $m$  belongs to zone  $z$ ,  
436 otherwise,  $Y_{mz} = 0$ . The objective is to ensure that  $\hat{T}_z$  is as close as possible to the corrected  
437 CDR trip productions for zone  $z$ . If  $\varphi_z$  denotes the corrected CDR trip productions for zone  $z$ ,  
438 the relationship between  $\varphi_z$  and  $\hat{T}_z$  can be expressed as follows;

$$\varphi_z = \hat{T}_z + \omega_z \quad (7)$$

439  
440 Where  $\omega_z$  is an error term which we assume follows a normal distribution with a mean of zero,  
441  $\omega_z \sim N(0, \sigma^2)$ <sup>4</sup>.  $P(\varphi_z)$  is then the likelihood of observing the CDR trip productions for zone  
442  $z$ , and, from Equation 7, this can be expressed as follows;

$$P(\varphi_z) = \frac{1}{\sqrt{2\pi\sigma^2}} \exp\left(\frac{-(\varphi_z - \hat{T}_z)^2}{2\sigma^2}\right) \quad (8)$$

443  
444  $P(\varphi_z)$  clearly depends on  $\hat{P}_{nt}$  given that  $\hat{T}_z$  is a function of  $\hat{P}_{mt}$ , which depends on the  
445 calculations of  $\hat{P}_{nt}$  as explained earlier. For each survey respondent in zone  $z$ , we need to  
446 maximise the probability of the chosen alternative and ensure that the probabilities of all the  
447 alternatives maximise  $P(\varphi_z)$ . Let  $t_n^o$  denote the number of trips observed for individual  $n$  in  
448 the household survey data, such that  $\hat{P}_{nt^o}$  gives the logit probability of the observed choice for  
449 individual  $n$ . The overall joint likelihood ( $L$ ) of the observed choices and the aggregate CDR  
450 trip productions across individuals is calculated as follows;

$$L = \prod_{n=1}^N \left[ \sum_{z=1}^Z H_{nz} (\hat{P}_{nt^o} * P(\varphi_z)) \right] \quad (9)$$

---

<sup>4</sup> The assumption of normality is based on its widespread use in the choice modelling literature in representing error terms (owing to the computational feasibility), though other distributions may be applicable as well.

$$= \left( \frac{1}{\sqrt{2\pi\sigma^2}} \right)^N \prod_{n=1}^N \left[ \sum_{z=1}^Z H_{nz} \left( \frac{\exp(\alpha\beta'_{t^o} X_n)}{\sum_{t^*} \exp(\alpha\beta'_{t^*} X_n)} * \exp \left( \frac{-(\varphi_z - \hat{T}_z)^2}{2\sigma^2} \right) \right) \right]$$

452

453 Where the dummy variable  $H_{nz} = 1$  if and only if survey respondent  $n$  belongs to zone  $z$ .

454 This is based on the assumption that  $\hat{P}_{nt}$  and  $P(\varphi_z)$  are independent. This is not unreasonable  
 455 given the sources of potential errors are very different (reporting errors in case of the HHS and  
 456 coarse resolution in case of the CDR) and there is no obvious source of correlation among the  
 457 two probabilities. Since products are difficult to differentiate, we obtain the log-likelihood ( $LL$ )  
 458 by applying logarithms to Equation 9 resulting in Equation 10.

459

$$LL = -\frac{N}{2} \log(2\pi) - N \log(\sigma) + \sum_{n=1}^N \sum_{z=1}^Z H_{nz} \left( \ln \left[ \frac{\exp(\alpha\beta'_{t^o} X_n)}{\sum_{t^*} \exp(\alpha\beta'_{t^*} X_n)} \right] - \frac{1}{2\sigma^2} (\varphi_z - \hat{T}_z)^2 \right) \quad (10)$$

460

461 Three parameter scaling scenarios are tested, and these are;

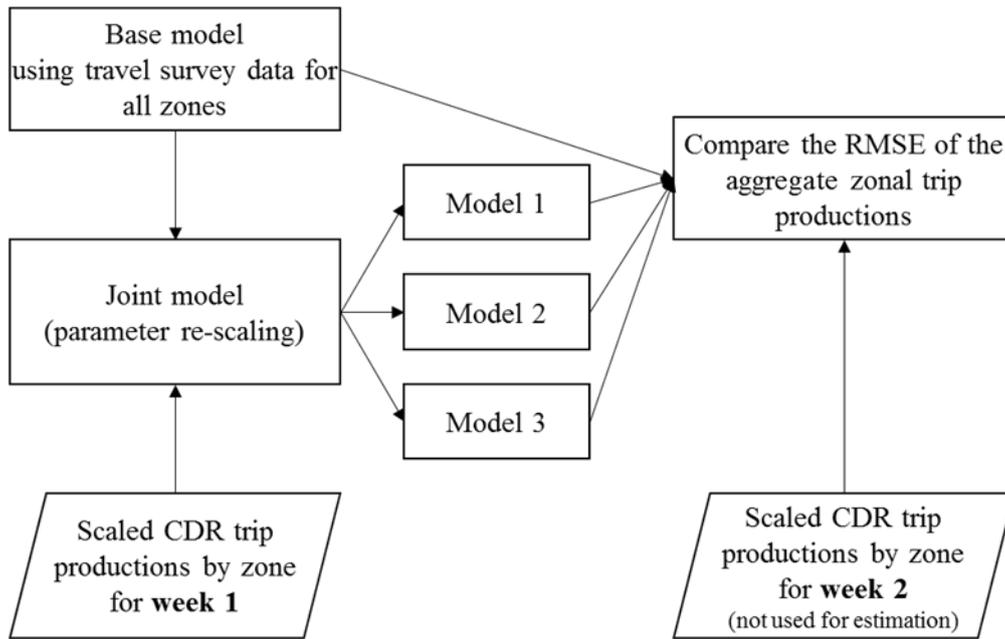
- Model 1 This specification applies the same  $\alpha$  scaling factor to the utility models of the different trip generation levels (see Equation 4), i.e.  $\alpha_t = \alpha, \forall t$ . The updated utility models have the same relative variable sensitivities as in the base model, albeit with different parameter scales.
- Model 2 This specification applies a different  $\alpha_t$  scaling factor to the utility model of each trip generation level. The updated utility models maintain the base model relative variable sensitivities for each particular trip generation level, however, the variable sensitivities across the different trip generation levels are adjusted with different parameter scales, and hence the relative values across levels change from the base model.
- Model 3 This specification applies a different  $\alpha_x$  scaling factor to each explanatory variable  $X$  (e.g. gender, age-group, and working status), however,  $\alpha_x$  does not change across the different trip generation levels. The updated utility models maintain the base model attribute-level relative sensitivities for a particular variable across the different trip generation levels, however, the inter-variable relative sensitivities are adjusted with different parameter scales.

462 **4.3 Model evaluation framework**

463 The performance of the joint models is evaluated in terms of both the temporal and the spatial  
 464 transferability as presented in Figures 6 and 7, respectively.

465 In terms of temporal transferability, the joint models associated with each parameter scaling  
 466 scenario are estimated using the zonal aggregate CDR trip productions for week 1. The  
 467 prediction capacities of the estimated joint models, as well as the base model are then compared  
 468 in terms of the root mean square errors with respect to the zonal aggregate CDR trip productions  
 469 for week 2 (see Figure 7).

470 In terms of spatial transferability, the study area zones are randomly divided into two groups.  
 471 The base and the joint models are then estimated using the data for one group of zones and  
 472 applied to the other group of zones (not used for estimation). The prediction capacities of the  
 473 models are then compared in terms of the predictive joint log-likelihoods, and the root mean  
 474 square errors with respect to the aggregate CDR trip productions of the application zones (see  
 475 Figure 8).



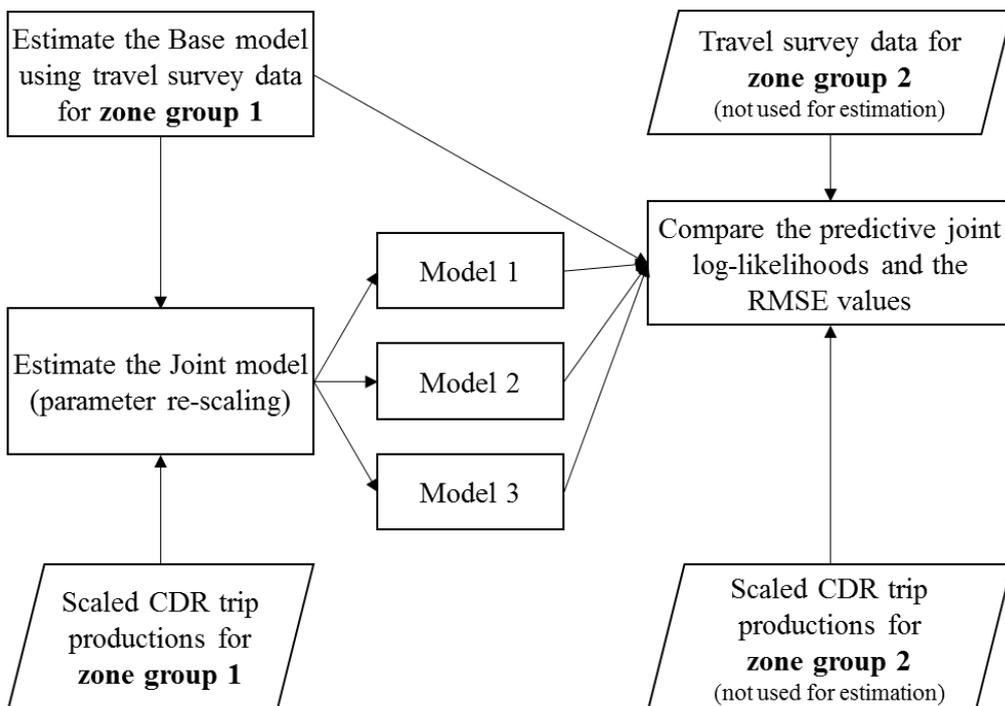
476

477

Figure 7: Temporal transferability framework

478

479



480

481

Figure 8: Spatial transferability framework

## 482 **5 Modelling results**

483 This section presents the final model specification, as well as the model estimation and  
484 validation results.

### 485 **5.1 Variable specification**

486 The dependent variable is the number of individual home-based trips (irrespective of the trip  
487 purpose). This is because we could not reliably infer the purposes of the CDR trips. Based on  
488 distributions in the data, the trip generation levels were grouped into 0, 1-2, 3-4, and 5+ trips  
489 per day. The explanatory variables considered for possible inclusion in the model are those that  
490 were used for population synthesis. The household-level variables (i.e. household size and  
491 type) were however not included in the final model as they led to unreasonable parameter signs,  
492 potentially due to their weak influence on individual trip-making decisions<sup>5</sup>. The final model  
493 specification thus contains the gender, the age-group, and the working status of the individuals,  
494 coded as dummy variables.

495 For model identification purposes, the parameters associated with the zero trip generation level  
496 were treated as the base (for all explanatory variables). Furthermore, male non-workers in the  
497 30-49 age-group were treated as the base demographic group, and their preferences are entirely  
498 explained by the alternative specific constants. Thus, the model parameter estimates represent  
499 the differential impact on utility with respect to the zero trip generation level and the base  
500 demographic group.

### 501 **5.2 Estimation results**

#### 502 *5.2.1 Base model*

503 We first estimated the base model to assess whether the parameter estimates are in line with  
504 the expected travel behaviour. The model results are presented in Table 6.

505

506

**Table 6: Base model results**

<b>Variable</b>	<b>Parameter</b>	<b>t-statistic</b>
<b>Alternative specific constants (ASCs)</b>		
1-2 trips	-0.2069	-7.46
3-4 trips	-1.0408	-24.56
5+ trips	-3.0859	-31.19
<b>Dummies specific to gender (base category is males)</b>		
<i>Females</i>		
1-2 trips	0.0870	3.94
3-4 trips	-0.2841	-7.95
5+ trips	-0.2654	-3.15
<b>Dummies specific to working-status (base category is non-workers)</b>		
<i>Workers</i>		
1-2 trips	0.4630	17.23
3-4 trips	0.9252	23.05

<sup>5</sup> The larger household sizes in Dhaka can often be attributed to the number of support staff members (e.g. cooks, cleaners, gardeners, housekeepers etc.) who stay and work full-time in the household. This is a potential contributing factor to the weak correlation between the numbers of people in a household and trip generation, which we appreciate is different in a more European/North American context.

5+ trips

1.1482

12.38

507

Table 6 cont'd

Variable	Parameter	t-statistic
<i>Students</i>		
1-2 trips	1.4079	46.47
3-4 trips	0.9381	17.13
5+ trips	-0.5333	-2.65
<b>Dummies specific to age-group (base category is the 30-49 years age-group)</b>		
<i>Age 1-9 years</i>		
1-2 trips	-1.6354	-50.69
3-4 trips	-3.1065	-36.73
5+ trips	-3.5549	-9.46
<i>Age 10-14 years</i>		
1-2 trips	-0.8143	-19.49
3-4 trips	-1.7635	-22.52
5+ trips	-1.9201	-6.00
<i>Age 15-19 years</i>		
1-2 trips	-0.6539	-16.22
3-4 trips	-0.9669	-15.71
5+ trips	-1.0077	-5.71
<i>Age 20-29 years</i>		
1-2 trips	-0.1457	-5.67
3-4 trips	-0.3249	-9.58
5+ trips	-0.3009	-4.02
<i>Age 50-59 years</i>		
1-2 trips	-0.1423	-4.12
3-4 trips	-0.2552	-5.92
5+ trips	-0.3721	-3.81
<i>Age 60+ years</i>		
1-2 trips	-0.2494	-5.63
3-4 trips	-0.3531	-6.14
5+ trips	-0.4853	-3.47
<b>Measures of fit</b>		
Number of observations	65419	
Log-likelihood at zero	-90689.99	
Log-likelihood at convergence	-64859.90	
Number of parameters	30	
Adjusted rho-square	0.2845	
Likelihood ratio	51660.10	
P value of the likelihood ratio	0.0000	

508

509 The alternative specific constants capture the underlying differential impact on utility with  
510 respect to the zero trip generation level. All the estimates are negative, and their magnitude  
511 increases with respect to the trip generation level. Keeping all other factors constant, this  
512 reflects a general tendency to make fewer trips, especially by the base category (i.e. male, non-  
513 workers, aged 30-49 years).

514 The parameter estimates for females represent the differential impact on utility with respect to  
515 males. For 1-2 trips, we obtain a positive parameter estimate, while for the higher trip  
516 generation levels, we obtain negative parameter estimates. The proportion of women working  
517 in the garments industry, one of the leading sectors in Dhaka, is 64-90% (ADB and ILO, 2016).  
518 This probably explains the positive parameter sign for 1-2 trips. Otherwise, males are more  
519 likely to make a higher number of trips compared to females, probably due to the average  
520 higher income levels of the former (BBS, 2012) and socio-cultural factors.

521 The parameter estimates for the working status variables (i.e. workers and students) represent  
522 the differential impact on utility with respect to non-workers. As observed, the parameters for  
523 workers are positive, and their magnitudes increase with respect to the trip generation level, an  
524 indication that workers generally make more trips compared to non-workers. On the other hand,  
525 the parameter estimates for students are positive for 1-2 and 3-4 trips, and negative for 5+ trips.  
526 This shows that students make more trips compared to non-workers only up to a reasonable  
527 level expected for school going individuals.

528 Similarly, the parameter estimates for the age-group variables represent the differential impact  
529 on utility with respect to the 30-49 years age-group. As observed, the parameter estimates for  
530 all the other age-groups are negative, an indication that they generally make fewer trips  
531 compared to the base age-group (30-49 years). The active working age of white-collar workers  
532 in Bangladesh typically ranges between 29 and 60 years (i.e. the latest age for completing  
533 tertiary education and the retirement age respectively (BBS, 2012)). It is therefore reasonable  
534 that persons in the 30-49 years age-group are more active travellers due to their economic  
535 vibrancy.

536 Finally, it is observed that the overall model (in terms of the likelihood ratio), as well as all the  
537 parameter estimates (in terms of the t-statistics) are statistically significant at the 99% level of  
538 confidence (see Ben-Akiva and Lerman, 1985 for details).

### 539 *5.2.2 Joint models*

540 As earlier mentioned, the parameters of the base model were fixed in the joint modelling  
541 framework, and only the scaling factors were estimated. Table 7 presents the estimated scaling  
542 factors and the measures of fit for all the three models for comparison purposes. Positive scaling  
543 factors were obtained for all the three models, an indication that the resultant coefficients in  
544 the scaled joint models have the same signs as those in the base model.

545 A comparison of the joint convergence log-likelihoods shows that Model 3 gives the best  
546 performance, followed by Model 2, and then Model 1. This is attributed to the flexibility of the  
547 parameter scaling framework. An important point to note is that all the three joint models  
548 perform better than the base model in terms of the joint log-likelihood.

549 As earlier mentioned, during model optimisation, we are basically dealing with a trade-off  
550 between disaggregate and aggregate model performance. Thus, the disaggregate log-likelihood  
551 of the joint models is a little worse than that of the base model. However, if the base model  
552 parameters are directly used to estimate the joint log-likelihood, it is observed that the model  
553 yields the worst performance.

554 The p-values of the likelihood ratios of the joint models with respect to the base model are all  
555 less than 0.01, an indication that the improvements in performance are statistically significant  
556 at the 99% confidence level beyond the advantages offered by the additional parameters (see  
557 Ben-Akiva and Lerman, 1985 for details).

558

559

**Table 7: Joint model scaling factors**

Description of scaling factor	Model 1		Model 2		Model 3	
	Estimate	t-stat	Estimate	t-stat	Estimate	t-stat
<b>Model 1</b>						
Uniform factor (applied to all the base model parameters)	1.3650	2280.16				
<b>Model 2</b> (Factors specific to trip generation level)						
1-2 trips			1.2716	131.39		
3-4 trips			1.4873	247.83		
5+ trips			1.1699	158.63		
<b>Model 3</b> (Factors specific to particular variables)						
Gender					1.5228	33.81
Working status					1.8148	105.16
Age-group					1.3262	120.70
ASCs					1.6023	171.51
<b>Measures of fit</b>						
Convergence LL at the disaggregate level	-66002.75		-65914.01		-67747.10	
Convergence LL at the aggregate level	-718560.40		-718377.10		-715805.30	
Joint convergence LL	-784563.20		-784291.20		-783552.40	
Base model convergence LL	-64859.90		-64859.90		-64859.90	
Base model LL at the aggregate level	-805093.10		-805093.10		-805093.10	
Base model joint convergence LL	-869953.00		-869953.00		-869953.00	
Likelihood ratio (joint model w.r.t the base model)	170780		171234		172801	
P value	0.0000		0.0000		0.0000	

560

### 561 5.3 Model evaluation in terms of transferability

562 The models based on the full sample have been presented in the previous section. To evaluate  
563 the stability and the predictive performance of the joint models as well as the base model, we  
564 compared their temporal and spatial transferability following the evaluation framework  
565 described in Section 4.3. Tables 8 and 9 present the measures of fit in terms of the temporal  
566 and the spatial transferability, respectively.

567

**Table 8: Temporal transferability**

	Measure	Base model	Model 1	Model 2	Model 3
<b>Week 1</b> (Estimation)	LL (disaggregate level)	-64859.90	-66024.40	-65940.80	-67850.40
	LL(aggregate level)	-805642.50	-719566.80	-719396.20	-716695.30
	Joint LL	-870502.40	-785591.20	-785337.00	-784545.70
<b>Week 2</b> (Application)	LL (disaggregate level)	-64859.90	-66024.40	-65940.80	-67850.40
	LL(aggregate level)	-804545.50	-717793.90	-717596.20	-715031.60
	Joint LL	-869405.40	-783818.30	-783537.00	-782882.00
	RMSE w.r.t CDR trips	43342.84	13547.09	13527.84	13328.49

568

569

**Table 9: Spatial transferability**

	Measure	Base model	Model 1	Model 2	Model 3
<b>Zone group 1</b> (Estimation)	LL (disaggregate level)	-26102.10	-26712.45	-26652.76	-27724.63
	LL(aggregate level)	-321381.60	-290869.40	-290725.20	-288898.10
	Joint LL	-347483.70	-317581.85	-317377.96	-316622.73
<b>Zone group 2</b> (Application)	LL (disaggregate level)	-38859.38	-39701.58	-39352.09	-41303.51
	LL(aggregate level)	-491580.30	-429017.00	-428604.80	-426638.20
	Joint LL	-530439.68	-468718.58	-467956.89	-467941.71
	RMSE w.r.t CDR trips	50626.73	13375.06	13274.68	13161.58
<b>Zone group 2</b> (Estimation)	LL (disaggregate level)	-38688.76	-39227.43	-39333.92	-40185.59
	LL(aggregate level)	-482400.40	-428113.30	-427818.70	-426238.10
	Joint LL	-521089.16	-467340.73	-467152.62	-466423.69
<b>Zone group 1</b> (Application)	LL (disaggregate level)	-26219.53	-26689.06	-26786.11	-27445.95
	LL(aggregate level)	-315772.10	-289862.10	-289890.20	-288799.10
	Joint LL	-341991.63	-316551.16	-316676.31	-316245.05
	RMSE w.r.t CDR trips	38776.13	13702.57	13758.49	13602.58

570

571 From Table 8, it is observed that the temporal transferability of the joint models is generally  
572 higher than that of the base model in terms of the joint log-likelihoods and the root mean square  
573 errors (RMSE) with respect to the zonal CDR trips. Among the three joint models, Model 3  
574 offers the best transferability, however, Model 2 gives the best prediction at the disaggregate  
575 level in both the estimation and the application contexts.

576 For spatial transferability, we tested both directions of model transfer. It may be noted that the  
577 general interpretation of the base model parameters for each group of zones did not change.  
578 From Table 9, it is again observed that the joint models are generally more transferrable

579 compared to the base model in terms of the joint log-likelihoods and the root mean square  
580 errors for both directions.

581 In this particular case, it is observed that Model 2 gave the best disaggregate prediction for the  
582 zone group 1 to 2 transfer direction, while Model 1 gave the best disaggregate prediction for  
583 the reverse transfer direction.

584 An important point worth mentioning is that the superior performance of the base model at the  
585 disaggregate level is expected as it was designed to fit the travel survey data alone, but as  
586 mentioned earlier, this could be prone to reporting errors and hence less dependable.

587 From the results, it is clear that Model 3 gives the best overall spatial and temporal  
588 transferability, however, the disaggregate performance of Models 1 and 2 as highlighted above  
589 shows that these parameter scaling approaches offer some benefits as well. These results  
590 present initial efforts to exploit the benefits of both household travel survey and mobile phone  
591 data to optimise the performance of travel behaviour models, and there is a need for further  
592 research using data from different contexts to investigate the different parameter scaling  
593 approaches in further detail.

#### 594 **5.4 Model comparison in forecasting**

595 To test the sensitivity of the models to forecasting, the base model and the different joint  
596 models have been applied to the 2019 household survey data and the predictive measures of  
597 fit for the different models have been compared. The following three performance indicators  
598 have been used in this regard:

- 599 - Root Mean Square Error (RMSE), which has been obtained by comparing the  
600 modelled and the actual total trip productions associated with the 2019 sample data  
601 for each TAZ using the base and joint model parameters (pre-estimated using the  
602 2010 data).
- 603 - Average probability of correct prediction, which has been obtained by computing the  
604 mean probability of success for the 2019 sample data using the pre-estimated base and  
605 joint model parameters (pre-estimated using the 2010 data).
- 606 - The predictive adjusted-rho square, which has been obtained using the adjusted rho-  
607 square equation below for the pre-estimated base and the joint models;

$$\rho_{adj}^2 = 1 - \frac{LL(F) - k}{LL(0)} \quad (11)$$

608 Where;  $k$  is the number of model parameters,  $LL(F)$  and  $LL(0)$  are the values of the log-likelihood  
609 function at convergence and at zero respectively.

610  
611 Table 11 summarises the calculated predictive measures of fit on the 2019 forecasting sample for the  
612 base model and the different joint models.

613  
614

**Table 11: Predictive measure of fit on the 2019 forecasting sample**

Measure	Base model	Model 1	Model 2	Model 3
Root Mean Square Error (RMSE)	228.6346	218.5843	218.5505	214.0239
Average probability of correct prediction	0.4269	0.4553	0.4537	0.4679
Predictive adjusted rho-square	0.3548	0.3836	0.3810	0.3806

615

616 From Table 11, it is observed that overall the joint models generally perform better than the  
617 base model in forecasting at both the aggregate and disaggregate levels. Among the three joint  
618 models, it is observed that Model 3 gives the best performance in terms of both the Root Mean  
619 Square Error and the average probability of correct prediction, while giving the least  
620 performance in terms of the predictive adjusted rho-square. However, from a forecasting point  
621 of view, aggregate performance is more critical, and Model 3 would offer more benefits.

## 622 **6 Summary and conclusions**

623 This paper started by highlighting the reporting errors and sampling bias associated with  
624 household travel survey data, and how these could lead to biased model parameters (e.g.  
625 Rolstad et al., 2011, Groves, 2006). The paper outlines the possible consequences of such issues  
626 in the context of trip generation, where the estimated models would misrepresent the  
627 distribution of the aggregate travel demand across zones.

628 Although traditional travel surveys are increasingly being replaced by smartphone based  
629 surveys, which alleviate the issue of misreporting of trips, issues with representativeness and  
630 sample size remain, as well as with encouraging respondents to provide a sufficiently long  
631 stream of data (cf. Calastri et al., 2019). On the other hand, while mobile phone call detail  
632 record (CDR) data is widely available, large in size and more representative, it is lacking  
633 information on core causal variables.

634 The paper demonstrates the feasibility of a joint modelling framework to find the best fit at the  
635 joint level (i.e. between the aggregate and disaggregate levels) by combining household travel  
636 survey, census, and CDR data. The census data is crucial in creating a bridge between the two  
637 other data sources. The joint modelling framework operates by adjusting the parameter scale(s)  
638 of a pre-estimated base model to jointly optimise the prediction accuracy with respect to the  
639 reported trips in travel survey data and the zonal aggregate trip productions derived from CDR  
640 data. Three different approaches of parameter scaling were investigated (i.e. uniform,  
641 alternative specific, and variable specific scaling corresponding to joint models 1, 2, and 3  
642 respectively). All the three joint models were found to have higher temporal and spatial  
643 transferability, as well as better forecasting performance compared to the base model which  
644 relies on household travel survey data alone, thus making them more reliable. Although  
645 variable specific scaling (Model 3) produced the best overall results, there is a need for further  
646 research using data from different contexts to investigate if this finding is universally  
647 applicable. In particular, in this case, we did not have any independent measure to confirm that  
648 either of the data represented the ground truth which prompted us to give equal weight to the  
649 two types of data. This may not be the case in all contexts. More work is also needed on how  
650 to specify the joint likelihood combining the information from the two types of data and  
651 investigating the impact of the distribution of the error term, potential spatial correlation, etc.

652 Although the proposed framework has been tested in the context of trip generation, it has  
653 potential benefits in improving the modelling of the other transport choices (such as mode  
654 choice, route choice, departure time choice etc.). We conclude that the results of this study  
655 serve as a proof-of-concept that mobile phone data can be fused with traditional data sources  
656 to improve the temporal and spatial transferability of models. This approach is particularly  
657 important in the context of developing countries where reliable traditional data sources are  
658 scarce, and models making use of low-cost passive data to enhance their temporal and spatial  
659 transferability are invaluable.

660

661 **References**

- 662 ADB & ILO 2016. Bangladesh: Looking beyond garments: Employment diagnostic study.  
663 Manila, Phillipines: Asian Development Bank and International Labour Organization.
- 664 Agyemang-Duah, K. & Hall, F. L. 1997. Spatial transferability of an ordered response model  
665 of trip generation. *Transportation Research Part A: Policy and Practice*, 31, 389-402.
- 666 Barthelemy, J. & Toint, P. L. 2013. Synthetic population generation without a sample.  
667 *Transportation Science*, 47, 266-279.
- 668 BBS 2012. Community Report: Dhaka Zila: June 2012. *Population and Housing Census*  
669 *2011*. Dhaka: Bangladesh Bureau of Statistics (BBS).
- 670 BBS 2013. District Statistics 2011 Dhaka. Dhaka: Bangladesh Bureau of Statistics.
- 671 Beckman, R. J., Baggerly, K. A. & McKay, M. D. 1996. Creating synthetic baseline  
672 populations. *Transportation Research Part A: Policy and Practice*, 30, 415-429.
- 673 Ben-Akiva, M. E. & Lerman, S. R. 1985. *Discrete choice analysis: theory and application to*  
674 *travel demand*, MIT press.
- 675 Bhat, C. R. & Pulugurta, V. 1998. A comparison of two alternative behavioral choice  
676 mechanisms for household auto ownership decisions. *Transportation Research Part*  
677 *B: Methodological*, 32, 61-75.
- 678 Bwambale, A., Choudhury, C. F. & Hess, S. 2017. Modelling trip generation using mobile  
679 phone data: A latent demographics approach. *Journal of Transport Geography*.
- 680 Bwambale, A., Choudhury, C. F. & Sanko, N. Modelling Car Trip Generation in the  
681 Developing World: The Tale of Two Cities. Transportation Research Board 94th  
682 Annual Meeting, 2015.
- 683 Calabrese, F., Di Lorenzo, G., Liu, L., & Ratti, C. (2011). Estimating origin-destination flows  
684 using mobile phone location data. *IEEE Pervasive Computing*, (4), 36-44.
- 685 Cárcamo, J. G., Vogel, R. G., Terwilliger, A. M., Leidig, J. P. & Wolffe, G. Generative  
686 models for synthetic populations. Proceedings of the Summer Simulation Multi-  
687 Conference, 2017. Society for Computer Simulation International, 7.
- 688 Casati, D., Müller, K., Fourie, P. J., Erath, A. & Axhausen, K. W. 2015. Synthetic population  
689 generation by combining a hierarchical, simulation-based approach with reweighting  
690 by generalized raking. *Transportation Research Record: Journal of the*  
691 *Transportation Research Board*, 107-116.
- 692 Chen, C., Bian, L. & Ma, J. 2014. From traces to trajectories: How well can we guess activity  
693 locations from mobile phone traces? *Transportation Research Part C: Emerging*  
694 *Technologies*, 46, 326-337.
- 695 Choupani, A.-A. & Mamdoohi, A. R. 2015. Population Synthesis in Activity-Based Models:  
696 Tabular Rounding in Iterative Proportional Fitting. *Transportation Research Record:*  
697 *Journal of the Transportation Research Board*, 1-10.

- 698 Çolak, S., Alexander, L. P., Alvim, B. G., Mehndiretta, S. R. & González, M. C. Analyzing  
699 Cell Phone Location Data for Urban Travel: Current Methods, Limitations and  
700 Opportunities. Transportation Research Board 94th Annual Meeting, 2015.
- 701 Deutsch, K., McKenzie, G., Janowicz, K., Li, W., Hu, Y., & Goulias, K. (2012). Examining  
702 the use of smartphones for travel behavior data collection. In The 13th International  
703 Conference on Travel Behavior Research Toronto, Toronto.
- 704 Farooq, B., Bierlaire, M., Hurtubia, R. & Flötteröd, G. 2013. Simulation based population  
705 synthesis. *Transportation Research Part B: Methodological*, 58, 243-263.
- 706 Ferrer López, S., & Ruiz Sánchez, T. (2014). Travel behavior characterization using raw  
707 accelerometer data collected from smartphones. *Procedia Social and Behavioral Sciences*,  
708 160, 140-149.
- 709 Gerpott, T. J. & Thomas, S. 2014. Empirical research on mobile Internet usage: A meta-analysis of  
710 the literature. *Telecommunications Policy*, 38, 291-310.
- 711 Greaves, S., Ellison, A., Ellison, R., Rance, D., Standen, C., Rissel, C., & Crane, M. (2015).  
712 A web-based diary and companion smartphone app for travel/activity surveys.  
713 *Transportation Research Procedia*, 11, 297-310.
- 714 Groves, R. M. 2006. Nonresponse rates and nonresponse bias in household surveys. *Public  
715 opinion quarterly*, 646-675.
- 716 GSM Association. 2017. *The Mobile Economy 2017* [Online]. Available:  
717 [https://www.gsmainelligence.com/research/?file=9e927fd6896724e7b26f33f61db5b9](https://www.gsmainelligence.com/research/?file=9e927fd6896724e7b26f33f61db5b9d5&download)  
718 [d5&download](https://www.gsmainelligence.com/research/?file=9e927fd6896724e7b26f33f61db5b9d5&download) [Accessed 04 November 2017].
- 719 Guo, J. & Bhat, C. 2007. Population synthesis for microsimulating travel behavior.  
720 *Transportation Research Record: Journal of the Transportation Research Board*, 92-  
721 101.
- 722 Huang, H., Gartner, G., Krisp, J. M., Raubal, M., & Van de Weghe, N. (2018). Location  
723 based services: ongoing evolution and research agenda. *Journal of Location Based  
724 Services*, 12(2), 63-93.
- 725 Iqbal, M. S., Choudhury, C. F., Wang, P. & González, M. C. 2014. Development of origin–  
726 destination matrices using mobile phone call data. *Transportation Research Part C:  
727 Emerging Technologies*, 40, 63-74.
- 728 Itsubo, S. and Hato, E., 2006. *Effectiveness of household travel survey using GPS-equipped  
729 cell phones and Web diary: Comparative study with paper-based travel survey* (No.  
730 06-0701).
- 731 Janzen, M., Müller, K. & Axhausen, K. W. Population Synthesis for Long-Distance Travel  
732 De-mand Simulations using Mobile Phone Data. 6th Symposium of the European  
733 Association for Research in Transportation (hEART 2017), 2017.
- 734 JICA 2010. Dhaka Urban Transport Network Development Study (DHUTS) in Bangladesh,  
735 Final Report. Dhaka: Japan International Cooperation Agency.

- 736 Kressner, J. D. 2017. Synthetic Household Travel Data Using Consumer and Mobile Phone  
737 Data. *Final Report for NCHRP IDEA Project 184*. Transportation Research Board.
- 738 Marschak, J. 1960. Binary Choice Constraints on Random Utility Indications. *In: ARROW,*  
739 *K. (ed.) Stanford Symposium on Mathematical Methods in the Social Science.*  
740 *Stanford, California: Stanford University Press.*
- 741 McFadden, D. 1974. Conditional logit analysis of qualitative choice behavior. *Frontiers in*  
742 *Econometrics*, 105-142.
- 743 Ortúzar, J. D. D. & Willumsen, L. G. 2011. *Modelling transport*, John Wiley & Sons.
- 744 Pan, C., Lu, J., Di, S. & Ran, B. 2006. Cellular-based data-extracting method for trip  
745 distribution. *Transportation Research Record: Journal of the Transportation*  
746 *Research Board*, 33-39.
- 747 Panigutti, C., Tizzoni, M., Bajardi, P., Smoreda, Z. & Colizza, V. 2017. Assessing the use of  
748 mobile phone data to describe recurrent mobility patterns in spatial epidemic models.  
749 *Royal Society open science*, 4, 160950.
- 750 Patterson, Z., & Fitzsimmons, K. (2016). Datamobile: Smartphone travel survey experiment.  
751 *Transportation Research Record*, 2594(1), 35-43.
- 752 Pettersson, P. & Schmöcker, J.-D. 2010. Active ageing in developing countries?—trip  
753 generation and tour complexity of older people in Metro Manila. *Journal of Transport*  
754 *Geography*, 18, 613-623.
- 755 Pritchard, D. R. & Miller, E. J. 2012. Advances in population synthesis: fitting many  
756 attributes per agent and fitting to household and person margins simultaneously.  
757 *Transportation*, 39, 685-704.
- 758 QGIS Development Team. 2018. *QGIS Geographic Information System* [Online]. Available:  
759 <https://qgis.org/en/site/> [Accessed 14 August 2018].
- 760 Rao, B., & Minakakis, L. (2003). Evolution of mobile location-based services.  
761 *Communications of the ACM*, 46(12), 61-65.
- 762 Rolstad, S., Adler, J. & Rydén, A. 2011. Response burden and questionnaire length: is shorter  
763 better? A review and meta-analysis. *Value in Health*, 14, 1101-1108.
- 764 Ros, O. G. C. & Albertos, P. G. 2016. D5.4 Enhanced Version of MATSim: Synthetic  
765 Population Module. *Innovative Policy Modelling and Governance Tools for*  
766 *Sustainable Post-Crisis Urban Development (INSIGHT)*. Madrid, Spain: INSIGHT  
767 Consortium.
- 768 Ryan, J., Maoh, H. & Kanaroglou, P. 2009. Population synthesis: Comparing the major  
769 techniques using a small, complete population of firms. *Geographical Analysis*, 41,  
770 181-203.
- 771 Safi, H., Assemi, B., Mesbah, M., & Ferreira, L. (2016). Trip detection with smartphone-  
772 assisted collection of travel data. *Transportation Research Record*, 2594(1), 18-26.

- 773 Safi, H., Assemi, B., Mesbah, M., Ferreira, L., & Hickman, M. (2015). Design and  
774 implementation of a smartphone-based travel survey. *Transportation Research*  
775 *Record*, 2526(1), 99-107.
- 776 Shin, D., Aliaga, D., Tunçer, B., Arisona, S. M., Kim, S., Zünd, D., & Schmitt, G. (2015).  
777 Urban sensing: Using smartphones for transportation mode classification. *Computers,*  
778 *Environment and Urban Systems*, 53, 76-86.
- 779 Stopher, P., FitzGerald, C. and Xu, M., 2007. Assessing the accuracy of the Sydney  
780 Household Travel Survey with GPS. *Transportation*, 34(6), pp.723-741.
- 781 Sun, L. & Erath, A. 2015. A Bayesian network approach for population synthesis.  
782 *Transportation Research Part C: Emerging Technologies*, 61, 49-62.
- 783 Vlassenroot, S., Gillis, D., Bellens, R., & Gautama, S. (2015). The use of smartphone  
784 applications in the collection of travel behaviour data. *International Journal of*  
785 *Intelligent Transportation Systems Research*, 13(1), 17-27.
- 786 Voas, D. & Williamson, P. 2000. An evaluation of the combinatorial optimisation approach  
787 to the creation of synthetic microdata. *International Journal of Population*  
788 *Geography*, 6, 349-366.
- 789 Vogel, N., Theisen, C., Leidig, J. P., Scripps, J., Graham, D. H. & Wolffe, G. 2015. Mining  
790 Mobile Datasets to Enable the Fine-Grained Stochastic Simulation of Ebola  
791 Diffusion. *Procedia Computer Science*, 51, 765-774.
- 792 White, J. & Wells, I. Extracting Origin Destination Information from Mobile Phone Data.  
793 Eleventh International Conference on Road Transport Information and Control (Conf.  
794 Publ. No. 486), March 2002 London. IET, pp. 30 - 34.
- 795 Wu, L., Yang, B., & Jing, P. (2016). Travel mode detection based on GPS raw data collected  
796 by smartphones: a systematic review of the existing methodologies. *Information*, 7(4),  
797 67.
- 798 Xiao, Y., Low, D., Bandara, T., Pathak, P., Lim, H. B., Goyal, D., ... & Ben-Akiva, M. (2012,  
799 January). Transportation activity analysis using smartphones. In 2012 IEEE  
800 Consumer Communications and Networking Conference (CCNC) (pp. 60-61). IEEE.
- 801 Xiao, G., Juan, Z., & Zhang, C. (2016). Detecting trip purposes from smartphone-based travel  
802 surveys with artificial neural networks and particle swarm optimization.  
803 *Transportation Research Part C: Emerging Technologies*, 71, 447-463.
- 804 Ye, X., Konduri, K., Pendyala, R. M., Sana, B. & Waddell, P. A methodology to match  
805 distributions of both household and person attributes in the generation of synthetic  
806 populations. 88th Annual Meeting of the Transportation Research Board,  
807 Washington, DC, 2009.
- 808 Zhang, D. (2018). Social-enabled Urban Data Analytics, Doctoral Dissertation, University of  
809 California Berkeley  
810 [https://digitalassets.lib.berkeley.edu/etd/ucb/text/Zhang\\_berkeley\\_0028E\\_17723.pdf](https://digitalassets.lib.berkeley.edu/etd/ucb/text/Zhang_berkeley_0028E_17723.pdf)  
811 [accessed 14.5.2020]

- 812 Zhao, F., Pereira, F.C., Ball, R., Kim, Y., Han, Y., Zegras, C. and Ben-Akiva, M., 2015.  
813 Exploratory analysis of a smartphone-based travel survey in  
814 Singapore. *Transportation Research Record: Journal of the Transportation Research*  
815 *Board*, 2(2494), pp.45-56.
- 816 Zhou, X., Yu, W., & Sullivan, W. C. (2016). Making pervasive sensing possible: Effective  
817 travel mode sensing based on smartphones. *Computers, Environment and Urban*  
818 *Systems*, 58, 52-59.
- 819 Zhu, Y. & Ferreira Jr, J. 2014. Synthetic population generation at disaggregated spatial scales  
820 for land use and transportation microsimulation. *Transportation Research Record*,  
821 2429, 168-177.
- 822 Zilske, M. & Nagel, K. 2014. Studying the accuracy of demand generation from mobile  
823 phone trajectories with synthetic data. *Procedia Computer Science*, 32, 802-807.
- 824 Zilske, M. & Nagel, K. 2015. A simulation-based approach for constructing all-day travel  
825 chains from mobile phone data. *Procedia Computer Science*, 52, 468-475.
- 826