

1 **DECOMPOSING URBAN TRAFFIC FLOWS: A MULTI-STAGE APPROACH TO**  
2 **MODEL HEAVY VEHICLE MOVEMENTS IN GREATER SYDNEY**

3

4

5

6 **Patrick Fernandez**

7 Ph.D. Candidate

8 School of Engineering

9 The Australian National University, Canberra ACT 2601, Australia

10 patrick.fernandez@anu.edu.au

11

12 **Surendra Reddy Kancharla, Ph.D.**

13 Postdoctoral Researcher

14 "Friedrich List" Faculty of Transport and Traffic Sciences

15 Technische Universität Dresden, Saxony, Germany, 01069

16 surendrareddy.kancharla@gmail.com

17

18 **Dung-Ying Lin, Ph.D.**

19 Associate Chairman and Professor

20 Department of Industrial Engineering and Engineering Management

21 National Tsing-Hua University (NTHU), Taiwan

22 dylin@ie.nthu.edu.tw

23

24 **S. Travis Waller, Ph.D., Corresponding Author**

25 Lighthouse Professor and Chair of Transport Modelling and Simulation

26 "Friedrich List" Faculty of Transport and Traffic Sciences

27 Technische Universität Dresden, Saxony, Germany, 01069

28 Professor

29 The Australian National University, Canberra ACT 2601, Australia

30 steven\_travis.waller@tu-dresden.de

31

32

33 Word Count: 4908 words + 1 table(s) × 250 = 5158 words

34

35

36

37

38

39

40 Submission Date: October 29, 2024

**1 ABSTRACT**

2 This study presents an innovative multi-stage methodology for decomposing urban traffic flows  
3 into light vehicle (LV) and heavy vehicle (HV) categories, addressing a critical gap in transporta-  
4 tion network analysis. Utilizing data from Greater Sydney's road network, we develop a compre-  
5 hensive approach comprising three main stages: origin-destination (OD) matrix estimation using  
6 RapidEx [1], quadratic programming optimization for HV/LV proportion estimation, and XGBoost  
7 regression for generalization. Our analysis examines associations between HV trips and various  
8 urban characteristics, including points of interest (POI), nightlight intensity, zonal population, and  
9 zonal area. This approach yields novel insights into urban logistics, commercial activities, and  
10 public transportation patterns. The XGBoost machine learning model, a key component of our  
11 methodology, achieves a test R-squared of 0.637 and identifies critical factors influencing HV  
12 proportions. Key findings include the significant impact of destination characteristics, particu-  
13 larly complex non-linear relationships between nightlight intensity and HV proportions, as well  
14 as, service-related activities and area size, on HV traffic patterns.

15

16 *Keywords:* urban transportation, traffic decomposition, heavy vehicles, light vehicles, origin-destination  
17 flows, machine learning

## 1 INTRODUCTION AND BACKGROUND

2 Urban transportation networks are complex systems that significantly influence economic effi-  
3 ciency, environmental sustainability, and quality of life in cities worldwide. These networks are  
4 shaped by diverse vehicle types and their movement patterns, with a crucial distinction between  
5 Heavy Vehicles (HV) and Light Vehicles (LV). HVs, encompassing trucks, buses, and coaches,  
6 constitute a significant portion of total traffic volume. Trucks usually account for up to 15% of  
7 traffic, while buses may represent approximately 3% [2]. However, these figures exhibit sub-  
8 stantial variability across different regions, influenced by factors such as local industrial activity,  
9 transportation infrastructure, and the extent of public transit systems.

10 The ability to distinguish between LV and HV movements is crucial for several reasons.  
11 It allows for a more detailed understanding of urban logistics and commercial activities, as HVs  
12 are often indicative of goods movement and industrial operations [3, 4]. It also provides insights  
13 into public transportation patterns as buses belong to the HV category. This decomposition by  
14 vehicle type enables more targeted approaches to traffic management, infrastructure planning, and  
15 policy-making. Despite constituting a small proportion of total traffic, HVs contribute disproportio-  
16 nately to various urban challenges. Commercial vehicles significantly impact traffic congestion  
17 and emissions [5, 6]. In terms of infrastructure wear, HVs are responsible for a disproportionate  
18 share of the road maintenance burden [7]. This underscores the need for accurate identification and  
19 disaggregation of HV movements to generate effective policies and manage their impact on urban  
20 areas.

21 Even though the importance of differentiating between vehicle types is recognized, achiev-  
22 ing a fine level of disaggregation in traffic flows has remained a significant challenge in transporta-  
23 tion research. The primary obstacles have been the lack of comprehensive flow data at the link  
24 level and the scarcity of disaggregated flow observations. While the increased use of sensors such  
25 as loop detectors and cameras has provided disaggregated flow data at select locations, estimating  
26 potential flow at every link of a network has been largely unattainable.

27 Recent advancements in transportation modeling tools, such as the RapidEx [1], have cre-  
28 ated new opportunities to address these research gaps. By leveraging travel time data from sources  
29 such as Google or TomTom, in conjunction with road network information from OpenStreetMap,  
30 RapidEx can now estimate potential flow at every link in a network, achieving travel times that  
31 closely align with observed data. This capability, combined with the tool's ability to generate OD  
32 patterns consistent with these flows, provides an opportunity for fine-grained traffic flow disaggre-  
33 gation.

34 However, the application of such advanced tools often faces challenges due to the lack of  
35 detailed, decomposed data needed to generate accurate OD patterns for all links in a network. Even  
36 in developed nations, the availability of such data frequently lacks the necessary granularity and  
37 network coverage [8]. While sensor deployment for obtaining decomposed flow data is increasing  
38 in some areas, it is not yet commonplace worldwide. This data scarcity highlights the need for  
39 models that can be developed in data-rich areas and then applied to other contexts where specific,  
40 detailed transportation data is limited.

41 To address this challenge, understanding the relationship between transportation and eco-  
42 nomic activity becomes crucial. This relationship is fundamental for urban planning and policy-  
43 making [3, 4], as it allows for the creation of transferable models that can estimate traffic patterns  
44 based on more commonly available economic indicators. While this applies to all types of vehicle  
45 movements, it is particularly relevant for HVs, which include both freight transport and public

1 transit. The movement of HVs, especially in freight transport, plays a vital role in bridging the gap  
2 between production and consumption locations. This necessitates the development of models that  
3 integrate economic activity, logistics decision-making, and traffic flows [9, 10]. By focusing on  
4 these relationships, we can create more comprehensive models that account for both LV and HV  
5 movements, providing a better understanding of urban transportation dynamics.

6 In recent years, there has been increased interest in freight transport research, driven by the  
7 need to better understand freight flows and volumes within transportation networks [10, 11]. How-  
8 ever, obtaining reliable freight trip generation data remains challenging, especially where large-  
9 scale freight transport surveys are infrequent [12, 13]. To address these limitations, researchers  
10 have explored various proxies for economic activity in trip generation models [14–18]. Among  
11 these proxies, POI data and nighttime light intensity data have gained significant attention as indi-  
12 cators of economic activity [17, 18]. POI data, representing entities with geolocation information,  
13 has shown strong potential in transport research [19, 20]. These studies demonstrate the effec-  
14 tiveness of such proxies in estimating patterns of socio-economic change and transport demand at  
15 various geographical scales [21].

16 To fully capitalize on these advancements and explore relationships between vehicle move-  
17 ment patterns and urban socio-economic indicators, machine learning techniques have emerged as  
18 a promising approach. In particular, XGBoost has shown impressive results in traffic prediction  
19 tasks [22–24]. Its ability to handle complex non-linear relationships and robustness to overfitting  
20 makes it particularly suitable for transportation data analysis, enabling the development of models  
21 that link socio-demographic factors and other proxies to vehicle type splits.

22 Our research leverages these techniques to develop an innovative approach for decompos-  
23 ing urban traffic flows into LV and HV categories. Focusing on Greater Sydney’s road network,  
24 we aim to identify relationships between vehicle movement patterns and urban socio-economic  
25 indicators. This study is motivated by the need to provide urban planners, policymakers, and trans-  
26 portation engineers with a more comprehensive understanding of the factors influencing HV and  
27 LV movements in urban areas. The significance of this research extends beyond the immediate case  
28 study. By developing a machine learning-based regression model that links socio-demographic fac-  
29 tors and other proxies to vehicle type splits, we create a framework that can be generalized to other  
30 regions or cities. Our approach addresses the scarcity of observed link-level decomposed flow data  
31 and offers a scalable solution for urban transportation analysis on a global scale, moving beyond  
32 traditional survey-based methods to provide a more accurate and transferable understanding of  
33 urban transportation patterns and their relationship to economic activities.

34 Major contributions of our study are as follows:

- 35 1. We develop a novel multi-stage methodology that integrates advanced network model-  
36 ing, optimization techniques, and machine learning to disaggregate traffic flows.
- 37 2. We provide a comprehensive analysis of the factors influencing HV proportions in ur-  
38 ban areas, offering a nuanced view of the complex dynamics governing urban freight  
39 movement.
- 40 3. We establish a generalizable framework for predicting vehicle type distributions across  
41 entire urban networks, enabling more accurate forecasting and scenario analysis for  
42 urban planning in diverse contexts.

43 By bridging the gap between aggregated traffic data and disaggregated vehicle flows, this research  
44 aims to enhance our understanding of urban transportation dynamics and provide practical tools  
45 for evidence-based decision making. The insights gained from this study have the potential to

1 make a noticeable difference towards improving urban mobility efficiency, reducing congestion,  
2 and supporting more sustainable urban development practices across various urban environments.

3 The remainder of this paper is organized as follows: Section 2 describes different data  
4 sources used. Section 3 describes detailed methodology. Section 4 presents results and discusses  
5 their implications for urban planning and transportation policy. Finally, Section 5 concludes with  
6 a summary of our findings and suggestions for future research directions.

## 7 DATA

8 The study leverages three comprehensive datasets that capture various aspects of the urban activity,  
9 and nightlight intensity in Sydney, Australia. These datasets provide a rich and diverse set of  
10 information, enabling a thorough analysis of the relationship between urban characteristics and the  
11 composition of traffic flows.

12 The first dataset, Points of Interest (POI), is sourced from SafeGraph, a leading provider  
13 of geospatial data. SafeGraph is renowned for its comprehensive and accurate POI data, which is  
14 widely used in academic research [25–27], government analyses, and business applications. The  
15 company employs rigorous data collection and validation processes, including both automated and  
16 manual checks, to ensure data quality. The POI data for the year 2023 describes the distribution  
17 of various amenities and services associated with POIs in the Sydney region. This information is  
18 crucial for understanding the spatial patterns of urban activities and their potential influence on  
19 traffic composition.

20 The visible infrared imaging radiometer suite (VIIRS) day/night band (DNB) nighttime  
21 lights data, maintained and provided by the Earth Observation Group<sup>1</sup>, is increasingly recognized  
22 as the gold standard for economic studies utilizing night-lights [28–30]. This dataset captures the  
23 intensity of artificial lighting during nighttime hours, which has been demonstrated to correlate  
24 strongly with economic activity and urbanization levels [31]. For this study, we employed the  
25 Annual VNL V2 [32] average radiance composite data for 2023. This dataset is part of the VIIRS  
26 nighttime light time series, which provides cloud-free average radiance values from nighttime  
27 VIIRS DNB data. The VNL V2 product offers several advantages over its predecessors, including  
28 improved filtering of aurora, stray light, and other ephemeral lights. We selected the 2023 data to  
29 ensure our analysis reflects the most recent available information on nighttime light patterns in our  
30 study area.

31 Lastly, the study incorporates traffic volume data from Transport for New South Wales for  
32 the year 2023. The Traffic Volume Viewer provided by the transport authority splits the traffic  
33 flow for specific links in the road network into HV and LV, offering valuable insights into their  
34 respective proportions and counts [33]. This data source has been utilized in published research  
35 studies, demonstrating its reliability and relevance in traffic analysis [34].

36

37 Table 1 presents summary statistics for key variables in the dataset, encompassing urban  
38 features, population characteristics, and traffic-related metrics. The statistics reveal several note-  
39 worthy patterns across the studied areas. Among urban amenities, services and retail establish-  
40 ments appear to be the most prevalent, with mean counts of 709.50 and 296.83, respectively. The  
41 high standard deviations for these variables suggest considerable variation in their spatial distribu-  
42 tion. Nightlight intensity, used as a proxy for economic activity, exhibits a mean value of 25.51

---

<sup>1</sup><https://eogdata.mines.edu/products/vnl/>

1 with a standard deviation of 24.19, indicating diverse levels of urbanization or economic development across the study area. The population variable also shows substantial variation, with a mean  
 2 of 27,659.27 and a standard deviation of 31,814.99, reflecting the heterogeneous distribution of  
 3 population densities within Sydney.  
 4

**TABLE 1 Descriptive Statistics of Key Variables**

	POI_education	POI_employment	POI_food	POI_medical	POI_retail
mean	21.47	36.07	198.20	184.22	296.83
std	25.61	62.15	221.02	230.21	331.23
min	0.00	0.00	2.00	1.00	4.00
25%	5.00	6.25	53.00	36.00	84.75
50%	11.00	17.50	119.50	101.00	210.00
75%	29.00	36.00	245.75	215.25	345.75
max	161.00	621.00	1455.00	1157.00	2139.00
	POI_services	POI_transit	Nightlight_intensity	Population	Area
mean	709.50	49.60	25.51	27659.27	20.07
std	863.77	53.07	24.19	31814.99	19.43
min	10.00	1.00	0.42	14.72	0.89
25%	224.25	16.00	7.55	7506.12	6.20
50%	403.50	31.00	20.68	16118.90	6.07
75%	857.50	55.75	33.80	33588.74	43.42
max	6868.00	322.00	127.02	157582.34	43.46

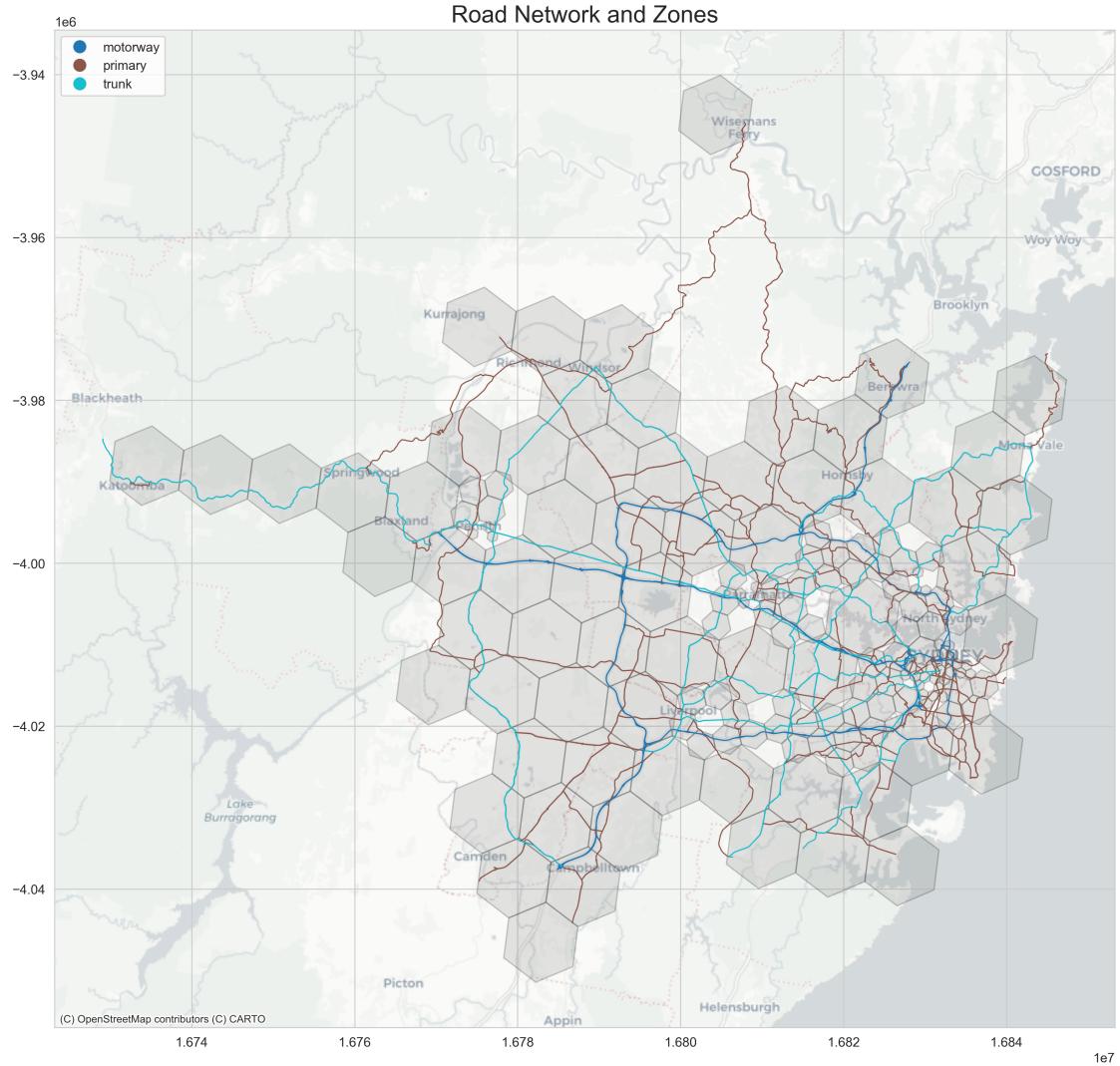
## 5 METHODOLOGY

6 This study presents a novel multi-stage approach to estimate the proportions of HVs and LVs in  
 7 OD matrices for transportation network analysis. Our methodology integrates advanced network  
 8 modeling techniques, optimization algorithms, and machine learning to provide a comprehensive  
 9 solution for disaggregating vehicle types in OD flows. The approach consists of three stages: OD  
 10 matrix estimation, optimization-based proportion estimation, and machine learning-based general-  
 11 ization.

## 12 OD Matrix Estimation

13 The first stage of our methodology focuses on estimating the OD matrix for the study area using  
 14 the RapidEx [1]. RapidEx is a powerful network modeling software that leverages OpenStreetMap  
 15 data to extract road network information that has been successfully applied to analyze similar road  
 16 networks [35]. To manage computational complexity while maintaining essential traffic corridors,  
 17 we focused on motorways, trunk roads, and primary roads. This network simplification strategy  
 18 allows for efficient processing without compromising the integrity of major HV and LV flow pat-  
 19 terns. Travel time data for each link was collected using the Google Maps API for a representative  
 20 weekday morning peak period (June 6, 2024, 7:00-9:00 AM). This time window was selected to  
 21 capture peak traffic conditions. Additionally, we gathered population data from WorldPOP [36]  
 22 and categorized POI information to inform the demand estimation process.

Figure 1 illustrates the road network of Greater Sydney, color-coded as motorways (dark blue), trunk roads (cyan), and primary roads (brown). Secondary and tertiary roads are omitted, as the selected network captures the majority of HV and LV flows. The overlay shows zones (158) based on the Uber H3 hexagonal hierarchical spatial index. These zones start at resolution 6 (approximately  $36 \text{ km}^2$  per hexagon) and are further subdivided up to resolution 8 (approximately  $0.9 \text{ km}^2$  per hexagon) based on POI density.



**FIGURE 1 Road Network and Zones of Greater Sydney**

The RapidEx tool employs a bi-level approach for OD matrix estimation. In the upper level, OD demand is initially estimated using a line search for total demand ranging from 5,000 to 1,500,000 vehicles, utilizing population and POI data proportions to distribute demand among OD pairs. The lower level consists of a traffic assignment problem, determining link-level flows and travel times based on the current OD estimate. A genetic algorithm then operates at the upper level, adjusting the OD demand based on a fitness function that measures the match between estimated and observed traffic flows or travel times at the lower level. This bi-level problem is solved

1 iteratively until a specified convergence criterion is met. For clarity, it's important to note that  
 2 the OD demand estimated by RapidEx is equivalent to Passenger Car Units (PCU), a standardized  
 3 measure of traffic flow accounting for various vehicle characteristics. For a more comprehensive  
 4 understanding of the methodology and its implementation, readers are encouraged to refer to [1].

## 5 Optimization-Based Proportion Estimation

6 The second stage of our methodology focuses on estimating the proportions of HVs and LVs for  
 7 each OD pair using an optimization-based approach. Since the OD demand estimated by RapidEx  
 8 is in PCU and we have the OD contribution to each link (derived during the traffic assignment  
 9 step), different proportions of HVs and LVs will not impact travel times as the PCU value remains  
 10 constant. The number of HVs and LVs can vary while maintaining the same PCU demand, as  
 11 multiple values satisfy the PCU for different HV and LV combinations.

12 The optimization model considers all OD pairs and observed flow [33] at a few links to  
 13 identify the optimal and unique proportions. By leveraging the OD matrix from the previous stage  
 14 and the contribution of each OD pair to link flows, we formulated a quadratic programming opti-  
 15 mization model to determine HV and LV proportions at the OD level. Since demand is constant  
 16 and the value of travel time will not be affected by these proportions, we eliminated the need to  
 17 solve the traffic assignment problem at every iteration, thus making the optimization process ex-  
 18 tremely fast.

19

20 The optimization problem is defined as follows:

21 Objective Function:

$$22 \min \sum_{l \in obs} \left( F_{l,obs}^H - \sum_{(i,j) \in OD} ODL_{i,j,l} \times p_{i,j}^H \right)^2 + \left( F_{l,obs}^L - \sum_{(i,j) \in OD} ODL_{i,j,l} \times p_{i,j}^L \right)^2 \quad (1)$$

23 Subject to:

24

$$25 p_{i,j}^H + p_{i,j}^L = 1 \forall (i,j) \in OD \quad (2)$$

$$26 p_{i,j}^H \in [0, 1] \forall (i,j) \in OD \quad (3)$$

$$27 p_{i,j}^L \in [0, 1] \forall (i,j) \in OD \quad (4)$$

28

29

30 Where  $F_{l,obs}^H$  and  $F_{l,obs}^L$  are the observed flows for HVs and LVs on link  $l$ , respectively.

31  $ODL_{i,j,l}$  represents the demand on link  $l$  contributed by OD pair  $(i,j)$ , and  $p_{ij}^H$  and  $p_{ij}^L$  are the  
 32 proportions of HVs and LVs for OD pair  $(i,j)$ .

33 To solve this optimization problem, we employed Gurobi 11, a commercial optimization  
 34 solver known for its efficiency in handling large-scale quadratic programming problems. To ensure  
 35 a high-quality solution, a relative optimality gap of 1e-6 is set as the stopping criterion.

36 The output of this stage is a set of HV and LV proportions for each OD pair in the study area.

37 These proportions provide a detailed understanding of the composition of traffic flows between  
 38 zones, enabling more accurate modeling and analysis of HV and LV movements in the network.

**1 Machine Learning-Based Generalization**

2 The final stage of our methodology focuses on developing a machine learning model to generalize  
3 the relationship between various socio-economic, land-use, and demographic factors and the  
4 proportions of HVs in OD flows. While the previous stage provides exact proportions at the OD  
5 pair level for the study area, this stage aims to build a model that can be applied to other regions  
6 or time periods where detailed data may not be available. We proposed an XGBoost regression  
7 model [37] to capture the complex relationships between these factors and the proportions of HVs.  
8 The proportion of LVs can be easily derived as the complement of the HV proportion (i.e., 1 - HV  
9 proportion).

10 To train the XGBoost model, we engineered a comprehensive set of features that cap-  
11 ture the factors influencing HV and LV proportions. These features included POI data and its  
12 sub-categories (employment centers, educational institutions, transit hubs, retail establishments,  
13 medical facilities, and service related), which serve as proxies for land-use patterns and economic  
14 activity; nightlight intensity data, an additional proxy for economic activity [38], capturing the in-  
15 tensity of human settlements and industrial areas; population data, representing the distribution of  
16 residential areas and potential trip generation/attraction zones; and zonal area, accounting for the  
17 size and spatial extent of each zone. These features were compiled for both origin and destination  
18 zones of each OD pair, providing a rich set of explanatory variables for the XGBoost model.

19 To improve model quality and reduce noise in the dataset, we applied several preprocessing  
20 steps. OD pairs with zero demand and those with an average annual daily traffic (AADT) for HVs  
21 less than 1 were excluded. All features underwent z-score normalization to prevent bias from  
22 high-magnitude variables. The dataset was partitioned into an 80% training set and a 20% testing  
23 set.

24 To optimize the performance of the XGBoost model, we conducted an extensive hyperpa-  
25 rameter tuning process using grid search. The grid search explored a vast space of 468,750 pa-  
26 rameter combinations, including key parameters such as maximum depth, minimum child weight,  
27 subsample, column sample by tree, learning rate, alpha (L1 regularization), lambda (L2 regu-  
28 larization), and the number of estimators. This exhaustive search allowed us to identify the best  
29 combination of hyperparameters that maximize the model's predictive accuracy and generalization  
30 ability.

31 The trained XGBoost model was then evaluated on the testing set to assess its performance  
32 in predicting HV proportions for unseen OD pairs. To gain further insights into the model's behav-  
33 ior and interpretability, we conducted residual analysis to check for heteroskedasticity and bias,  
34 performed feature importance analysis to identify the most influential factors in determining HV  
35 proportions, and created partial dependence plots to visualize the relationship between key features  
36 and the target variable.

37 The machine learning-based generalization stage offers several benefits for extending the  
38 HV proportion estimates to other regions or time periods. The XGBoost model captures the com-  
39 plex non-linear relationships between the input features and the target variable, allowing for ac-  
40 curate predictions even in the absence of detailed flow data. Furthermore, the model's ability to  
41 handle large datasets and its robustness to outliers and noise make it suitable for application in  
42 diverse urban contexts.

43 The output of this stage is a fully trained XGBoost model that can predict HV proportions  
44 for any OD pair, given the corresponding input features. This model serves as a powerful tool for  
45 transportation planners and policymakers, enabling them to assess the impact of different scenarios

- 1 on the distribution of HVs in the network and make informed decisions regarding infrastructure  
2 investments, traffic management strategies, and environmental policies.

### 3 RESULTS & DISCUSSION

4 The efficacy of our multi-stage approach is demonstrated through rigorous validation and perfor-  
5 mance metrics. The OD matrix estimation showed high accuracy, with link flows within 4% error  
6 on links with available count data, and travel times within 10% of observed values for 85% of  
7 the links. This level of precision in the initial stage provides a solid foundation for subsequent  
8 analyses.

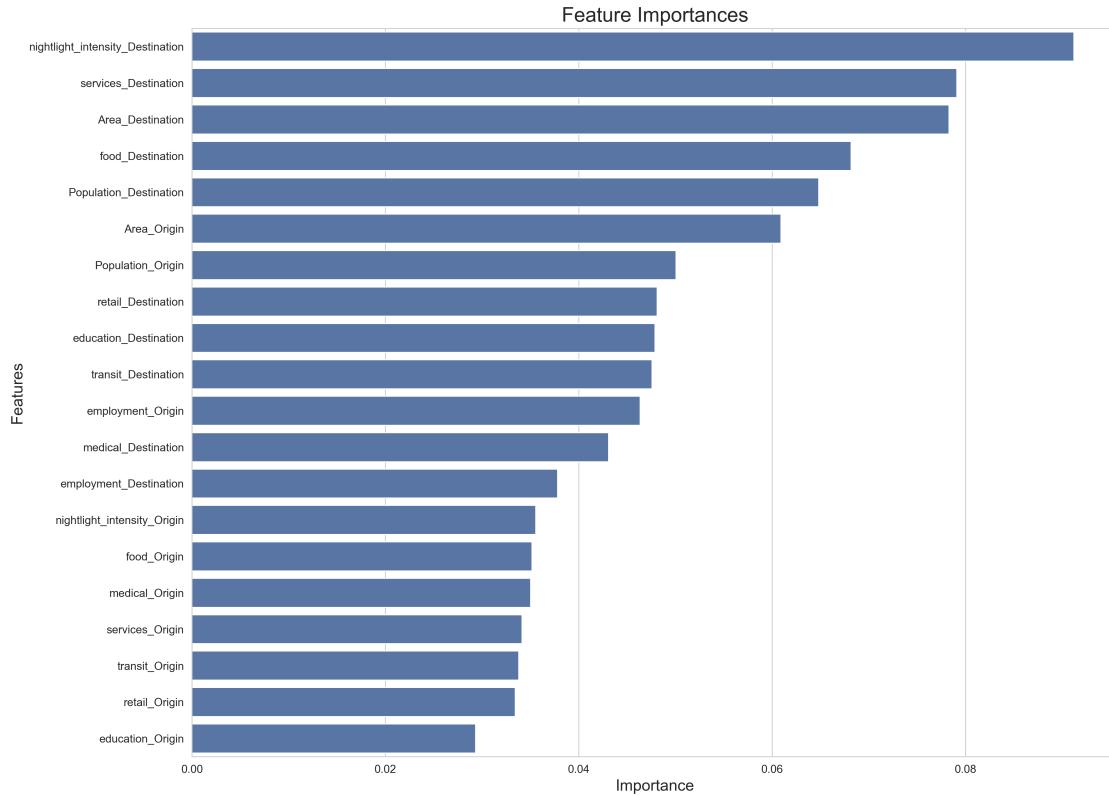
9 Our XGBoost model, optimized through extensive hyperparameter tuning, achieved robust  
10 performance metrics. The optimal configuration included 1200 estimators, a maximum depth of  
11 5, minimum child weight of 4, subsample and column sample by tree both at 0.8, learning rate of  
12 0.06, and regularization parameters alpha and lambda at 0.3 and 9, respectively. This configuration  
13 was determined using five-fold cross-validation, with mean squared error (MSE) as the primary  
14 selection metric, ensuring a balance between predictive accuracy and model generalizability.

15 The final model demonstrated strong predictive capability, achieving a training  $R^2$  of 0.931,  
16 a test  $R^2$  of 0.637, and a cross-validation  $R^2$  of 0.566 (standard deviation: 0.015). The Mean  
17 Absolute Error (MAE) of 0.137 further underscores the model's accuracy. These performance  
18 metrics indicate that our model captures a significant portion of the variance in HV proportions  
19 across the network, providing reliable insights into the factors influencing urban freight movement  
20 patterns.

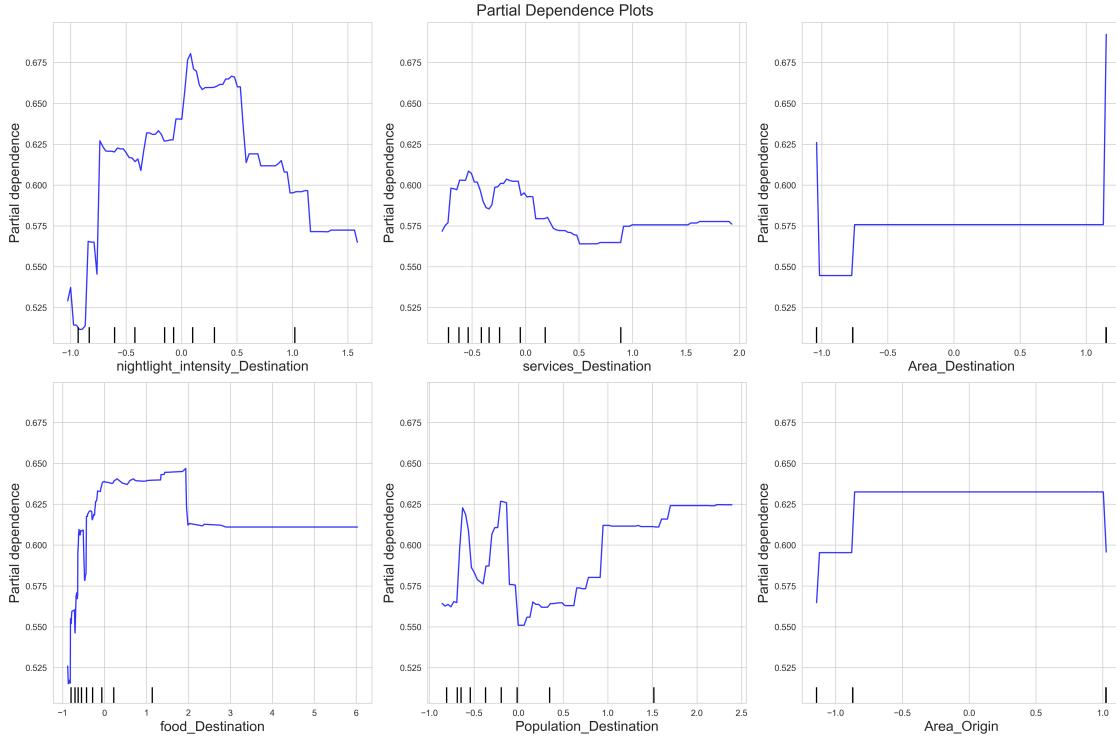
21 Our analysis reveals significant insights into the factors influencing HV proportions in the  
22 origin-destination matrices. The machine learning model proposed demonstrates a robust ability to  
23 capture complex patterns in transportation flows, as evidenced by the feature importance analysis  
24 and partial dependence plots. These findings provide a in-depth understanding of the relationships  
25 between urban characteristics and HV traffic, offering valuable information for urban planning and  
26 transportation policy decisions.

27 The feature importance analysis, presented in Figure 2, identified key drivers of HV pro-  
28 portions. The top five most influential features are nightlight\_intensity\_Destination (0.0912), ser-  
29 vices\_Destination (0.0791), Area\_Destination (0.0783), food\_Destination (0.0681), and Popula-  
30 tion\_Destination (0.0648). This hierarchy of features underscores the critical role of destination  
31 characteristics in determining HV traffic patterns. The presence of services and area-related fea-  
32 tures at both origins and destinations highlights the model's ability to capture the bi-directional  
33 nature of freight transportation. Notably, the inclusion of nightlight intensity as a significant factor  
34 demonstrates the model's capacity to incorporate proxy indicators of economic activity, enhancing  
35 its predictive power.

36 The partial dependence plots, shown in Figure 3, provide valuable insights into the rela-  
37 tionships between specific features and HV proportions. Nightlight\_intensity\_Destination demon-  
38 strates a complex, non-linear relationship with HV proportion. It shows an initial increase, fol-  
39 lowed by a decline, and then another rise. This pattern might reflect the varying nature of nighttime  
40 activities in different urban contexts, with some levels of nightlight intensity corresponding to in-  
41 dustrial areas that generate more HV traffic, while others might indicate residential areas with less  
42 freight movement. The services\_Destination plot shows an initial sharp increase in HV proportion  
43 as service-related activities increase, followed by a more gradual rise and eventual plateau. This  
44 suggests that areas with even a moderate level of service activities generate significantly more HV

**FIGURE 2 Feature Importances**

1 traffic compared to areas with very low service presence. However, beyond a certain threshold, ad-  
 2 ditional service destinations do not substantially increase HV proportions. This pattern aligns with  
 3 the intuitive understanding that service-oriented areas attract freight movement, but there may be  
 4 a saturation point in terms of the logistics needs of these areas. The Area\_Destination plot exhibits  
 5 a step-like pattern, indicating discrete jumps in HV proportion as the destination area increases.  
 6 This could reflect zoning or land-use policies that cluster industrial or logistics activities in specific  
 7 area ranges, leading to abrupt increases in HV traffic at certain thresholds. The food\_Destination  
 8 plot reveals a sharp initial increase in HV proportion, followed by a gradual decline. This could  
 9 indicate that areas with moderate food-related destinations attract more HV traffic, possibly for  
 10 restocking, while areas with very high concentrations might rely more on frequent, smaller deliv-  
 11 eries using light commercial vehicles. Population\_Destination plot displays a fluctuating pattern,  
 12 indicating a complex relationship between population and HV proportion. This could reflect the  
 13 interplay between residential areas and the commercial/industrial zones that generate HV traffic,  
 14 with certain population densities corresponding to mixed-use areas that attract more freight move-  
 15 ment. Finally, the Area-Origin shows a simpler step function, suggesting that the size of the origin  
 16 area has a more straightforward impact on HV proportion, with larger areas generally associated  
 17 with higher HV traffic. This could reflect the tendency of larger origin areas to house more di-  
 18 verse economic activities that generate freight movement. These nuanced relationships highlight  
 19 the complex interplay of urban features in determining HV traffic patterns, underscoring the value  
 20 of our machine learning approach in capturing these intricate dynamics.

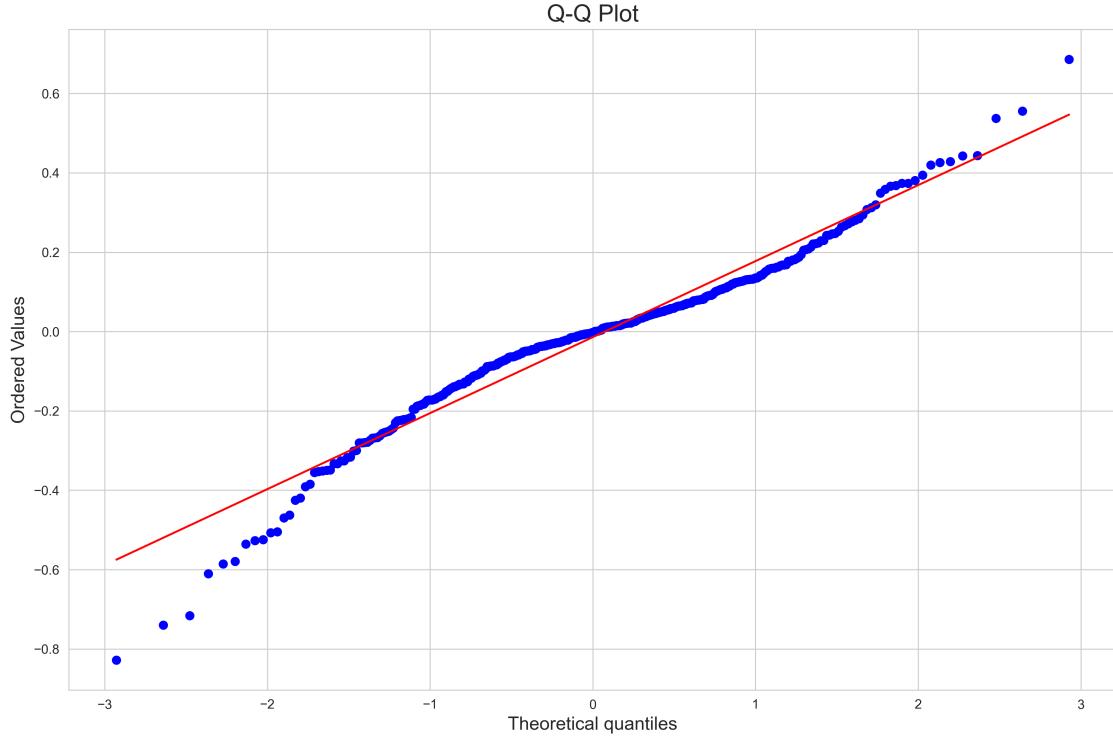


**FIGURE 3 Partial Dependence Plots**

1        The model's performance is further validated by the Q-Q plot presented in Figure 4. This  
 2 plot compares the distribution of the model's residuals to a normal distribution, demonstrating  
 3 a close alignment with the diagonal line. It indicates that the residuals closely follow a normal  
 4 distribution, which is a desirable characteristic for regression models and supports the reliability  
 5 of our predictions.

6        The residual plot, shown in Figure 5, provides additional insight into the model's perfor-  
 7 mance. While the residuals are generally distributed around the zero line, there is a noticeable  
 8 pattern of increasing spread as the predicted values increase. This pattern suggests the presence  
 9 of heteroskedasticity in the model. The observed heteroskedasticity in our model, while notewor-  
 10 thy, should be interpreted within the context of real-world transportation patterns. It is important  
 11 to recognize that large HV proportions are relatively rare in urban transportation networks, with  
 12 LVs typically dominating traffic flows in most areas. This reality of transportation systems has  
 13 significant implications for our model's performance and interpretation.

14        The increasing variability in residuals at higher predicted values likely corresponds to these  
 15 less common, high-HV proportion scenarios. Given the rarity of such cases, the increased vari-  
 16 ability in predictions for these instances is less likely to significantly impact the model's overall  
 17 utility for transportation planning and policy-making. The majority of predictions, which fall in the  
 18 more common lower and moderate HV proportion ranges, demonstrate more consistent accuracy.  
 19 Furthermore, this pattern in the residuals provides valuable insights into the nature of HV traffic  
 20 prediction. The greater variability for high-HV proportion scenarios may reflect the more complex  
 21 and diverse factors influencing these unusual cases. For instance, areas with exceptionally high  
 22 HV proportions might be subject to unique local conditions which introduce additional complexity

**FIGURE 4 Q-Q Plot**

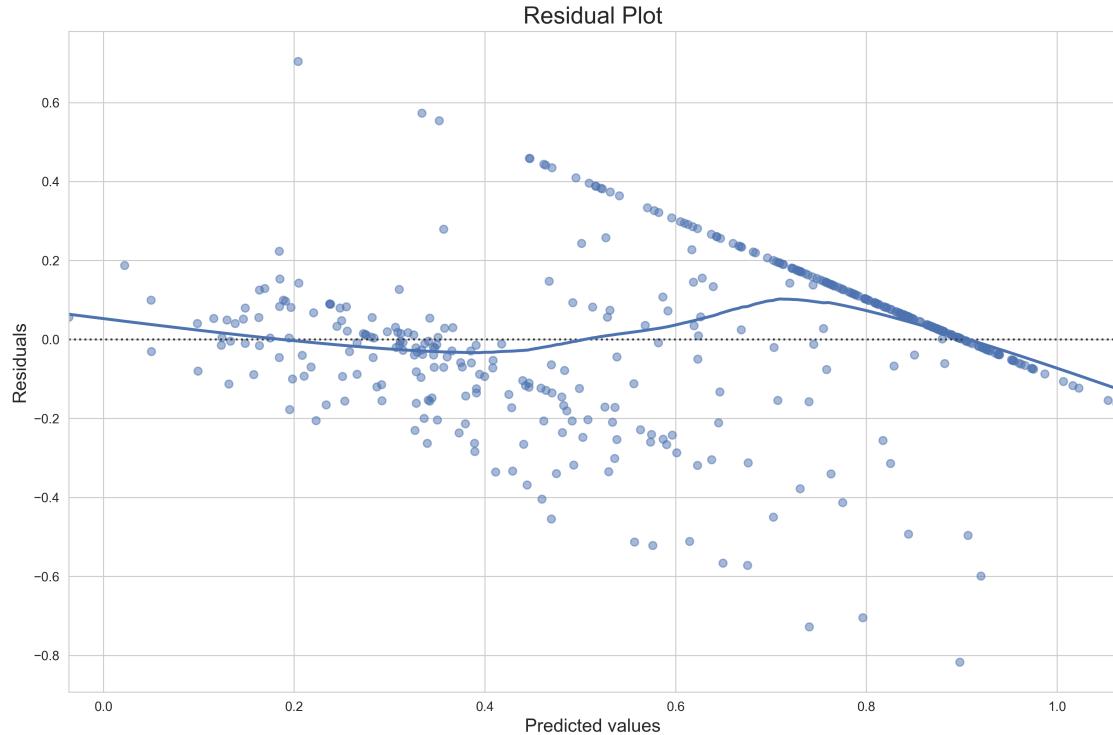
1 not captured by our general set of predictors.

2 Despite this limitation, the model's strong performance in predicting the more common  
 3 scenarios of lower to moderate HV proportions remains significant. The mean of residuals (-  
 4 0.0158) being close to zero indicates minimal overall bias in the predictions, further supporting the  
 5 model's reliability for typical urban transportation contexts. While our model exhibits some het-  
 6 eroskedasticity, particularly for the rare high-HV proportion scenarios, it successfully captures key  
 7 relationships between urban features and HV proportions for the majority of more common cases.  
 8 This characteristic aligns well with the practical needs of urban transportation planning, where  
 9 accurate predictions for typical scenarios are most crucial. The model's performance, combined  
 10 with its ability to highlight areas requiring specialized attention, provides a robust foundation for  
 11 informed decision-making in urban freight management and transportation policy.

## 12 CONCLUSION

13 This study introduces an innovative approach to decomposing urban traffic flows into light vehi-  
 14 cle (LV) and heavy vehicle (HV) categories, addressing a critical gap in transportation network  
 15 analysis. Our three-stage methodology, combining OD matrix estimation, quadratic programming  
 16 optimization, and XGBoost regression, provides a robust framework for estimating HV and LV  
 17 proportions across Greater Sydney's road network. By leveraging diverse data sources, we have  
 18 proposed a comprehensive model that achieves a test  $R^2$  of 0.637, demonstrating its effectiveness  
 19 in capturing complex urban mobility patterns.

20 Our analysis reveals intricate, non-linear relationships between urban features and HV  
 21 proportions. The top five influential features identified are nightlight\_intensity\_Destination, ser-



**FIGURE 5 Residual Plot**

1 vices\_Destination, Area\_Destination, food\_Destination, and Population\_Destination. Notably,  
 2 nightlight intensity at the destination shows a complex pattern with initial increase, decline, and  
 3 subsequent rise, reflecting varying nighttime activities across urban contexts. Services\_Destination  
 4 exhibits a sharp initial increase followed by a plateau, suggesting a saturation point in HV traffic  
 5 generation. Area\_Destination displays a step-like pattern, potentially indicating the impact of zon-  
 6 ing policies on HV traffic. These nuanced relationships underscore the complex interplay of urban  
 7 features in determining HV traffic patterns and highlight the value of our machine learning ap-  
 8 proach in capturing these intricate dynamics.

9 Expanding this framework to diverse global urban contexts would test its generalizability  
 10 and potentially yield insights for creating more equitable and sustainable transportation systems.  
 11 By incorporating socio-economic factors and developing policy recommendations, this approach  
 12 could become a powerful tool for comprehensive urban planning. It would bridge the gap between  
 13 transportation efficiency and social equity in our evolving smart cities, enabling the application of  
 14 Mobility as a Resource (MaaR) concepts. This enhanced model could support decision-makers in  
 15 optimizing urban mobility solutions that are not only efficient but also socially inclusive and en-  
 16 vironmentally sustainable, aligning with the broader goals of smart city initiatives and sustainable  
 17 urban development.

18 Future research could address the limitations identified in this study, particularly focusing  
 19 on improving predictions for high-HV proportion scenarios. Developing specialized models or  
 20 incorporating additional features specific to areas with high HV traffic could enhance the frame-  
 21 work's accuracy across all traffic conditions. This might involve integrating more detailed data  
 22 on industrial activities, zoning regulations, or specific infrastructure designed for HVs. Addition-

ally, extending the model to capture dynamic scenarios by incorporating time-dependent variables could provide insights into daily or seasonal variations in HV proportions.

### 3 AUTHOR CONTRIBUTIONS

4 The authors confirm contribution to the paper as follows: study conception and design: PF, SRK,  
5 STW; data collection: PF, SRK; analysis and interpretation of results: PF, SRK, DYL, STW; draft  
6 manuscript preparation: PF, SRK, DYL, STW. All authors reviewed the results and approved the  
7 final version of the manuscript.

### 8 ACKNOWLEDGEMENTS

9 The authors extend their gratitude to Anna Sotnikova for her assistance in matching observed flows  
10 to the exact links, which significantly contributed to the quality of this research.

### 11 REFERENCES

- 12 1. Waller, S. T., S. Chand, A. Zlojutro, D. Nair, C. Niu, J. Wang, X. Zhang, and V. V. Dixit,  
13 Rapidex: a novel tool to estimate origin–destination trips using pervasive traffic data. *Sus-  
tainability*, Vol. 13, No. 20, 2021, p. 11171.
- 15 2. Evgenikos, P., G. Yannis, K. Folla, R. Bauer, K. Machata, and C. Brandstaetter, Character-  
16 istics and causes of heavy goods vehicles and buses accidents in Europe. *Transportation  
research procedia*, Vol. 14, 2016, pp. 2158–2167.
- 18 3. Dablanc, L., City distribution, a key element of the urban economy: guidelines for practi-  
19 tioners. In *City distribution and urban freight transport*, Edward Elgar Publishing, 2011.
- 20 4. Taniguchi, E., R. G. Thompson, and T. Yamada, Recent trends and innovations in mod-  
21 elling city logistics. *Procedia-Social and Behavioral Sciences*, Vol. 125, 2014, pp. 4–14.
- 22 5. Figliozi, M. A., The impacts of congestion on commercial vehicle tour characteristics and  
23 costs. *Transportation research part E: logistics and transportation review*, Vol. 46, No. 4,  
24 2010, pp. 496–506.
- 25 6. Hartgen, D. T., M. G. Fields, A. L. Layzell, and E. S. Jose, How employers view traffic  
26 congestion: Results of National Survey. *Transportation research record*, Vol. 2319, No. 1,  
27 2012, pp. 56–66.
- 28 7. Song, Y., X. Wang, G. Wright, D. Thatcher, P. Wu, and P. Felix, Traffic volume prediction  
29 with segment-based regression kriging and its implementation in assessing the impact of  
30 heavy vehicles. *IEEE Transactions on Intelligent Transportation Systems*, Vol. 20, No. 1,  
31 2018, pp. 232–243.
- 32 8. Stopher, P. R. and S. P. Greaves, Household travel surveys: Where are we going? *Trans-  
33 portation Research Part A: Policy and Practice*, Vol. 41, No. 5, 2007, pp. 367–381.
- 34 9. Caspersen, E., An explorative approach to freight trip attraction in an industrial urban area.  
35 *City Logistics 3: Towards Sustainable and Liveable Cities*, 2018, pp. 249–268.
- 36 10. Sánchez-Díaz, I., J. Holguín-Veras, and X. Wang, An exploratory analysis of spatial effects  
37 on freight trip attraction. *Transportation*, Vol. 43, 2016, pp. 177–196.
- 38 11. Holguín-Veras, J., M. Jaller, L. Destro, X. Ban, C. Lawson, and H. S. Levinson, Freight  
39 generation, freight trip generation, and perils of using constant trip rates. *Transportation  
Research Record*, Vol. 2224, No. 1, 2011, pp. 68–81.
- 41 12. Pani, A. and P. K. Sahu, Planning, designing and conducting establishment-based freight

- surveys: A synthesis of the literature, case-study examples and recommendations for best practices in future surveys. *Transport Policy*, Vol. 78, 2019, pp. 58–75.
13. Middela, M. S. and G. Ramadurai, Spatial Seemingly Unrelated Regression Models for Freight Trip Generation by Vehicle Type: Application to the Chennai Metropolitan Area in India. *Transportation Research Record*, Vol. 2676, No. 4, 2022, pp. 380–392.
14. Giuliano, G., S. Kang, Q. Yuan, N. Hutson, and M. T. Center, *The freight landscape: using secondary data sources to describe metropolitan freight flows*, 2015.
15. Kang, S., Q. Yuan, N. Hutson, G. Giuliano, and M. T. Center, *The freight landscape: using secondary data sources to describe metropolitan freight flows*, 2010.
16. Alho, A. R. and J. d. A. e Silva, Analyzing the relation between land-use/urban freight operations and the need for dedicated infrastructure/enforcement—Application to the city of Lisbon. *Research in Transportation Business & Management*, Vol. 11, 2014, pp. 85–97.
17. Mellander, C., J. Lobo, K. Stolarick, and Z. Matheson, Night-time light data: A good proxy measure for economic activity? *PloS one*, Vol. 10, No. 10, 2015, p. e0139779.
18. Yan, G., L. Zou, and Y. Liu, The Spatial Pattern and Influencing Factors of China's Night-time Economy Utilizing POI and Remote Sensing Data. *Applied Sciences*, Vol. 14, No. 1, 2024.
19. Jiang, S., A. Alves, F. Rodrigues, J. Ferreira Jr, and F. C. Pereira, Mining point-of-interest data from social networks for urban land use classification and disaggregation. *Computers, Environment and Urban Systems*, Vol. 53, 2015, pp. 36–46.
20. Yuan, J., Y. Zheng, and X. Xie, Discovering regions of different functions in a city using human mobility and POIs. In *Proceedings of the 18th ACM SIGKDD international conference on Knowledge discovery and data mining*, 2012, pp. 186–194.
21. Gao, S., K. Janowicz, and H. Couclelis, Extracting urban functional regions from points of interest and human activities on location-based social networks. *Transactions in GIS*, Vol. 21, No. 3, 2017, pp. 446–467.
22. Zhang, Y. and A. Haghani, A gradient boosting method to improve travel time prediction. *Transportation Research Part C: Emerging Technologies*, Vol. 58, 2015, pp. 308–324.
23. Xu, J., M. Saleh, and M. Hatzopoulou, A machine learning approach capturing the effects of driving behaviour and driver characteristics on trip-level emissions. *Atmospheric Environment*, Vol. 224, 2020, p. 117311.
24. Zafar, N. and I. Ul Haq, Traffic congestion prediction based on Estimated Time of Arrival. *PloS one*, Vol. 15, No. 12, 2020, p. e0238200.
25. Jiao, J., M. Bhat, and A. Azimian, Measuring travel behavior in Houston, Texas with mobility data during the 2020 COVID-19 outbreak. *Transportation letters*, Vol. 13, No. 5-6, 2021, pp. 461–472.
26. Li, Z., H. Ning, F. Jing, and M. N. Lessani, Understanding the bias of mobile location data across spatial scales and over time: a comprehensive analysis of SafeGraph data in the United States. *Plos one*, Vol. 19, No. 1, 2024, p. e0294430.
27. Prestby, T., J. App, Y. Kang, and S. Gao, Understanding neighborhood isolation through spatial interaction network analysis using location big data. *Environment and Planning A: Economy and Space*, Vol. 52, No. 6, 2020, pp. 1027–1031.
28. Bennett, M. M. and L. C. Smith, Advances in using multitemporal night-time lights satellite imagery to detect, estimate, and monitor socioeconomic dynamics. *Remote Sensing of Environment*, Vol. 192, 2017, pp. 176–197.

- 1 29. Gibson, J., S. Olivia, G. Boe-Gibson, and C. Li, Which night lights data should we use in  
2 economics, and where? *Journal of Development Economics*, Vol. 149, 2021, p. 102602.
- 3 30. Price, N. and P. M. Atkinson, Global GDP Prediction With Night-Lights and Transfer  
4 Learning. *IEEE Journal of Selected Topics in Applied Earth Observations and Remote  
5 Sensing*, Vol. 15, 2022, pp. 7128–7138.
- 6 31. Chen, X. and W. D. Nordhaus, Using luminosity data as a proxy for economic statistics.  
7 *Proceedings of the National Academy of Sciences*, Vol. 108, No. 21, 2011, pp. 8589–8594.
- 8 32. Elvidge, C. D., M. Zhizhin, T. Ghosh, F.-C. Hsu, and J. Taneja, Annual time series of  
9 global VIIRS nighttime lights derived from monthly averages: 2012 to 2019. *Remote Sens-  
10 ing*, Vol. 13, No. 5, 2021, p. 922.
- 11 33. New South Wales Roads and Maritime Services, *Traffic Volume Viewer*.  
12 [http://www.rms.nsw.gov.au/about/corporate-publications/statistics/  
13 traffic-volumes/aadt-map/index.html#z=6](http://www.rms.nsw.gov.au/about/corporate-publications/statistics/traffic-volumes/aadt-map/index.html#z=6), 2023, [Online; accessed 6-June-  
14 2024].
- 15 34. Infante, W. and J. Ma, Clustering and Previous Visit Dependency Technique for Electric  
16 Vehicle Station Visits. In *2018 IEEE PES Innovative Smart Grid Technologies Conference  
17 Europe (ISGT-Europe)*, 2018, pp. 1–5.
- 18 35. Waller, S. T., M. Qurashi, A. Sotnikova, L. Karva, and S. Chand, Analyzing and modeling  
19 network travel patterns during the Ukraine invasion using crowd-sourced pervasive traffic  
20 data. *Transportation research record*, Vol. 2677, No. 10, 2023, pp. 491–507.
- 21 36. Bondarenko, M., D. Kerr, A. Sorichetta, and A. J. Tatem, *Census/projection-disaggregated  
22 gridded population datasets for 189 countries in 2020 using Built-Settlement Growth  
23 Model (BSGM) outputs*. <https://www.worldpop.org/>, 2020.
- 24 37. Chen, T. and C. Guestrin, Xgboost: A scalable tree boosting system. In *Proceedings of the  
25 22nd ACM SIGKDD international conference on knowledge discovery and data mining*, 2016,  
26 pp. 785–794.
- 27 38. Elvidge, C. D., K. E. Baugh, S. J. Anderson, P. C. Sutton, and T. Ghosh, The Night  
28 Light Development Index (NLDI): a spatially explicit measure of human development  
29 from satellite data. *Social Geography*, Vol. 7, No. 1, 2012, pp. 23–35.