

Commuting Flow Prediction using OpenStreetMap Data

Kuldip Singh Atwal^{1*}, Taylor Anderson¹, Dieter Pfoser¹
and Andreas Züfle²

¹George Mason University, Geography and Geoinformation
Science, Fairfax, VA 22030, United States.

²Emory University, Department of Computer Science, Atlanta,
GA 30322, United States.

*corresponding:katwal@gmu.edu.

Abstract

Accurately predicting commuting flows is crucial for sustainable urban planning and preventing disease spread due to human mobility. While recent advancements have produced effective models for predicting these recurrent flows, the existing methods rely on datasets exclusive to a few study areas, limiting the transferability to other locations. This research broadens the applicability of state-of-the-art commuting flow prediction models by employing features from freely accessible and globally available OpenStreetMap data. We show that the prediction accuracy of several state-of-the-art models using open data is comparable to location-specific and proprietary data. Our experiments indicate that consistent with theoretical and analytical models, building types, distance, and population are the determining characteristics for mobility related to commuting. Furthermore, our experiments show that predicted flows closely match ground truth flows. It helps establish the practical relevance of flow prediction models for real-world applications such as urban planning and epidemiology.

Introduction

Understanding how individuals routinely move from one place to another is as challenging as it is significant [1, 2]. Commuting flow prediction estimates

30 the number of people moving between regions in a geographic area based on
31 descriptive features, such as population [3], distance to other locations [4], and
32 land use type [5]. Commuting flow prediction is helpful in many applications,
33 such as understanding migration patterns [6, 7], urban planning [8, 9], and
34 epidemiology [10, 11]. Considering that commuting flows vary little from day
35 to day [12, 13], the goal is typically to predict a set of static flows where
36 each flow represents the average number of daily commuters between origin-
37 destination pairs, i.e., home and work locations [14, 15]. Therefore, similar to
38 other approaches [9, 16], we define the term flow prediction as the task of
39 predicting repetitive static flows rather than forecasting flows along a series of
40 points in time using historical data, which is a time series problem.

41 Analytical flow prediction approaches include spatial interaction models
42 such as the gravity model [17] and its extensions, including the radiation
43 model [18–20], the intervening opportunities model [21, 22], and the competing
44 migrants model [23]. Each model proposes different characteristics to predict
45 accurate flows. For example, the gravity model assumes that the flow between
46 locations is a function of two main characteristics: (i) the population at both
47 locations and (ii) the distance between them. In another example, the inter-
48 vening opportunities model replaces distance with the number of opportunities
49 at the destination location that satisfy the trip objective [24]. Thus, when pre-
50 dicting commuter flows, the “opportunity” in question might be the number
51 of commercial businesses.

52 More recently, machine learning models for commuting flow prediction far
53 outperform the traditional mathematical approaches when comparing the pre-
54 dicted flows with ground truth [16, 25–28]. These models leverage machine
55 learning approaches that can more flexibly incorporate different features of
56 the origin-destination and can capture complex and non-linear relationships in
57 the data [29–31]. Many studies use spatiotemporal characteristics to address
58 the flow prediction problem using neural networks [32–35], which can also be
59 combined with ordinary differential equations [36]. A current state-of-the-art
60 model, the Geo-contextual Multitask Embedding Learner (GMEL) [9] learns
61 commuting flows based on origin-destination features and their spatial con-
62 texts. GMEL uses 65 features derived from the 2015 NYC Primary Land Use
63 Tax Lot Output (PLUTO)[37] dataset. In another example, the ConvGCN-RF
64 model [38] uses convolutional neural network, graph convolutional network,
65 and a random forest regressor to predict the commuting flow based on origin-
66 destination features related to land use, as well as the residential and working
67 population for homogeneous spatial units in the region of Beijing, China.
68 Spadon et al. [39] derive 22 urban features from datasets provided by the
69 Brazilian Institute of Geography and Statistics (IBGE) to predict intercity
70 commuting in Brazil.

71 Despite the ability of such models to accurately predict flows, these high-
72 performing models use a large number of input features derived from location-
73 specific data sets that are not available outside of the study area. It makes
74 the use of the model in other data-poor study regions challenging. In addition,

75 given the variety of different input features used across models, it is difficult
76 to compare models independent of the used data.

77 Our goal in this research is to address the limitations that restrict the appli-
78 cability of current commuting flow prediction models to arbitrary study areas.
79 More precisely, we assess the effectiveness of these models by employing a min-
80 imal set of input features obtained from a globally accessible dataset called
81 OpenStreetMap (OSM) [40]. Moreover, since numerous models are assessed
82 using high-level metrics, such as Root Mean Square Error (RMSE), Coeffi-
83 cient of Determination (R^2), and Common Part of Commuters (CPC), which
84 provide limited insight into the model's ability to replicate authentic patterns
85 intrinsic to commuting flows, we investigate the degree to which these models
86 prove valuable in predicting significant mobility flows at different scales. The
87 extensive analysis of flows explains some of the underlying phenomena driving
88 commuting mobility. Motivated by features used in previous theoretical work,
89 including the gravity model and intervening opportunities model, we consider
90 three characteristics to address the flow prediction problem: building types,
91 distance, and population. Specifically, we extract nine input features from open
92 data based on these characteristics that potentially drive commuters' mobility,
93 as follows:

- 94 • The number (count), density, and area of residential and non-residential
95 buildings, respectively (six features),
- 96 • Region population and population density (two features), and
- 97 • Distance between census tracts (one feature)

98 The feature generation leverages existing work on using a machine learning
99 approach to classify OSM building types [41] beyond the information avail-
100 able in OSM. Additionally, we use Open Source Routing Machine (OSRM),
101 an OSM-based routing API [42], to generate trip duration between all pairs
102 of regions used to represent distance. Using these features, we first provide a
103 fair comparison of different models for predicting commuter flows. Our first
104 case study focuses on New York City (NYC), USA, at the census tract gran-
105 ularity, where we compare two state-of-the-art models, including GMEL [9]
106 and Deep Gravity [27], and eXtreme Gradient Boosting (XGBoost) and ran-
107 dom forests (RF) as out-of-the-box models commonly used for commuting
108 prediction [25, 26, 39]. The 2015 Longitudinal Employer-Household Dynam-
109 ics (LEHD) Origin-Destination Employment Statistics (LODES) data [43] is
110 used to evaluate the effectiveness of our approach. We compare model per-
111 formance using OSM-derived features with region-specific features unavailable
112 outside the study area. Finally, we demonstrate the inherent flexibility of using
113 OSM-derived features by predicting commuting flows for Fairfax County, USA.
114 Results from both case studies validate the intuitive understanding that the
115 destination flows, commuters going to workplaces, are concentrated in a few
116 places.

Table 1: Notations used in the study

Notation	Meaning
$A = \{a_1, \dots, a_n\}$	The study region
a_i	A subregion of the study region
n	The number of subregions
T_{ij}	The ground truth commuter flow from Region a_i to Region a_j
\widehat{T}_{ij}	The estimated commuter flow from Region a_i to Region a_j
d_{ij}	Spatial distance between two subregions
$O_i = \sum_j T_{ij}$	The total outflow of region a_i (to any other region)
$I_i = \sum_j T_{ji}$	The total inflow of region a_i (to any other region)
$\widehat{O}_i = \sum_j \widehat{T}_{ij}$	The estimated outflow of region a_i (to any other region)
$\widehat{I}_i = \sum_j \widehat{T}_{ji}$	The estimated inflow of region a_i (to any other region)

117 Results

118 Results show that we can get accurate flow predictions between census tracts
 119 using features derived from open data, and population, building type, and
 120 distance are the significant characteristics driving commuting mobility. The
 121 evidence from experiments at multiple scales suggests our approach produces
 122 meaningful mobility patterns while providing notable insights into the com-
 123 muting flows. Before presenting our findings, we briefly define the commuting
 124 flow prediction problem.

125 Problem Definition

126 The commuting flow prediction problem can be defined as follows. Table 1
 127 summarizes the used notations.

128 **Definition 1** (Commuting Flow Prediction). Let A denote a study region
 129 partitioned into n smaller regions (a_1, \dots, a_n) , such as census tracts in the
 130 United States. For each region a_i , let f_i denote a corresponding set of features,
 131 and for each pair of regions a_i, a_j , let d_{ij} denote a distance measure between
 132 regions. Given these features and distance, the task is to predict the commuting
 133 flow T_{ij} for each pair of regions $a_i, a_j \in A$.

134 Benchmark Results

135 Using OSM data and the same set of derived features for New York City
 136 (NYC), Table 2 provides the commuting flow prediction accuracy for state-of-
 137 the-art models GMEL [9] and Deep Gravity [27], and out-of-the-box models
 138 XGBoost [44] and RF [45]. To evaluate model performance, we use the RMSE
 139 [46], the Coefficient of Determination R^2 [47], and the Common Part of
 140 Commuters metric [48].

The RMSE is defined as follows:

$$RMSE(A) = \sqrt{\frac{\sum_{a_{ij}} (\hat{T}_{ij} - T_{ij})^2}{n}}$$

141 where A is the NYC study region, \hat{T}_{ij} is the predicted commuting flow (c.f.
 142 Definition 1), T_{ij} is the ground truth flow obtained for NYC using LODES
 143 data, and n is the number of census tracts of NYC.

144 RMSE values are notoriously difficult to interpret. For example, it is not
 145 clear to what degree a prediction with an RMSE of 2.279 is accurate. As
 146 such, we also provide the Coefficient of Determination R^2 and Common Part
 147 of Commuters (CPC) to provide an additional evaluation of model accuracy.
 148 Although the R^2 is well known and measures the fraction of variance explained
 149 by the model, the Common Part of Commuters (CPC) is less known. Thus,
 150 we define CPC, as follows:

$$CPC(A) = \frac{2 \sum_{a_{ij}} \min(\hat{T}_{ij}, T_{ij})}{\sum_{a_{ij}} \hat{T}_{ij} + \sum_{a_{ij}} T_{ij}}$$

151 CPC is 0 when predicted and ground truth flows do not overlap and 1 when
 152 both are identical [49].

153 Based on the results presented in Table 2, GMEL has the lowest RMSE and
 154 highest CPC and R^2 in comparison to XGBoost, Deep Gravity, and RF. Note
 155 that the two state-of-the-art models, GMEL and Deep Gravity, are originally
 156 implemented to predict commuting flow using a different set of input features,
 157 making them difficult to compare. Therefore, in order to evaluate the perfor-
 158 mance of the models independent of the data, the models are implemented
 159 using the same set of input features derived from OSM. The experiment shows
 160 that GMEL is the best-performing model compared to other models using the
 161 same features.

Table 2: Evaluation of different flow prediction models using OSM data

Model	RMSE	CPC	R^2
GMEL	2.279	0.495	0.535
XGBoost	3.125	0.261	0.111
Deep Gravity	3.144	0.325	0.078
RF	3.228	0.218	0.051

162 Comparative Analysis

163 Given our results showing that GMEL is the best-performing model, we next
 164 compare the performance of the originally proposed GMEL model, which
 165 leverages the PLUTO dataset [37] available only for New York City, with
 166 the performance of GMEL using globally available OSM data. To distinguish

167 between the two, we call the original model GMEL-PLUTO and our approach
 168 GMEL-OSM throughout the rest of the paper. In other words, GMEL-PLUTO
 169 uses region-specific PLUTO data for flow prediction, while GMEL-OSM uses
 170 features derived from OSM data.

Table 3: Comparison of OSM and PLUTO data using GMEL model for NYC

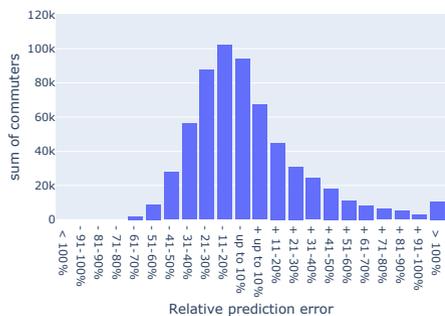
Features	RMSE	CPC	R^2
GMEL-OSM	2.279	0.495	0.535
GMEL-PLUTO	2.084	0.536	0.611

171 Table 3 shows that a comparable level of prediction accuracy can be
 172 achieved overall when using features derived from globally accessible and freely
 173 available OSM data. The R^2 value indicates that the three characteristics
 174 account for an 53.5% variation in commuting flows. Additionally, GMEL-OSM
 175 utilizes a smaller set of features to achieve accuracy close to GMEL-PLUTO
 176 with 65 features.

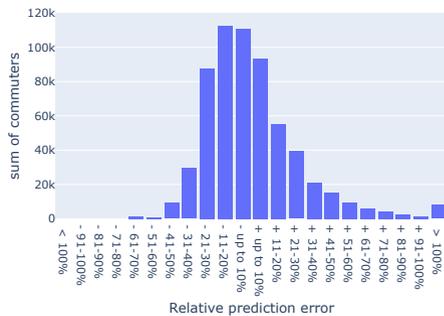
177 To better understand the ability of the models to capture meaningful mobil-
 178 ity patterns beyond aggregate metrics, we also evaluate the predicted sum of
 179 outgoing commuters from an origin location a_i denoted as $\hat{O}_i = \sum_j \hat{T}_{ij}$, which
 180 we call *outflows*, and the predicted sum of incoming commuters to a destina-
 181 tion location a_i denoted as $\hat{I}_i = \sum_j \hat{T}_{ji}$, which we call *inflows*. The \hat{O}_i and
 182 \hat{I}_i for each region a_i stemming from the GMEL-OSM and GMEL-PLUTO
 183 predictions are then compared to the ground truth values $O_i = \sum_j T_{ij}$ and
 184 $I_i = \sum_j T_{ji}$ derived from LODES data for NYC.

185 Figure 1 shows the distribution of relative prediction errors for the out-
 186 flows $\frac{O_i - \hat{O}_i}{O_i}$ and the inflows $\frac{I_i - \hat{I}_i}{I_i}$ for GMEL-OSM (Figures 1a and 1c) and for
 187 GMEL-PLUTO (Figure 1b and 1d). We observe that GMEL-OSM is compa-
 188 rable with GMEL-PLUTO to predict outflows, but performs somewhat weaker
 189 for inflows. It is likely due to the nature of commuting flows, with inflows being
 190 limited to a small group of destination census tracts (cf. discussion in the Data
 191 Section). Even so, the results show the practicality of predicted flows compared
 192 to ground truth data. Out of those census tracts where flow is over-predicted
 193 by more than 100%, many have a commuting flow count of 10 individuals or
 194 fewer. It indicates that our approach is capable of predicting real-world com-
 195 muting mobility at the tract level, where the flow count is generally more than
 196 10.

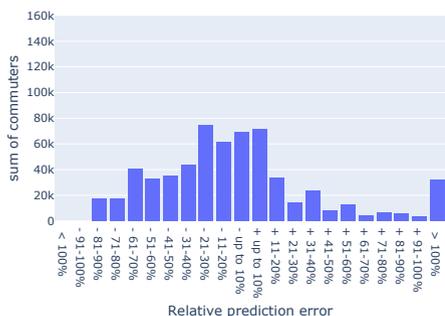
197 To assess the accuracy of the predicted inflows and outflows for census
 198 tracts, Figure 2 shows scatter plots comparing the ground truth flows against
 199 the predicted flows using GMEL-OSM (Figures 2a and 2c) and GMEL-PLUTO
 200 (Figures 2b and 2d). Both models tend to overestimate inflows that are smaller
 201 in the real world and underestimate large inflows, as indicated by the points
 202 that fall above and below the identity line. Likewise, both models also tend
 203 to overestimate smaller outflows. Again, while both models produce similar



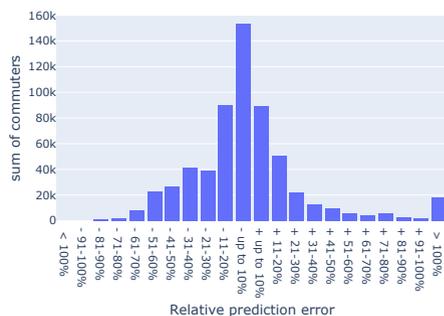
(a) Percentage of under or overestimation of NYC commuters' outflows using GMEL-OSM.



(b) Percentage of under or overestimation of NYC commuters' outflows using GMEL-PLUTO.



(c) Percentage of under or overestimation of NYC commuters' inflows using GMEL-OSM.



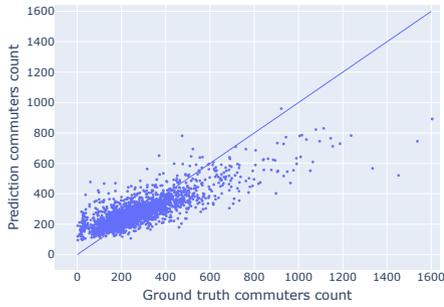
(d) Percentage of under or overestimation of NYC commuters' inflows using GMEL-PLUTO.

Fig. 1: Comparison of GMEL-OSM and GMEL-PLUTO commuters under or overestimation in NYC flows.

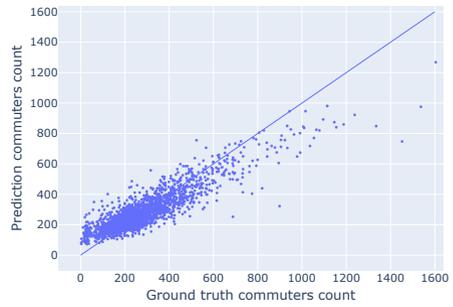
204 results for outflows, GMEL-PLUTO (65 custom feature model) seems to per-
 205 form better when predicting the inflows, essentially confirming the results of
 206 Figure 1 at a more granular level.

207 We note that the maximum number of commuters going to a census tract
 208 is much higher than coming from a home location, which is consistent in both
 209 prediction models and the ground truth. It indicates that the inflows are much
 210 denser to specific census tracts or workplaces. We investigate and explain this
 211 phenomenon in our Data Section.

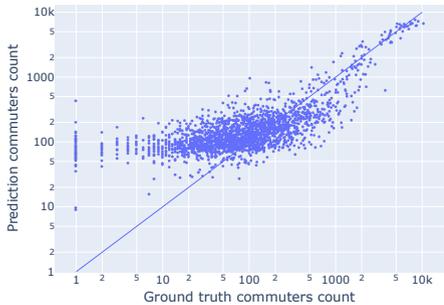
212 We can also map the differences between predicted and ground truth out-
 213 flows as presented in Figure 3 and inflows presented in Figure 4. Positive
 214 relative prediction errors indicate over-prediction and are depicted in shades of
 215 blue colors. In contrast, negative percentages indicate under-prediction and are
 216 shown in shades of red. Green shows a prediction largely matching the ground
 217 truth flows. Note that the large tracts in the south of the study area are mostly



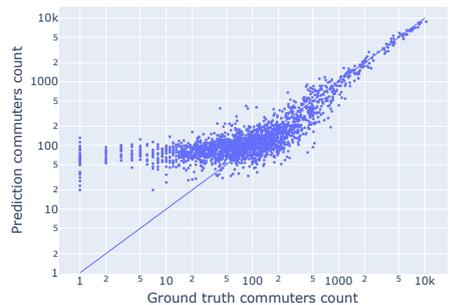
(a) Comparison of NYC commuters' outflows using GMEL-OSM with ground truth.



(b) Comparison of NYC commuters' outflows using GMEL-PLUTO with ground truth.



(c) Comparison of NYC commuters' inflows using GMEL-OSM with ground truth (log-log scale).



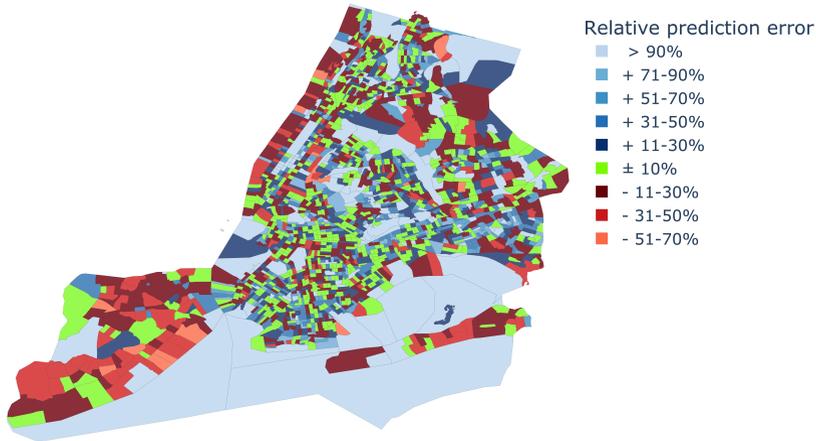
(d) Comparison of NYC commuters' inflows using GMEL-PLUTO with ground truth (log-log scale).

Fig. 2: Comparison of GMEL-OSM and GMEL-PLUTO commuters with ground truth in NYC flows.

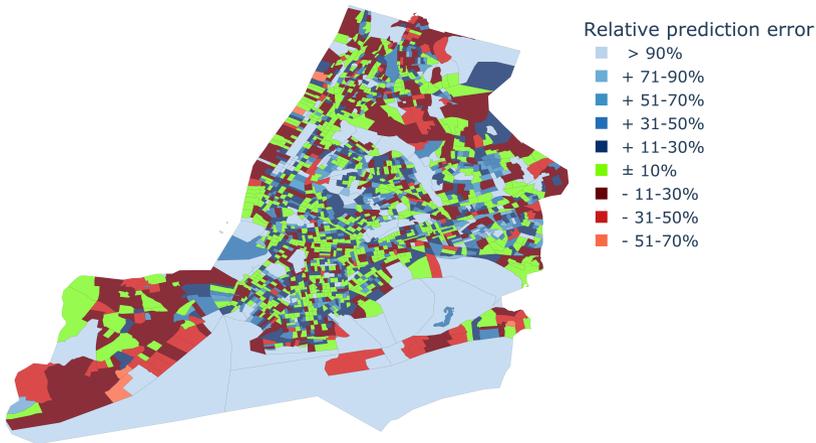
218 comprised of water, thus having small in and outflows. As a result, minor flow
 219 prediction errors for these census tracts provide high relative percentage errors
 220 and as such are shown as large light blue areas.

221 Upon comparing Figures 3 and 4, we can see that GMEL-OSM and GMEL-
 222 PLUTO flow predictions are very similar in terms of the relative prediction
 223 error. Both approaches have less success in predicting destination flows. It is
 224 once again likely due to the large number of features used in GMEL-PLUTO
 225 that are likely better at capturing the inflows to destination census tracts. We
 226 discuss steps that we may take to address this in future work in the Discussion
 227 Section.

228 To better understand the utility of predicted commuter flows, we also
 229 performed experiments focusing on a single origin (destination) tract to under-
 230 stand how well models can capture the distribution of destination (origin)
 231 tracts to (from) this tract. For this purpose, we select the census tract hav-
 232 ing the median outflow (GeoID: 36047037300, denoted as the *Origin Median*)



(a) Relative errors of NYC outflows using GMEL-OSM.

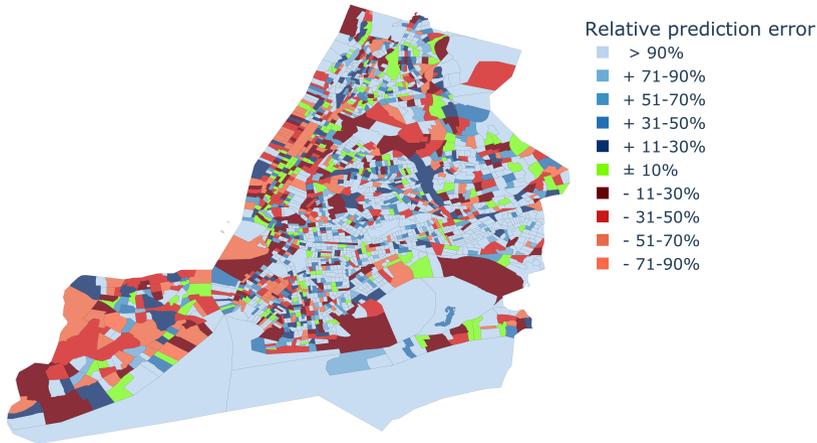


(b) Relative errors of NYC outflows using GMEL-PLUTO.

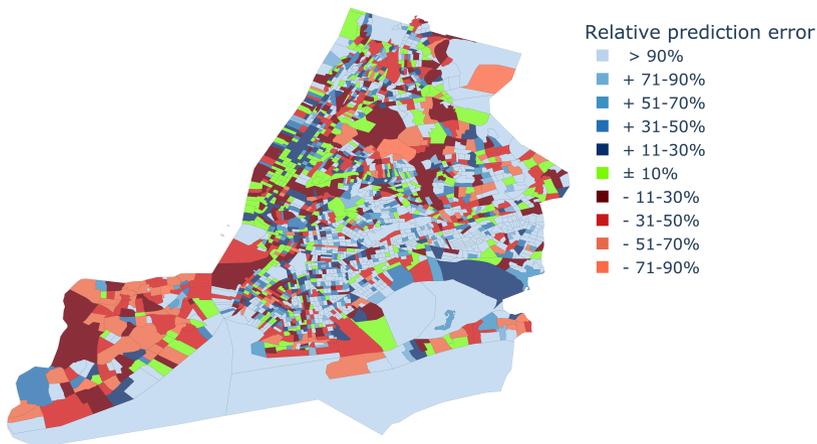
Fig. 3: Comparison of GMEL-OSM and GMEL-PLUTO in NYC outflows. Plotly version 5.13.0 was used to generate the maps.

233 and the census tract having the median inflow (GeoID 36005024800, denoted
 234 as the *Destination Median*). We use these two census tracts to evaluate (i)
 235 the distribution of outflows from the Origin Median to understand how well
 236 the models can understand where people commute to (from one specific cen-
 237 sus tract) and (ii) the distribution of inflows from the Destination Median
 238 to understand how well our models can capture the distribution of where people
 239 commute from (to one specific census tract).

240 Table 4 shows the results of these experiments. Out of all 448 census tracts
 241 in the NYC study region included in the test set, 354 tracts have a zero com-
 242 muting flow from the Origin Median. The remaining 94 census tracts having



(a) Relative errors of NYC inflows using GMEL-OSM.



(b) Relative errors of NYC inflows using GMEL-PLUTO.

Fig. 4: Comparison of GMEL-OSM and GMEL-PLUTO in NYC inflows. Plotly version 5.13.0 was used to generate the maps.

243 non-zero commuting flows capture a total of 244 commuters. Using GMEL-
 244 OSM, we have 332 predicted zero commuting flows and 116 predicted non-zero
 245 commuting flow. Out of the predicted 116 predicted non-zero flows, 48 match
 246 with the 94 ground truth non-zero flows. Out of the 332 predicted zero flows,
 247 286 match with the 354 ground truth flows. It yields an overall 74.5% accuracy
 248 in predicting whether any census tract has a non-zero flow from the Origin
 249 Median. Note that we round predictions to the nearest integer for this exper-
 250 iment, such as that a predicted zero flow is equivalent to a predicted flow of
 251 less than 0.5 individuals. We observe that for GMEL-PLUTO, the accuracy

Table 4: Single origin and destination census tract predictions

Census Tract	Approach	Zero Flows Count (Matching)	Non-Zero Flows Count (Matching)	Sum of Commuters
Origin Median	Ground Truth	354 (354)	94 (94)	244
	GMEL-OSM	332 (286)	116 (48)	212
	GMEL-PLUTO	345 (304)	103 (53)	201
Destination Median	Ground Truth	411 (411)	46 (46)	81
	GMEL-OSM	418 (393)	39 (21)	43
	GMEL-PLUTO	427 (398)	30 (17)	32

252 is higher at 79.6%, indicating that the model can better predict destination
 253 flows by leveraging PLUTO data.

254 Similarly, by considering only the Destination Median as a single destina-
 255 tion, GMEL-OSM and GMEL-PLUTO matched 90.5% and 90.8%, respec-
 256 tively, out of 457 origin tracts in the test set. We observe that the destination
 257 median has a relatively small number of only 81 incoming commuters in the
 258 ground truth. It is explained by the long-tail distribution of inflows, which we
 259 further investigate and explain in the Data Section.

260 Overall, we observe that while GMEL-OSM and GMEL-PLUTO provide
 261 very accurate flow predictions when aggregated to census tracts, the prediction
 262 of individual origin-destination flows remains challenging. The reason is that
 263 the vast majority of origin-destination flows are zero and among the non-zero
 264 flows, most flows are less than five individuals. Despite these small numbers,
 265 which correspond to rare events of individual origin-destination commutes,
 266 both GMEL-OSM and GMEL-PLUTO give good results.

267 Based on the results presented so far, we can conclude that there are
 268 marginal gains in performance by using a large number of region-specific fea-
 269 tures using GMEL-PLUTO, and we can achieve similar results with a small
 270 set of features derived from open data that is globally available. To examine
 271 whether GMEL-OSM is usable in other regions, we trained and tested the
 272 model for Fairfax County in Virginia and compared the predicted flows with
 273 the LODES data as ground truth. Note that we cannot compare GMEL-OSM
 274 with GMEL-PLUTO because the latter approach uses NYC-specific data,
 275 which is publicly unavailable for Fairfax.

276 Histograms in Figure 5 show the relative percentage errors of outflows and
 277 inflows at the tract level compared to the ground truth. Figure 6 demonstrates
 278 the trend of flow prediction for outflows and inflows, respectively. We observe
 279 that the model performance in Fairfax, VA is comparable, if not better than the
 280 NYC case study using GMEL-PLUTO. Based on the histograms, it appears
 281 that the commuting inflows for Fairfax are easier to predict and less extreme
 282 than in NYC.

283 Additionally, we trained GMEL-OSM using NYC data and tested the pre-
 284 trained model to predict the commuting flows for Fairfax to determine whether
 285 the model is useful in locations where training commuting flow data (obtained
 286 for the U.S. from LODES data) is not available. Table 5 shows that the model



(a) Percentage of under or overestimation of Fairfax commuters' outflows using GMEL-OSM.

(b) Percentage of under or overestimation of Fairfax commuters' inflows using GMEL-OSM.

Fig. 5: Commuters under or overestimation using GMEL-OSM for Fairfax.

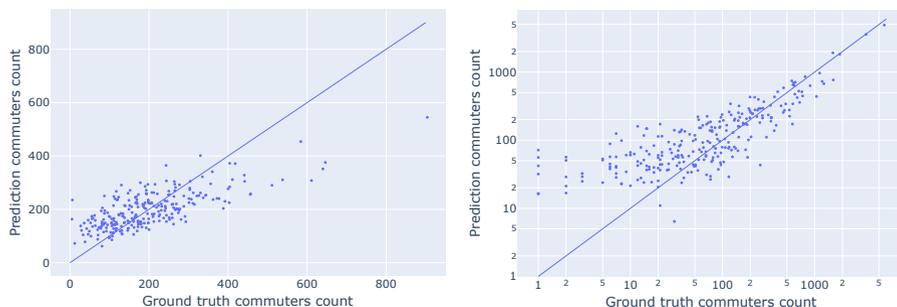
287 trained in NYC and transferred to Fairfax provides acceptable results by
 288 explaining 62.1% of the variation in the commuting flows of Fairfax, compared to
 289 to 70.2% using the model that was trained using Fairfax LODS data.

Table 5: Comparison of GMEL-OSM in Fairfax using transfer learning

Training data	RMSE	CPC	R^2
Fairfax	6.476	0.643	0.702
NYC	7.427	0.572	0.621

290 Discussion

291 Results for the two study areas show that commuting flows can be accurately
 292 predicted using features derived from OSM data, which is globally available
 293 and freely accessible. Comparative results reveal that GMEL-OSM achieves
 294 accuracy close to region-specific GMEL-PLUTO, which outperforms other
 295 state-of-the-art models but cannot be used outside NYC due to a lack of input
 296 data for other regions. The learning framework of GMEL-OSM relies on geo-
 297 graphic contextual information [50] for predicting commuting flows between
 298 origin-destination pairs of subregions. Our findings suggest that the OSM
 299 data captures the contextual information very well for the origin and desti-
 300 nation locations, providing a rich and effective source of input features for
 301 GMEL-OSM. Besides aggregated results, the in-depth analysis demonstrates
 302 the usefulness of the predicted flows for urban planning [51], disease trans-
 303 mission [52, 53], and other applications [54, 55]. We find that inflows are
 304 concentrated in a few destinations while outflows are more evenly distributed,
 305 validating the intuition that people commute to a few workplaces and reside in
 306 dispersed locations. Our analysis shows that GMEL-OSM effectively captures



(a) Comparison of Fairfax commuters' outflows using GMEL-OSM with ground truth. (b) Comparison of Fairfax commuters' inflows using GMEL-OSM with ground truth (log-log scale).

Fig. 6: Comparison of GMEL-OSM commuters prediction with ground truth for Fairfax flows.

307 this divergent phenomenon, matching the trend of outflows and inflows in the
 308 ground truth. Additionally, we also illustrate that the number of residential
 309 and non-residential buildings in census tracts plays a crucial role in predicting
 310 commuters' mobility. Our results indicate that building types, distance, and
 311 population are the essential characteristics driving commuting mobility.

312 While the population can be estimated at a fine-grained scale using OSM
 313 data [56, 57], for simplicity, we utilized the U.S. Census data as a proxy for
 314 this. In future work, we plan to extend our proposed approach for generat-
 315 ing population features, alleviating the need for census data. To investigate
 316 the explainability of the input features, we might explore a unified mechanism
 317 for interpreting predictions such as SHapley Additive exPlanations (SHAP)
 318 [58]. It would help us understand which features are useful for better commut-
 319 ing flow predictions, potentially leading to more suitable feature selection for
 320 improving the performance of our approach. Where we found relatively weaker
 321 prediction accuracy for the destination flows, there is an opportunity to exam-
 322 ine what features might improve this aspect of the predictions. Prior work
 323 shows the effectiveness of points of interest (PoIs) [59] and land use [60, 61]
 324 for predicting flows. Therefore, we would explore types of PoIs and land use
 325 as other characteristics driving mobility. Finally, our transfer learning results
 326 for Fairfax County show promise for future work in which we would plan to
 327 apply our approach to regions where LODES or equivalent commuting data is
 328 not publicly unavailable, potentially outside the U.S.

329 Methods

330 Models

331 We aim to predict commuting flows from three characteristics operationalized
 332 using globally available and openly accessible data. Therefore, we examine

333 four models including GMEL, Deep Gravity, XGBoost, and random forest
334 (RF), comparing their performance using the same set of features derived
335 from OSM. GMEL employs graph representation learning by using the graph
336 attention network (GAT) framework for capturing the geographic contextual
337 information from the nearby regions for commuting flow predictions. Given the
338 potentially unique characteristics of the regions, it uses two GATs separately
339 for origin and destination locations. As described in the proposed model [9],
340 we used one hidden layer and an embedding size of 128 as hyperparameters for
341 GMEL-OSM. Deep Gravity utilizes deep neural networks to generate mobility
342 flows using features retrieved from OSM and census data [27]. The main fea-
343 tures include land use, points of interest, road networks, and the population
344 of the study region. XGBoost is a regression tree gradient boosting model,
345 a highly scalable learning system capable of efficiently handling sparse data
346 and supporting multicore parallel computing for quick model exploration [44].
347 XGBoost has been shown to outperform traditional mathematical gravity and
348 radiation models for commuting flow prediction using U.S. Census data [25].
349 Random forests are the ensemble of individual tree predictions averaged for
350 regression problems and the prediction with maximum votes selected for clas-
351 sification problems [45]. Compared to the gravity model and artificial neural
352 networks, the accuracy for the random forest is higher for predicting commut-
353 ing flows in NYC in previous work [26]. As described in Results Section, we
354 evaluate the comparative performance of these models for our approach using
355 the parameters and configurations prescribed in the proposed studies.

356 Data

357 We use real-world commuting flows obtained from the Longitudi-
358 nal Employer-Household Dynamics (LEHD) Origin-Destination Employment
359 Statistics (LODES) 2015 dataset [43, 62] as ground truth for training and test-
360 ing the models. LODES data captures the raw number of commuters between
361 two regions at the census block level, and we aggregated it at the census tract
362 level.

363 Across the 2,168 NYC census tracts, there are $2168^2 = 4,700,224$ pair-
364 wise flows, of which 905,837 are non-zero with a total of 3,031,641 commuters.
365 Similarly, across the 263 Fairfax County census tracts, there are a possible
366 69,169 flows out of which 34,366 are non-zero flows, capturing 259,792 com-
367 muters. Unlike prior work [9, 12, 26], we include flows that are zero in the
368 ground truth LODES data. While LODES data does not explicitly include
369 zero flows in their data, the omitted flows between a pair of census tracts are
370 implicitly assumed to be zero values, which are missing from the evaluation
371 of prior work [9, 12, 26]. However, omitting such flows creates biased models
372 that learn that any pair of origin-destination census tracts must always have
373 at least a flow count of one commuter. Our experiments include all pairs of
374 census tracts, including zero flows, eliminating the bias. In other words, we
375 add zero flows to training and test sets of all evaluated models to allow a fair
376 evaluation. We note that due to this difference, the quantitative results we

377 report in the aggregated metrics in the Results Section (such as Table 2) are
 378 generally lower than reported in prior work, as our results include cases of
 379 flows where models predict a non-zero flow instead of a zero flow count in the
 380 ground truth. For training and testing, we split the flows into a 60% training
 381 set, a 20% validation set, and a 20% test set.

382 Table 6 presents the descriptive statistics for the NYC and Fairfax County
 383 LODES outflows O_i and inflows I_i aggregated at the tract level. We notice
 384 a much higher standard deviation of the inflow of commuters in both study
 385 regions. The maximum count of commuters for the inflows also highlights the
 386 significant difference in variance. Furthermore, the 3rd quantile values in both
 387 cases show the skewness in the distribution of commuters. These results demon-
 388 strate the concentrated nature of inflows in comparison to outflows, where the
 389 majority of commuters move to a small set of destination census tracts. There-
 390 fore, as our results suggest, it is much harder to predict the commuters' count
 391 for inflows.

Table 6: Descriptive statistics of ground truth data Data

Study Area	Flow Type	Mean	Standard Deviation	Min	25%	Median	75%	Max
NYC	Outflows	280	176	4	168	244	350	1604
	Inflows	280	817	1	34	81	190	10243
Fairfax	Outflows	197	120	5	111	173	255	904
	Inflows	197	482	1	21	67	180	5702

392 OSM is an open-source collaborative project that provides free access to
 393 geographic data collected by volunteers at the global level [40]. The OSM
 394 data is structured as a set of elements such as nodes, ways, and relations that
 395 represent points of interest, polylines or polygons, and more complex shapes
 396 consisting of relationships between simple elements. Tags of key and value pairs
 397 can describe all the elements. For instance, a polygon can be tagged with the
 398 key as building and value as a residential, describing a residential building.
 399 This way, OSM data provides extensive coverage of points, buildings, roads,
 400 parking lots, and many other types of geographic information via editable
 401 maps. The OSM data we used for this work consists of 1,090,752 NYC and
 402 204,671 Fairfax building footprints.

403 Features

404 The features used in the models for predicting the flows are derived from OSM
 405 and the 2010 U.S. Census data [63]. Previous work shows that building types
 406 are missing from a vast majority of OSM data, and the spatial and non-spatial
 407 features of the data can be used to categorize buildings into residential or
 408 non-residential types [41]. We use this classification method to label the OSM
 409 buildings data and derive six input features for our study. In the first step of
 410 data preparation, we classify buildings for NYC and Fairfax. And in the second

411 step, we calculate the count, area, and density of two building types for each
412 census tract, resulting in six features.

413 We use population and the population density for each tract as two more
414 input features. Although population estimates can be derived from OSM fea-
415 tures in the same way [56, 57], we use census data as a proxy for this approach.
416 Finally, we obtain the trip duration between the centroids of census tracts
417 using Open Source Routing Machine (OSRM) [42] and use it as the edge fea-
418 ture for the geo-adjacency network of GMEL-OSM. OSRM also relies on the
419 maps from the OSM road network for calculating the shortest paths between
420 O-D pairs.

421 Data availability

422 Data are available from OSF at <https://osf.io/sxzar/>

423 Code availability

424 The code is available in a GitHub repository at https://github.com/heykuldip/commuting_flows_prediction
425

426 References

- 427 [1] Gonzalez, M.C., Hidalgo, C.A., Barabasi, A.-L.: Understanding individual
428 human mobility patterns. *nature* **453**(7196), 779–782 (2008)
- 429 [2] Alessandretti, L., Aslak, U., Lehmann, S.: The scales of human mobility.
430 *Nature* **587**(7834), 402–407 (2020)
- 431 [3] Rong, C., Li, T., Feng, J., Li, Y.: Inferring origin-destination flows from
432 population distribution. *IEEE Transactions on Knowledge and Data*
433 *Engineering* **35**(1), 603–613 (2021)
- 434 [4] Levinson, D.M.: Accessibility and the journey to work. *Journal of*
435 *transport geography* **6**(1), 11–21 (1998)
- 436 [5] Layman, C.C., Horner, M.W.: Comparing methods for measuring excess
437 commuting and jobs-housing balance: empirical analysis of land use
438 changes. *Transportation Research Record* **2174**(1), 110–117 (2010)
- 439 [6] Luca, M., Barlacchi, G., Lepri, B., Pappalardo, L.: A survey on deep
440 learning for human mobility. *ACM Computing Surveys (CSUR)* **55**(1),
441 1–44 (2021)
- 442 [7] Jiang, R., Cai, Z., Wang, Z., Yang, C., Fan, Z., Chen, Q., Tsubouchi, K.,
443 Song, X., Shibasaki, R.: Deepcrowd: A deep model for large-scale citywide
444 crowd density and flow prediction. *IEEE Transactions on Knowledge and*
445 *Data Engineering* **35**(1), 276–290 (2021)

- 446 [8] Rodrigue, J.-P.: The Geography of Transport Systems. Routledge, UK
447 (2020)
- 448 [9] Liu, Z., Miranda, F., Xiong, W., Yang, J., Wang, Q., Silva, C.: Learning
449 geo-contextual embeddings for commuting flow prediction. In: Proceed-
450 ings of the AAAI Conference on Artificial Intelligence, vol. 34, pp. 808–816
451 (2020)
- 452 [10] Eubank, S., Guclu, H., Anil Kumar, V., Marathe, M.V., Srinivasan, A.,
453 Toroczkai, Z., Wang, N.: Modelling disease outbreaks in realistic urban
454 social networks. *Nature* **429**(6988), 180–184 (2004)
- 455 [11] Balcan, D., Colizza, V., Gonçalves, B., Hu, H., Ramasco, J.J., Vespignani,
456 A.: Multiscale mobility networks and the spatial spreading of infectious
457 diseases. *Proceedings of the national academy of sciences* **106**(51), 21484–
458 21489 (2009)
- 459 [12] Yang, Y., Herrera, C., Eagle, N., González, M.C.: Limits of predictability
460 in commuting flows in the absence of data for calibration. *Scientific reports*
461 **4**(1), 1–9 (2014)
- 462 [13] Schneider, C.M., Belik, V., Couronné, T., Smoreda, Z., González, M.C.:
463 Unravelling daily human mobility motifs. *Journal of The Royal Society*
464 *Interface* **10**(84), 20130246 (2013)
- 465 [14] Kung, K.S., Greco, K., Sobolevsky, S., Ratti, C.: Exploring universal pat-
466 terns in human home-work commuting from mobile phone data. *PloS one*
467 **9**(6), 96180 (2014)
- 468 [15] Masucci, A.P., Serras, J., Johansson, A., Batty, M.: Gravity versus radia-
469 tion models: On the importance of scale and heterogeneity in commuting
470 flows. *Physical Review E* **88**(2), 022812 (2013)
- 471 [16] Yin, G., Huang, Z., Bao, Y., Wang, H., Li, L., Ma, X., Zhang, Y.: Convgen-
472 rf: A hybrid learning model for commuting flow prediction considering
473 geographical semantics and neighborhood effects. *GeoInformatica*, 1–21
474 (2022)
- 475 [17] Zipf, G.K.: The $p^{-1} p^{-2/d}$ hypothesis: on the intercity movement of
476 persons. *American sociological review* **11**(6), 677–686 (1946)
- 477 [18] Alonso, W.: The system of intermetropolitan population flows. Institute
478 of Urban and Regional Development, University of California (1971)
- 479 [19] Simini, F., González, M.C., Maritan, A., Barabási, A.-L.: A universal
480 model for mobility and migration patterns. *Nature* **484**(7392), 96–100
481 (2012)

- 482 [20] Ren, Y., Ercsey-Ravasz, M., Wang, P., González, M.C., Toroczkai, Z.:
483 Predicting commuter flows in spatial networks using a radiation model
484 based on temporal ranges. *Nature communications* **5**(1), 5347 (2014)
- 485 [21] Stouffer, S.A.: Intervening opportunities: a theory relating mobility and
486 distance. *American sociological review* **5**(6), 845–867 (1940)
- 487 [22] Kotsubo, M., Nakaya, T.: Kernel-based formulation of intervening oppor-
488 tunities for spatial interaction modelling. *Scientific reports* **11**(1), 950
489 (2021)
- 490 [23] Stouffer, S.A.: Intervening opportunities and competing migrants. *Journal*
491 *of regional science* **2**(1), 1–26 (1960)
- 492 [24] Schneider, M.: Gravity models and trip distribution theory. *Papers in*
493 *Regional Science* **5**(1), 51–56 (1959)
- 494 [25] Morton, A., Piburn, J., Nagle, N.: Need a boost? a comparison of
495 traditional commuting models with the xgboost model for predict-
496 ing commuting flows (short paper). In: 10th International Conference
497 on Geographic Information Science (GIScience 2018) (2018). Schloss
498 Dagstuhl-Leibniz-Zentrum fuer Informatik
- 499 [26] Pourebrahim, N., Sultana, S., Niakanlahiji, A., Thill, J.-C.: Trip distri-
500 bution modeling with twitter data. *Computers, Environment and Urban*
501 *Systems* **77**, 101354 (2019)
- 502 [27] Simini, F., Barlacchi, G., Luca, M., Pappalardo, L.: A deep gravity model
503 for mobility flows generation. *Nature communications* **12**(1), 6576 (2021)
- 504 [28] Yao, X., Gao, Y., Zhu, D., Manley, E., Wang, J., Liu, Y.: Spatial origin-
505 destination flow imputation using graph convolutional networks. *IEEE*
506 *Transactions on Intelligent Transportation Systems* **22**(12), 7474–7484
507 (2020)
- 508 [29] Wang, Y., Yin, H., Chen, H., Wo, T., Xu, J., Zheng, K.: Origin-destination
509 matrix prediction via graph convolution: a new perspective of passenger
510 demand modeling. In: *Proceedings of the 25th ACM SIGKDD Interna-*
511 *tional Conference on Knowledge Discovery & Data Mining*, pp. 1227–1235
512 (2019)
- 513 [30] Koca, D., Schmöcker, J.D., Fukuda, K.: Origin-destination matrix esti-
514 mation by deep learning using maps with new york case study. In: 2021
515 7th International Conference on Models and Technologies for Intelligent
516 Transportation Systems (MT-ITS), pp. 1–6 (2021). IEEE
- 517 [31] Rong, C., Feng, J., Ding, J.: Goddag: Generating origin-destination flow

- 518 for new cities via domain adversarial training. *IEEE Transactions on*
519 *Knowledge and Data Engineering* (2023)
- 520 [32] Zhang, J., Zheng, Y., Qi, D.: Deep spatio-temporal residual networks for
521 citywide crowd flows prediction. In: *Proceedings of the AAAI Conference*
522 *on Artificial Intelligence*, vol. 31 (2017)
- 523 [33] Liang, Y., Ouyang, K., Sun, J., Wang, Y., Zhang, J., Zheng, Y., Rosen-
524 blum, D., Zimmermann, R.: Fine-grained urban flow prediction. In:
525 *Proceedings of the Web Conference 2021*, pp. 1833–1845 (2021)
- 526 [34] Zeng, H., Peng, Z., Huang, X., Yang, Y., Hu, R.: Deep spatio-temporal
527 neural network based on interactive attention for traffic flow prediction.
528 *Applied Intelligence*, 1–12 (2022)
- 529 [35] Robinson, C., Dilkina, B.: A machine learning approach to modeling
530 human migration. In: *Proceedings of the 1st ACM SIGCAS Conference*
531 *on Computing and Sustainable Societies*, pp. 1–8 (2018)
- 532 [36] Zhou, F., Li, L., Zhang, K., Trajcevski, G.: Urban flow prediction with
533 spatial–temporal neural odes. *Transportation Research Part C: Emerging*
534 *Technologies* **124**, 102912 (2021)
- 535 [37] NYC Department of City Planning: PLUTO and MapPLUTO. Accessed:
536 4/6/2023. [https://www.nyc.gov/site/planning/data-maps/open-data/
537 dwn-pluto-mappluto.page](https://www.nyc.gov/site/planning/data-maps/open-data/dwn-pluto-mappluto.page)
- 538 [38] Yin, G., Huang, Z., Bao, Y., Wang, H., Li, L., Ma, X., Zhang, Y.: Convgcnr-
539 rf: A hybrid learning model for commuting flow prediction considering
540 geographical semantics and neighborhood effects. *GeoInformatica* **27**(2),
541 137–157 (2023)
- 542 [39] Spadon, G., Carvalho, A.C.d., Rodrigues-Jr, J.F., Alves, L.G.: Recon-
543 structing commuters network using machine learning and urban indica-
544 tors. *Scientific reports* **9**(1), 11801 (2019)
- 545 [40] OpenStreetMap contributors: OpenStreetMap. Accessed: 4/6/2023
546 (2023). <https://www.openstreetmap.org>
- 547 [41] Atwal, K.S., Anderson, T., Pfoser, D., Züfle, A.: Predicting building types
548 using openstreetmap. *Scientific Reports* **12**(1), 19976 (2022)
- 549 [42] Luxen, D., Vetter, C.: Real-time routing with openstreetmap data. In:
550 *Proceedings of the 19th ACM SIGSPATIAL International Conference on*
551 *Advances in Geographic Information Systems*, pp. 513–516 (2011)
- 552 [43] U.S. Census Bureau: LODES data 2015 (<https://lehd.ces.census.gov/>)

- 553 [data/](#)). Accessed: 4/6/2023
- 554 [44] Chen, T., Guestrin, C.: Xgboost: A scalable tree boosting system. In: Pro-
555 ceedings of the 22nd Acm Sigkdd International Conference on Knowledge
556 Discovery and Data Mining, pp. 785–794 (2016)
- 557 [45] Breiman, L.: Random forests. *Machine learning* **45**, 5–32 (2001)
- 558 [46] Hancock, G.R., Freeman, M.J.: Power and sample size for the root mean
559 square error of approximation test of not close fit in structural equation
560 modeling. *Educational and Psychological Measurement* **61**(5), 741–758
561 (2001)
- 562 [47] Chicco, D., Warrens, M.J., Jurman, G.: The coefficient of determination
563 r-squared is more informative than smape, mae, mape, mse and rmse in
564 regression analysis evaluation. *PeerJ Computer Science* **7**, 623 (2021)
- 565 [48] Lenormand, M., Huet, S., Gargiulo, F., Deffuant, G.: A universal model
566 of commuting networks (2012)
- 567 [49] Lenormand, M., Bassolas, A., Ramasco, J.J.: Systematic comparison of
568 trip distribution laws and models. *Journal of Transport Geography* **51**,
569 158–169 (2016)
- 570 [50] Feng, J., Li, Y., Lin, Z., Rong, C., Sun, F., Guo, D., Jin, D.: Context-
571 aware spatial-temporal neural network for citywide crowd flow predic-
572 tion via modeling long-range spatial dependency. *ACM Transactions on*
573 *Knowledge Discovery from Data (TKDD)* **16**(3), 1–21 (2021)
- 574 [51] Zeng, J., Zhang, G., Rong, C., Ding, J., Yuan, J., Li, Y.: Causal learn-
575 ing empowered od prediction for urban planning. In: Proceedings of
576 the 31st ACM International Conference on Information & Knowledge
577 Management, pp. 2455–2464 (2022)
- 578 [52] Viboud, C., Børnstad, O.N., Smith, D.L., Simonsen, L., Miller, M.A.,
579 Grenfell, B.T.: Synchrony, waves, and spatial hierarchies in the spread of
580 influenza. *science* **312**(5772), 447–451 (2006)
- 581 [53] Ferguson, N.M., Cummings, D.A., Fraser, C., Cajka, J.C., Cooley, P.C.,
582 Burke, D.S.: Strategies for mitigating an influenza pandemic. *Nature*
583 **442**(7101), 448–452 (2006)
- 584 [54] Li, M.-H., Chen, B.-Y., Li, C.-T.: A hybrid method with gravity model and
585 nearest-neighbor search for trip destination prediction in new metropol-
586 itan areas. In: 2022 IEEE International Conference on Big Data (Big
587 Data), pp. 6553–6560 (2022). IEEE

- 588 [55] Delventhal, M.J., Kwon, E., Parkhomenko, A.: Jue insight: How do cities
589 change when we work from home? *Journal of Urban Economics* **127**,
590 103331 (2022)
- 591 [56] Bast, H., Storandt, S., Weidner, S.: Fine-grained population estimation.
592 In: *Proceedings of the 23rd SIGSPATIAL International Conference on*
593 *Advances in Geographic Information Systems*, pp. 1–10 (2015)
- 594 [57] Bakillah, M., Liang, S., Mobasheri, A., Jokar Arsanjani, J., Zipf, A.: Fine-
595 resolution population mapping using openstreetmap points-of-interest.
596 *International Journal of Geographical Information Science* **28**(9), 1940–
597 1963 (2014)
- 598 [58] Lundberg, S.M., Lee, S.-I.: A unified approach to interpreting model
599 predictions. *Advances in neural information processing systems* **30** (2017)
- 600 [59] Cai, M., Pang, Y., Sekimoto, Y.: Spatial attention based grid repre-
601 sentation learning for predicting origin–destination flow. In: *2022 IEEE*
602 *International Conference on Big Data (Big Data)*, pp. 485–494 (2022).
603 IEEE
- 604 [60] Lee, M., Holme, P.: Relating land use and human intra-city mobility. *PloS*
605 *one* **10**(10), 0140152 (2015)
- 606 [61] Horner, M.W.: Spatial dimensions of urban commuting: a review of
607 major issues and their implications for future geographic research. *The*
608 *Professional Geographer* **56**(2), 160–173 (2004)
- 609 [62] Credit, K., Arnao, Z.: A method to derive small area estimates of linked
610 commuting trips by mode from open source lodes and acs data. *Environ-*
611 *ment and Planning B: Urban Analytics and City Science* **50**(3), 709–722
612 (2023)
- 613 [63] U.S. Census Bureau: Population data 2010 ([https://data.census.gov/
614 table](https://data.census.gov/table)). Accessed: 4/6/2023

615 Acknowledgments

616 This work is supported by National Science Foundation Grant No. 2109647
617 titled "Data-Driven Modeling to Improve Understanding of Human Behavior,
618 Mobility, and Disease Spread".

619 This project was supported by resources provided by the Office of Research
620 Computing at George Mason University (URL: <https://orc.gmu.edu>) and
621 funded in part by grants from the National Science Foundation (Awards
622 Number 1625039 and 2018631).

623 **Author contributions statement**

624 K.S.A, T.A, A.Z, and D.P. designed the study. K.S.A, T.A, A.Z, and D.P.
625 performed the analyses. K.S.A, T.A, A.Z, and D.P. conceived the experi-
626 ments, K.S.A conducted the experiments. K.S.A, T.A, A.Z, and D.P. wrote
627 and reviewed the manuscript.

628 **Competing interests**

629 The authors declare no competing interests.

630 **Additional information**

631 **Correspondence** and requests for materials should be addressed to K.S.A.