



HAL
open science

Scaling traffic variables from sensors sample to the entire city at high spatiotemporal resolution with machine learning: applications to the Paris megacity

Xavier Bonnemaizon, Philippe Ciais, Chuanlong Zhou, Simon Ben Arous, Steven J Davis, Nicolas Megel

► To cite this version:

Xavier Bonnemaizon, Philippe Ciais, Chuanlong Zhou, Simon Ben Arous, Steven J Davis, et al.. Scaling traffic variables from sensors sample to the entire city at high spatiotemporal resolution with machine learning: applications to the Paris megacity. *Environmental Research: Infrastructure and Sustainability*, 2024, 4 (3), pp.035010. 10.1088/2634-4505/ad6bbf . hal-04688444

HAL Id: hal-04688444

<https://hal.science/hal-04688444v1>

Submitted on 5 Sep 2024

HAL is a multi-disciplinary open access archive for the deposit and dissemination of scientific research documents, whether they are published or not. The documents may come from teaching and research institutions in France or abroad, or from public or private research centers.

L'archive ouverte pluridisciplinaire **HAL**, est destinée au dépôt et à la diffusion de documents scientifiques de niveau recherche, publiés ou non, émanant des établissements d'enseignement et de recherche français ou étrangers, des laboratoires publics ou privés.

PAPER • OPEN ACCESS

Scaling traffic variables from sensors sample to the entire city at high spatiotemporal resolution with machine learning: applications to the Paris megacity

To cite this article: Xavier Bonnemaizon *et al* 2024 *Environ. Res.: Infrastruct. Sustain.* **4** 035010

View the [article online](#) for updates and enhancements.

You may also like

- [Inequalities in the production and use of cement and concrete, and their consequences for decarbonisation and sustainable development](#)
Alastair T M Marsh, Rachel Parker, Anna L Mdee et al.
- [Repurposing coal plants into thermal energy storage—a techno-economic assessment in the Indian context](#)
Serena Patel, Dharik Mallapragada, Karthik Ganesan et al.
- [Carbon storage in the built environment: a review](#)
Stavroula Bjānesøy, Antti Kinnunen, Hulda Einarsdóttir et al.

ENVIRONMENTAL RESEARCH
INFRASTRUCTURE AND SUSTAINABILITY

PAPER



Scaling traffic variables from sensors sample to the entire city at high spatiotemporal resolution with machine learning: applications to the Paris megacity

OPEN ACCESS

RECEIVED
4 March 2024REVISED
12 July 2024ACCEPTED FOR PUBLICATION
6 August 2024PUBLISHED
28 August 2024

Original content from this work may be used under the terms of the [Creative Commons Attribution 4.0 licence](#).

Any further distribution of this work must maintain attribution to the author(s) and the title of the work, journal citation and DOI.

Xavier Bonnemaizon^{1,*} , Philippe Ciais¹ , Chuanlong Zhou¹, Simon Ben Arous², Steven J Davis³ and Nicolas Megel⁴¹ Laboratoire des Sciences du Climat et de l'Environnement, LSCE/IPSL, CEA-CNRS-UVSQ, Université Paris-Saclay, Gif-sur-Yvette, France² Kayrros SAS, Paris, France³ Department of Earth System Science, University of California Irvine, Irvine, CA, United States of America⁴ NEXQT SAS, Paris, France

* Author to whom any correspondence should be addressed.

E-mail: xavier.bonnemaizon@lsce.ipsl.fr**Keywords:** road transportation, traffic monitoring, carbon dioxide emissions, COVID-19, bottom-up approach, machine learning, urban areasSupplementary material for this article is available [online](#)**Abstract**

Road transportation accounts for up to 35% of carbon dioxide and 49% of nitrogen oxides emissions in the Paris region. However, estimates of city traffic patterns are often incomplete and of coarse spatio-temporal resolution, even where extensive networks of sensors exist. This study uses a machine learning approach to analyze data from 2086 magnetic road sensors across Paris, generating a detailed dataset of hourly traffic flow and road occupancy covering 6846 road segments from 2018 to 2022. Our model captures flow and occupancy with a symmetric mean absolute percentage error of 37% and 54% respectively, providing high-resolution insights into traffic patterns. These insights allow for the creation of a comprehensive map of hourly transportation patterns in Paris, offering a robust framework for assessing traffic variables for each significant road link in the city. The model's ability to incorporate an emission factor based on the mean speed of the vehicle fleet, derived from flow and occupancy data, holds promise for developing a detailed CO₂ and pollutant inventory. This methodology is not limited to Paris; it can be applied to other urban centers with similar data availability, highlighting its potential as a versatile tool for sustainable urban monitoring.

1. Introduction

The road transport sector is a major contributor to greenhouse gas emissions in France [1]. Despite efforts to reduce its emissions in recent decades, the sector has seen little improvement, with the COVID-19 crisis providing a temporary drop followed by a rebound. With the European Parliament announcing a ban on the sale of carbon dioxide (CO₂)-emitting vehicles by 2035 [2], the need for decarbonization of the transport sector has become even more pressing. Cities are hotspots of traffic emissions and, in addition to their long term CO₂ emissions reduction goals and policies, they face air-quality challenges related to pollutants (typically nitrogen oxides (NO_x), carbon monoxide (CO) and particulate matter (PM)) co-emitted with CO₂ by on-road combustion engines [3, 4]. In this context, we focus on mapping traffic variables that control transport emission in the mega-city of Paris. According to the city official statistics, its road traffic represents 3018 million of vehicle.km with an associated 1.15 megatons of CO₂ emitted in 2018 [5]. The reduction of vehicle.km and emissions were respectively about 26% and 36% compared to 2004, the other relative reduction of emissions being explained by improved motorisation of the city's vehicle fleet.

To better understand how emissions from road transport are controlled by traffic variables, methods that go beyond city-wide average estimates are needed. In particular, high temporal resolution maps of variables are required. The IPCC recommends the use of top-down methods based for instance on fuel sales for national inventories [6], but such methods only give city-wide changes, may fail to capture sudden perturbations like the COVID-19 crisis or local changes, and can be problematic for cities where vehicles buying fuel outside emit in the city area. The COPERT methodology [7] suggested that approaches based on vehicle-kilometers and traveling speeds are preferable to methods based on fuel consumption. Pinto *et al* [3] provide a large review of traffic variables estimation used for computing on-road transportation emissions, distinguishing top-down and bottom-up approaches, static and dynamical application of an emission factor. They stressed the need of having high-quality and local dynamic inventories to support policymakers to develop relevant strategies for reducing emissions. Here we present a new bottom-up approach that aims to provide continuous maps of traffic flow and occupancy at hourly time step for the main streets of the city of Paris by upscaling point scale measurements collected from road sensors at a limited number of locations [8].

2. Previous studies

2.1. Mobility and traffic variables estimation

Mobility data can be gathered from sources like surveys, *in-situ* sensor measurements [9], and activity proxies like geolocation data of individual people or vehicles. Lenormand *et al* [10] explored the connections between regular surveys, cell phone data, and Twitter (X) posts, and found potential biases in Twitter data related to user age. In a Dallas-specific study, Xu *et al* [11] combined location-based data and surveys to investigate how vehicle drivers respond to traffic congestion. While surveys and census data cover a large area with little detail, they remain valuable for understanding movement patterns over extensive regions, as demonstrated in Île-de-France by Hörl and Balac [12]. During the COVID-19 pandemic, big data approaches relying on geolocation however demonstrated considerable promise to understand sudden mobility changes [13].

Sensor-based *in-situ* measurements of vehicle speed, numbers and even types from cameras or counting devices provide direct measurements of key traffic variables—such as flow, occupancy, density, and speed—on specific road segments at a specific point, but they do not cover all the streets. In less-monitored areas, where sensors are scarce, survey data proves complementary for estimating traffic [14]. Moreover, some types of sensors such as inductive loop detectors may produce non-representative data for instance during lane closures or openings. Additionally, sensor maintenance problems or malfunctions can lead to gaps in the data, emphasizing the need for methods to gap-fill sensor data and predict traffic variables on unmonitored road segments. Tarunesh and Chung [15] proposed the use of neural networks to address these issues. Xing *et al* [16] conducted a comprehensive review of methods for predicting missing traffic information, citing various machine learning examples for this predictive task. They proposed to distinguish three categories of research applications: the estimation of traffic under different scenarios of missing data, fusion with different types of detectors, and use of different data types (for instance mobile phone data and GPS). Our study falls in the first category as we aim to scale up traffic variables on non monitored roads with *in-situ* sensor-based measurements at fixed locations.

2.2. From traffic variables to road transport emissions

The estimation of emissions from road transportation with traffic variables has been a subject of extensive research, in particular with the emergence of the use of new activity proxies such as geolocation datasets. Uncertainties arise from the fact that some activity data are not directly related to traffic as highlighted in recent studies [17]. For instance, Guevera *et al* [18] utilized a combination of Google Mobility Index [19] which describes people's time spent in different locations rather than road traffic, along with reports from national transport agencies, in order to compute changes in pollutant emissions on a daily and country level during the COVID-19 crisis in Europe. Huo *et al* [20] estimated CO₂ emissions in many cities for road transportation, using a daily TomTom congestion index data [21] averaged at city scale, without information about the area being covered by the index and with a simple model calibrated only for one city using aggregates of sensor data.

Biswal *et al* [22] computed pollutant emissions at an hourly level in Delhi, with a speed-flow traffic model and car flow measurements at 72 locations between 08:00–14:00. The COPERT [23] method was used to determine mean speed dependent emission factors for pollutants. Li *et al* [24] applied the MOVES [25] model in Beijing for these emission factors with a speed-flow model that takes into account Greenshields Speed-Density hypothesis. These studies illustrate the diversity of bottom-up methodologies to estimate emissions from road transportation sensors, but they also highlight the need for careful consideration of data sources, spatial scales, and modeling techniques to reduce uncertainties in emissions inventories at city and road link scale.

2.3. The city of Paris and its road sensors data used in this study

The city of Paris has a dense network of roads and high vehicle traffic with a peripheral motorway around the entire city and intra-city roads with different importances. Since 2002, the local mobility agency using this sensor data has reported a significant decrease of vehicle-kilometers [26]. Nevertheless, both pollution and greenhouse gasses emissions remain a concern for the city as new measures are being instituted to further limit vehicle usage [27]. In the city of Paris like in many other cities, changes in mobility due to COVID-19, new remote working habits, the construction of recent infrastructure such as cycle paths and low emission zones have changed traffic patterns since 2018, and yet, these changes are poorly understood [28] because of the lack of continuous traffic data covering all roads.

The city of Paris traffic monitoring system is based on a network of magnetic sensors beneath the roads, measuring two key parameters of vehicle traffic: the flow (Q , number of vehicles per hour) and the occupancy (O , the percentage of time when the sensor was covered by a vehicle). The data is available through the city open data platform [8] with hourly aggregation and is used by *Poste Central d'exploitation Lutèce* to monitor the traffic and plan infrastructure works. This network of sensors covers 3350 road segments and only samples a small subset of the 14 010 total road segments reported by OpenStreetMap [29] (OSM). Therefore, a model is needed to upscale the sparse sensors' data to map space and time patterns of Q and O within the entire city.

Agence de la mobilité is using the information from this network of sensors to release quarterly bulletins on mobility habits [26]. *Air Paris*, the regional air quality agency, also uses sensor data to calibrate [30] an emission model for monitoring air quality in the area. Combined with Uber speed data [31], this dataset was also employed by Mahajan et al [32] who used a transfer learning method to predict the road flow from Paris data in the city of Madrid in 2019. The Uber speed data is no longer available since October 2023 [33]. The lack of continuous availability of activity proxies such as from Uber and Google stresses the necessity to develop robust methods based on publicly available data to build a historical traffic inventory and derive the associated emissions.

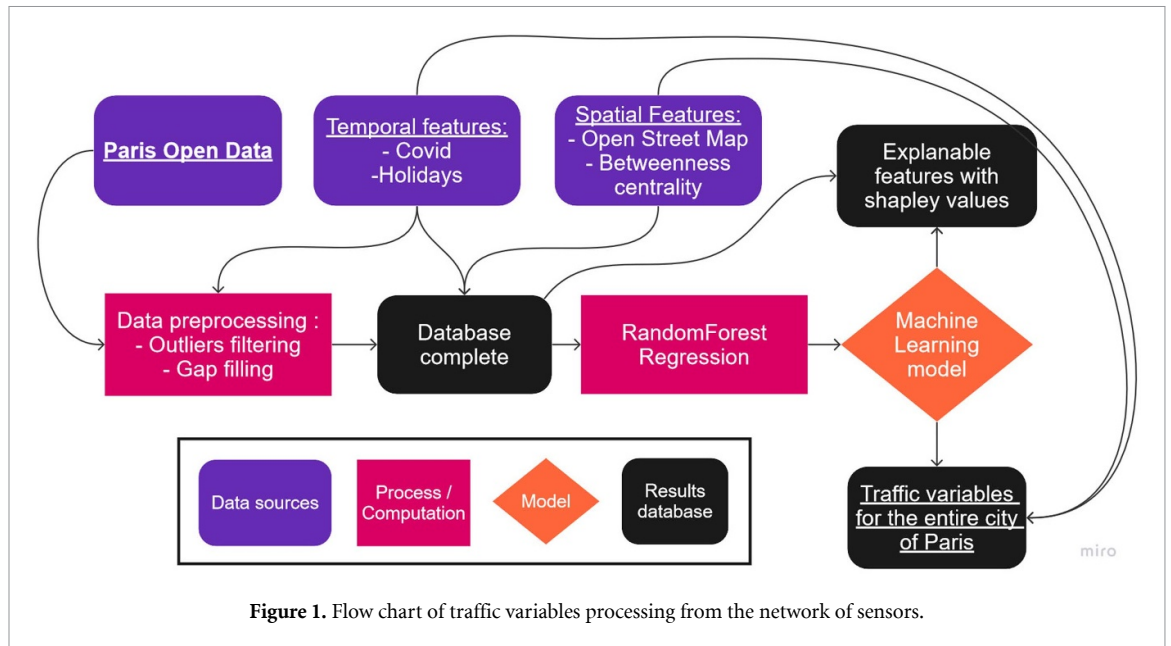
3. Research question

Although some of the previous studies used only one traffic variable like the flow of vehicles as a predictor for traffic emission change, combining flow (Q) and Occupancy (O) can be used to compute speed. Intuitively, the duration of vehicles' presence on a sensor directly reflects their density or concentration. When the concentration of vehicles increases, more vehicles pass through a sensor, resulting in longer time coverage readings. For this reason, the Occupancy is a proxy to measure the concentration level K , the number of vehicles per kilometer, that leads to congestion. The two variables O and K share a proportional relationship if we assume that the length of vehicles is uniform [34] (see supplementary information 3). Subsequently, because the ratio between Q and K provides the mean harmonic speed of the vehicle fleet on the road segment containing a sensor during hourly periods, it follows that the ratio of Q and O exhibits a near proportional relationship with the spatial average speed on a road segment (the relation between these variables is illustrated figure 3). Quantifying the relationship between Q and O is therefore essential to subsequently derive information on vehicle speed, which is needed to calculate emissions.

In the following, we aim to combine the road link scale sensor data from Paris available only on a limited number of roads with a new machine learning model to derive a high-resolution road-level mapping of Q and O for all the main roads of the whole city of Paris at hourly scale. This approach is particularly valuable to understand how perturbations may have a different impact on traffic between road types, such as the perturbations experienced by the city due to COVID and recent limitations on car usage. We will leave the construction and application of a relevant carbon dioxide emission factor for future work. Our approach aims to predict variables on main road segments of the urban network. Extrapolating traffic patterns to smaller roads not covered by the sensors used to train the model may introduce bias. To avoid this, additional monitoring data focused on these smaller road segments is needed. However, according to our discussion with PC Lutèce, such data is not publicly available. Consequently, our model is derived only for about 50% of the 14 010 road segments and will not cover small living streets in dense residential districts, for which a low traffic is expected.

Our first research question is 'How can we derive the hourly traffic flow and occupancy for all main road segments of a city from sparse local sensor data?'. This question deals with data extrapolation and homogeneity within a city. We aim at predicting flow and occupancy on an hourly and road-link level using the network of point-scale sensor data in Paris [8] for the 2018–2022 period. Our model is built to learn the spatio-temporal patterns of these variables using exogenous features such as hours or lanes number.

Our second research question is: 'How can Machine Learning models be used as a tool to find relevant attributes explaining changes in traffic variables?'. By providing insights on the influence of various attributes



to predict the Q and O , our machine learning model can then be used in an explanatory way to help understand the relative importance of time-related features such as holidays or health crises, as well as the effect of spatial features like infrastructure. This information can be later used for policymakers as it brings knowledge about influential variables and how they could be adjusted for possible solutions to mitigate emissions.

In the following, we present our methodology for preprocessing and analyzing traffic sensor data (section 2), evaluate the method performances, and discuss our main findings in terms of model explanative features (section 3). We conclude with a discussion of the implications of our findings for policymakers and future research directions.

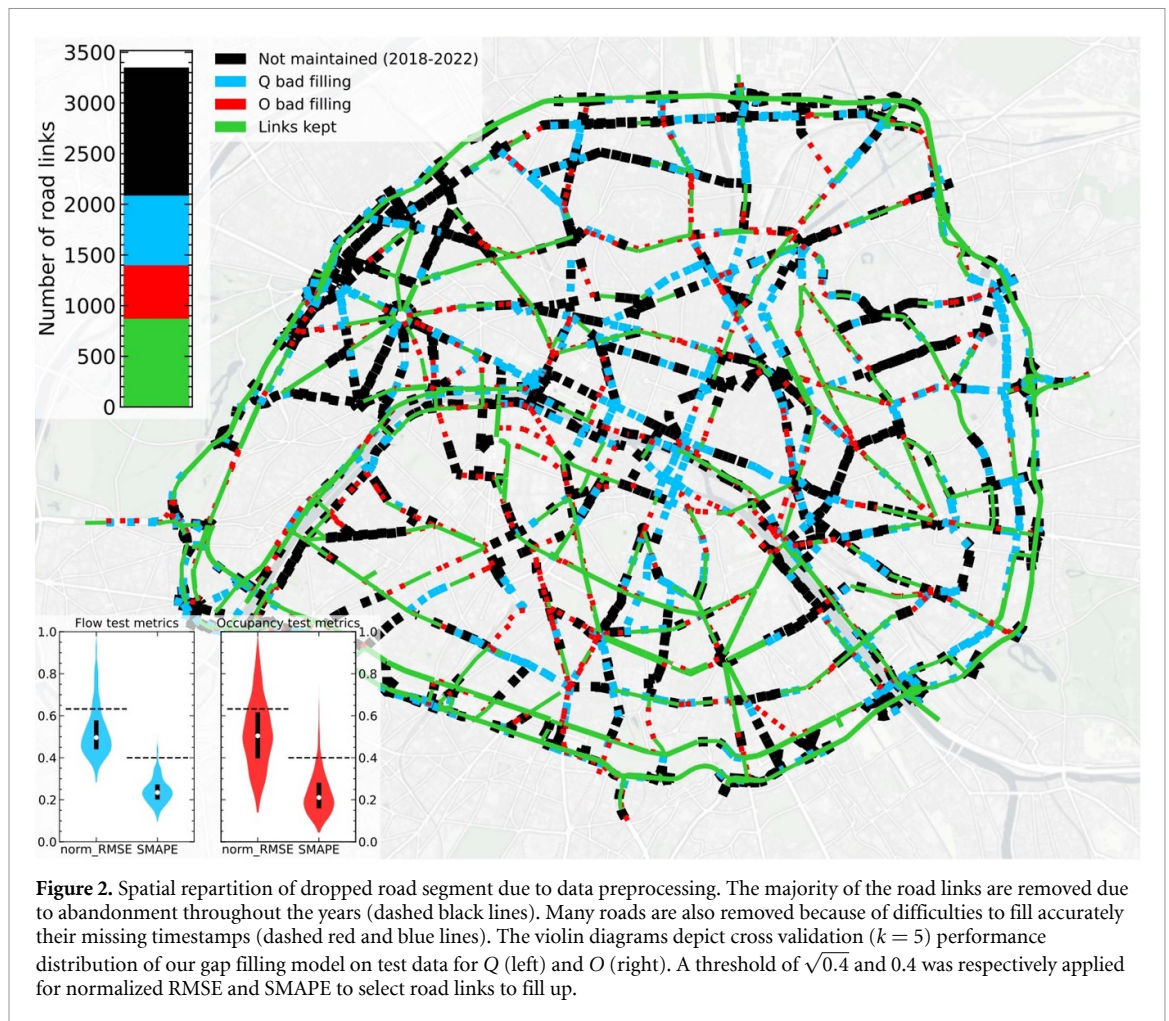
4. Data and methods

Figure 1 provides a graphical summary about the methodology of this study. The Paris open data sensor data is filtered and gap filled in order to be used as an input of a machine learning model predicting traffic on main roads and for explaining observed changes.

4.1. Data presentation and preprocessing

The sensor data from the Paris open data platform [8] comes from magnetic sensors beneath roads and monitors the hourly city traffic in real time. These data have outlier values and missing timestamps (point wise missing but mostly linear wise missing [16]) due to maintenance, unusual behavior like a vehicle parking on a sensor, or construction works that destroy the magnetic loop. The city of Paris is slowly converting those magnetic sensors to cameras that can be more reliable as well as more precise with distinction between vehicle types. Not all the sensors are maintained on a permanent basis: only 2086 road links were present in the entire 2018–2022 dataset out of the 3350 initial locations with a geographical reference. Moreover, sensors come with a variable rate of missing values, as described in Supplementary Information section 2. We observed in figure 13 showing the data density that there are less gaps for peripheral roads (located in the outer ring) compared to roads within the city, which might be explained by the involvement of another entity responsible for maintenance and monitoring of these peripheral road sensors. We therefore developed a method to remove outliers and fill missing data rather than dropping out a sensor with low coverage or many outliers for the whole period, in order to achieve better results and generalization for our machine learning traffic upscaling model (figure 1).

Outliers are identified by analyzing the distribution of Q and O for each road link and hour. Specifically, for Q , the distribution is fitted using a normal distribution, and values outside the 0.5th and 99.5th percentile range are excluded. Regarding the occupancy variable (O), a Gumbel distribution [35] is fitted to the data, and only values within the 3rd and 97th theoretical percentiles are retained. Data filtering based on the choice of these distributions gives the best results in terms of average maximum likelihood for fitting the corresponding data. The years 2018, 2019, and 2022 were treated separately from 2020 and 2021 that were strongly impacted by the COVID-19 crisis. For those two years only, values above the high threshold were



eliminated. As a result, our removal of outliers filtered 1.5% and 5% of the available data for Q and O respectively.

4.2. Filling missing timestamps

Both Q and O traffic variables in Paris exhibit temporal variations that are predominantly driven by a diurnal cycle of activities. Q tends to have a relatively flat peak throughout the day, while O displays sharper peaks during the morning and evening rush hours on busy days, especially for peripheral road, a highway that circles around the city (supplementary information section 1 figure 12). Other temporal features such as weekdays, seasons, holidays, and crises also show up in the temporal behavior of both traffic variables (see for instance figure 4) but to a lesser extent. The magnitude of diurnal changes is dependent on road link-specific characteristics and significance.

It is also clear that the variability illustrated in figure 12 is higher among intracity roads than the peripheral, indicative of their more diverse traffic patterns. Furthermore, these roads may have been subject to various interventions during the time range of the study (e.g. building of bus lanes, cycle paths, closures), influencing traffic conditions and potentially resulting in increased variability of traffic.

To fill the data gaps, we used random forest regression models that include as predictors hour, weekday, month, year, and the average stringency index during the COVID-19 crisis years [36]. We built one random forest model for each road link. Firstly, we performed the prediction of missing values for the flow (Q) variable. Subsequently, we extended this prediction to the occupancy (O) variable, utilizing the flow variable as a feature. Our results (section 3 figure 2) demonstrate a significant enhancement in the performance of occupancy prediction by incorporating the flow variable as a feature.

4.3. Spatio-temporal features used as predictor of traffic variables

For predicting Q and O across all the road links with machine learning models (figure 1), we chose a number of relevant features, which are described below and summarized in table 1. We have in total 11 temporal features including COVID-19 stringency, and 4 spatial features.

Table 1. Summary of features used to train the model.

Type of feature	Description	Source
Temporal attributes	Hour	Computed from timestamp value
	Week day	
	Month	
	Year	
COVID-19	Stringency index	[36]
Holidays	French bank holidays	[37]
	French school holidays (5 types)	
OpenStreetMap attributes	Lanes: lanes number	[29]
	Speed_kph: speed limitation (km/h)	
	Highway: indicates the relative importance of the road link	
Betweenness Centrality	Normalized number of shortest paths going through the edge	[38]

4.3.1. Characteristics of each road

Three characteristics of each road link used to predict Q and O are the number of lanes, the speed limitation, and the road type category (table 1). Unfortunately, no information on such characteristics comes from the Paris open data platform apart from a line string geometry. Thus, the geometry characteristics were matched with an independent road geospatial database to obtain the characteristics of each road link. As in Mahajan *et al* [32], we matched the geometries to OpenStreetMap [29] data using PyTrack [39] rather than SharedStreets [40] which seems to be no longer maintained. This Python toolkit uses Hidden Markov Model [41] to detect the most probable path of points retrieved from the Paris open data geometry through the OSM network. An example of the map matching is represented in supplementary information section 5. As a result, we obtain valuable characteristics, including road category, lane counts, and speed limits for each road segment. Road segments which lack values for the lanes parameter are filled in with the mean number of lanes for their respective road category.

4.3.2. Betweenness centrality

The betweenness centrality metric [38] quantifies the frequency with which an edge is utilized to connect two nodes by their shortest path within the network. An edge (here road link) of a network (here road network) is one of the connections between the nodes (here road intersections) of the network. We hypothesized that this feature can help to predict Q and O and computed it for each road link utilizing the OSM network. Each road link was weighted using the theoretical time to cross it using the speed limitation and length attributes from OSM. Roads with high betweenness centrality serve as critical connecting points between different parts of the network, making them attractive routes for vehicles to cross [42]. To ensure fair evaluation, we employed a 500-meter buffer around the city, thereby avoiding potential bias against outer edges, including the peripheral areas to calculate the betweenness centrality as given by:

$$BC(e) = \sum_{u \neq v} \frac{\sigma_{uv}(e)}{(N-1)(N-2)}$$

u, v are a couple of nodes from the graph.

$\sigma_{uv}(e)$ returns 0 if the shortest path from u to v does not cross the edge e and else 1.

N the total number of nodes.

Equation (1): normalized betweenness centrality definition.

4.3.3. Temporal profiles of human activity

Days off and periods of leave have an impact on people's mobility and were used as a predictor (table 1). Holidays were separated into bank holidays when most businesses and administrations are closed (11 d per year in France) and school holidays (5 different periods per year) when some households may stop working but not necessarily. This information was retrieved from the open data platform of the government [37] and used as a predictor for mapping the traffic across all Paris' streets.

4.3.4. COVID-19 induced changes in traffic variables

During the course of our analysis period, the COVID-19 crisis strongly affected France and Paris. The government took strong measures such as curfew, business and school closures, as well as lockdowns for

instance from 17th March to 11 May 2020, which strongly impacted people's mobility and reduced traffic patterns. To investigate the COVID-19 perturbation, we use the database of governments response from Hale *et al* [36] which provides a government stringency index. This index consists of an average of scores derived from 9 indicators:

- Closing of schools, workplaces, public transport.
- Cancelation of public events.
- Limits on gatherings size.
- Orders of confinement.
- Restriction on internal movements between regions and international travel.
- Presence of public information campaign.

4.4. Model selection and evaluation

We built a model of target labels Q and O using all the predictors of table 1. We chose ensemble methods for our machine learning model because they are effective for tabular data and can leverage various features. Boosting-based approaches, like XGBoost, can outperform deep learning models with less tuning required [43]. While we implemented XGBoost, we also tried Random Forest, a related technique that uses randomization. We found comparable results with much less computational expense, which is why we chose Random Forest for the study. For comparison, Mahajan *et al* [32] used LSTM and XGBoost algorithms. However, LSTM requires past data for predictions, which does not align with our goal of generalizing predictions for non-monitored road segments.

Random forest regression is a widely-used technique that deals well with correlated features and big datasets. Training was conducted on all road links available after the gap filling method, without road type distinctions, and with hourly records available for the flow and occupancy variables spanning the period from 2018 to 2022. The results are presented separately for peripheral roads and intra-city roads, as they exhibit distinct magnitudes and form well-separated clusters.

Performance evaluation of the random forest predictions was conducted using the normalized root mean squared error (RMSE), and symmetric mean absolute percentage error (SMAPE) for comparison with other studies [32]. The standard deviation-based normalized RMSE is closely linked to the R^2 score as it represents the ratio between the variation not explained by the model versus the overall variation in the data. SMAPE is a scale-independent error metric that assesses the relative accuracy of predictions, making it suitable for cases where the absolute magnitude of the target variable differs significantly among sensor types

$$\text{RMSE}_{\text{normalized}} = \frac{\text{RMSE}}{\sigma_x} = \sqrt{1 - R^2}$$

$$\text{SMAPE} = \frac{2}{n} \sum_{i=1}^n \frac{|x_i - \hat{x}_i|}{|x_i| + |\hat{x}_i|}$$

Equation (2): metrics to assess model performances.

Performance evaluation was also conducted using the decomposition of mean squared error (MSE) into three components [44]: standard bias (SB), standard deviation error (SDSD), and lack of correlation (LCS). These components sum up to the value of the MSE, and thus explain the origins of the quadratic error

$$\text{SB} = (\bar{x} - \bar{\hat{x}})^2$$

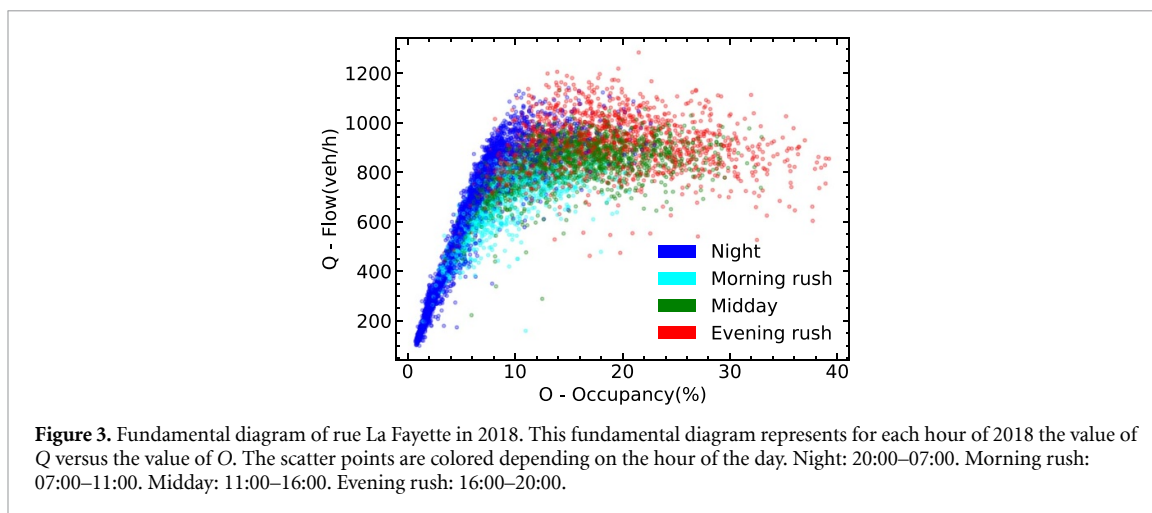
$$\text{SDSD} = (sd(x) - sd(\hat{x}))^2$$

$$\text{LCS} = 2sd(x)sd(\hat{x})(1 - \rho(x, \hat{x}))$$

$$\text{MSE} = \text{SB} + \text{SDSD} + \text{LCS}.$$

Equation (3): decomposition of the MSE.

To determine the optimal number of estimators and max depth, we performed a grid search analysis using k -fold cross-validation, specifically with k set to 5. The couple of parameters (number of trees, max depth) yielding the best performance on the cross-validated results was determined to be respectively (20, 40) for Q and (20, 30) for O . The number of estimators did not have a significant influence on the performances of the prediction, in opposition to the max depth parameter. Gradient boosting did not give significant improvements compared to random forests, and was more computationally intensive.



5. Results and discussion

5.1. Gap filling of road sensors time series

In this section, we first evaluate the performance of our temporal gap filling method for each road. To do so, we compute a 5-fold cross validation in terms of normalized RMSE and SMAPE defined in equation (2). Only roads with a normalized RMSE below $\sqrt{0.4}$ (corresponding to $R2 > 0.6$) together with a SMAPE below 0.4, and at least 15 000 records (nearly 2 years of data) are kept as target for the machine learning model. Figure 2 shows the test prediction performances using violin diagrams of the two metrics. Out of the 2086 sensors that were maintained through 2018–2022, 1397 showed satisfactory performances on Q predictions according to the RMSE and SMAPE thresholds. For the gap filling of the O variable, 871 road links were kept at the end of the same process. Figure 2 displays the spatial repartition of road segments that were not kept as not predictable enough to be gap filled up by our methodology. These 871 road links with good quality gap-filling were used for the training of our machine learning model.

5.2. Traffic flux and vehicles occupancy relationships

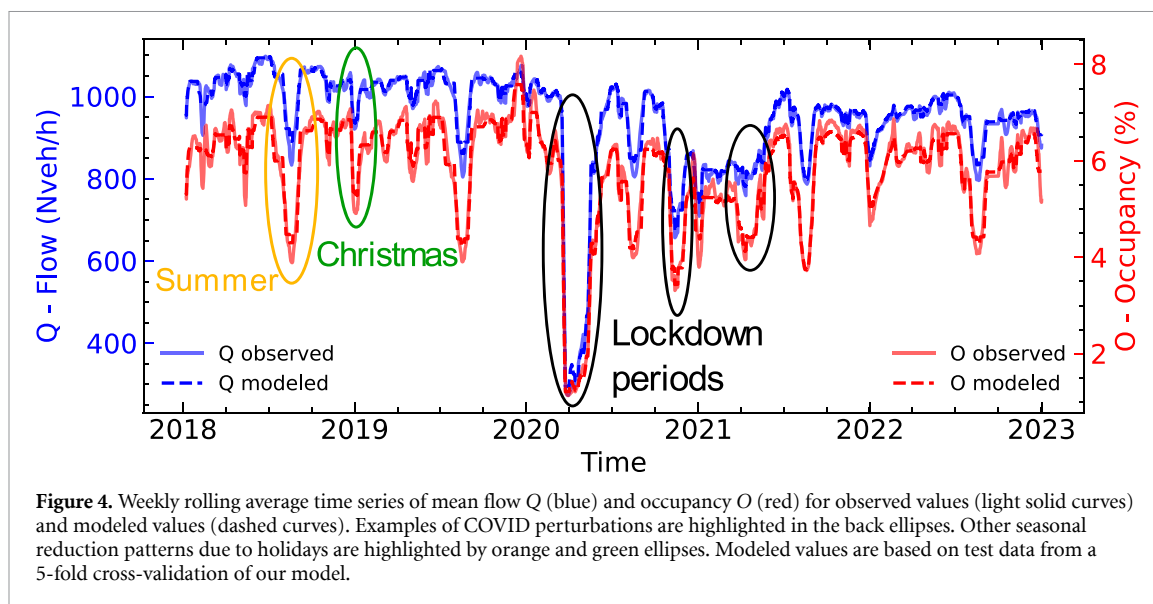
We should note that Q and O are strongly linked through traffic processes, which defines the so-called ‘fundamental diagrams’ [34, 45, 46]. The fundamental diagram represents an emerging relation between the flow Q and the congestion K , for which O is a good proxy. An example of fundamental diagram is shown in figure 3 from one sensor located in Rue Lafayette (inside the city). The left-hand part of the figure shows Q increasing with O quasi linearly. The corresponding slope is linked to the speed limitation as this left-hand part of the fundamental diagram represents the free-flow regime. The right-hand part shows that Q is stable or decreases when O further increases above a critical O threshold of about 12%, denoting a congested traffic regime where average speed decreases when O further increases. It is important to note that each road link has a specific fundamental diagram, and that this diagram evolves through the day with different congestion levels, as evidenced by different relationships between night, morning, midday and evening (figure 3).

5.3. Observed traffic changes from 2018 to 2022 at city scale and ability of the model to capture them

To present the results of the model that predicts Q and O in each road link, we first give a qualitative description of the main changes that were observed in these two variables, which can guide the evaluation of the model for its ability to capture those changes. Both Q and O are subject to periodical cycles with the influence of weekends, holidays and seasons, and experienced a strong drop during the COVID-19 period. This behavior is illustrated in figure 4 for the average of Q and O of the Paris roads. The most important perturbation clearly occurred during the first lockdown during the COVID-19 crisis. Smaller perturbations such as other lockdowns, summer or Christmas holidays can be observed as well in the data presented in figure 4.

Figure 4 also displays the same information simulated by our model on a 5-fold cross validation. The model accurately replicates weekly averages of Q and O over the whole study period, with a normalized RMSE of 0.105 ($R2 = 0.989$) for Q and 0.173 ($R2 = 0.970$) for O . Noteworthy changes due to the pandemic and the seasonal holiday effects are well represented by the model. Lower performances for the occupancy variable are explained by the failure to capture short term intra seasonal variations as depicted by figure 4.

Figure 5 displays changes in the average of observed Q and O relative to 2019 for subsequent years. Only values corresponding to working days (Monday to Friday) were selected for illustration. In 2020, a large



decline in traffic variables by about 10%–25% for Q , and 10%–30% for O is observed during busy hours, attributed to the first COVID-19 lockdown (17 March to 11 May) and other restrictions that followed. These COVID 19 induced reductions in traffic are particularly pronounced during night time hours due to the curfew measures. In 2021, there was a partial recovery, but lower values of Q and O persisted during night time hours due to continuing curfew restrictions (until 9 June). In 2022, a global recovery in traffic is evident with the disappearance of all COVID-19 mobility constraints, but there is an intriguing pattern of persisting lower Q and O values during morning (07:00–10:00) and evening (17:00–22:00) rush hours than in 2019. This reduction in morning traffic compared to pre-crisis conditions may be attributed to the widespread adoption of remote working [47].

The average temporal patterns depicted in figures 4 and 5 are faithfully replicated by our machine learning model, confirming its reliability for computing spatial averages over time. Notably the comparison of hourly median values between observations and the model shows a normalized RMSE of 0.033 ($R^2 = 0.999$) for Q and 0.077 ($R^2 = 0.994$) for O . Other quantiles of the distribution are also faithfully replicated by our model as depicted in figure 5. This underscores the model’s minimal bias in predictions as we aggregate roads.

5.4. Machine learning model performances for predicting traffic variables across main road links

To answer our first research question about whether our machine learning model (based on the predictors described in table 1 can simulate the sensors observed Q and O time series at road level, we analyze in table 2 and table 3 the model cross validation performances using the SMAPE, RMSE error metrics and the decomposition of the MSE into additive terms. Overall, the mean RMSE on test data is 33% for Q and 81% for O , respectively. The R^2 score is 0.89 for Q and 0.34 for O . In terms of SMAPE, our model predicts values within the rate of 60% even for the O variable which exhibits lower performances. The intricacy clusters demonstrate poorer results than the peripheral, occasionally reaching an RMSE above 100% in the case of O prediction for 38% of the roads. While normalized quadratic errors on O variables can become important, the SMAPE score suggests that these errors might be reasonable in terms of absolute values.

The previous study conducted by Mahajan *et al* [19] employed a XGBoost model to predict Q for the year 2019, utilizing comparable data. In our study, spanning a duration of five years and encompassing perturbations associated with the COVID-19 crisis, we observed in table 2 substantial improvement in SMAPE compared to their model (0.35 ± 0.03 versus 0.52 ± 0.05 on test data).

While the temporal error magnitudes, respectively around 37% and 54% of SMAPE for Q and O , remain relatively consistent throughout the years, increased local uncertainty becomes evident during the COVID-19 outbreak, possibly because our use of a stringency index for the entire France did not allow the capture of local reductions in traffic across road segments (figure 6). The normalized RMSE displays more pronounced fluctuations compared to the SMAPE (figure 6), with particularly significant deteriorations during the COVID-19 period (normalized RMSE exceeding 1 in some cases). This misfit can be attributed to the lower traffic flow and occupancy levels during the COVID crisis, when the model encounters limitations

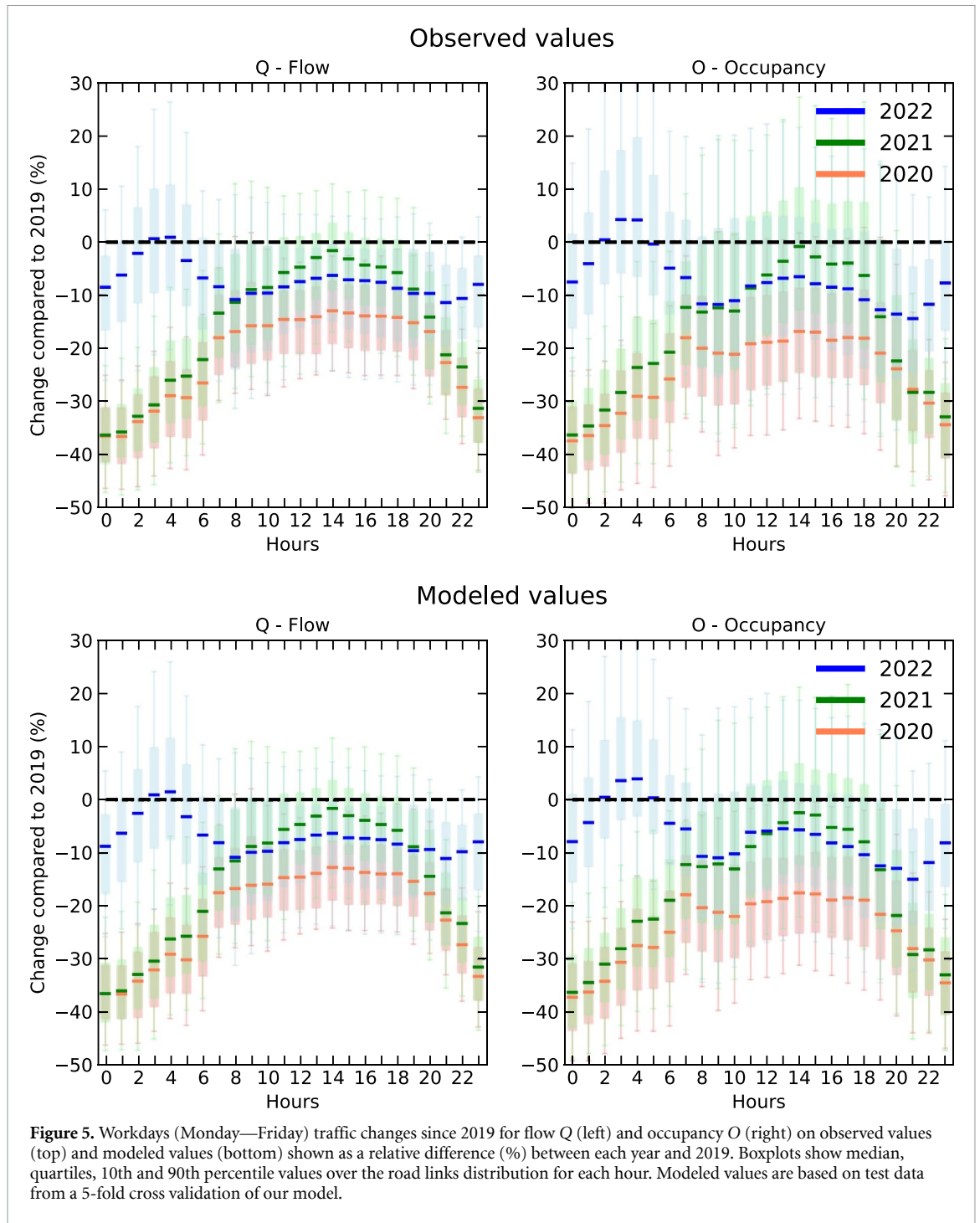


Table 2. Summary of performances result on Q target variable based on a 5-fold cross validation (mean ± max deviation from mean).

Type	Metric	Train data	Test data
All	Normalized RMSE	0.1 ± 0.0	0.33 ± 0.03
	SMAPE	0.11 ± 0.0	0.35 ± 0.03
Peripheral	Normalized RMSE	0.17 ± 0.01	0.48 ± 0.07
	SMAPE	0.07 ± 0.0	0.19 ± 0.02
Intracity	Normalized RMSE	0.16 ± 0.01	0.73 ± 0.09
	SMAPE	0.11 ± 0.0	0.38 ± 0.02

in accurately discerning traffic patterns due to reduced data volume. In this context, where individual random behaviors hold increased significance, metrics that emphasize variation explainability experience more pronounced penalization in contrast to those based on absolute percentage errors.

Table 3. Summary of performances result on *O* target variable based on a 5-fold cross-validation (mean \pm max deviation from mean).

Type	Metric	Train data	Test data
All	Normalized RMSE	0.32 \pm 0.01	0.81 \pm 0.12
	SMAPE	0.21 \pm 0.0	0.53 \pm 0.04
Peripheral	Normalized RMSE	0.31 \pm 0.01	0.71 \pm 0.07
	SMAPE	0.13 \pm 0.0	0.31 \pm 0.02
Intracity	Normalized RMSE	0.37 \pm 0.01	0.96 \pm 0.16
	SMAPE	0.22 \pm 0.0	0.56 \pm 0.04

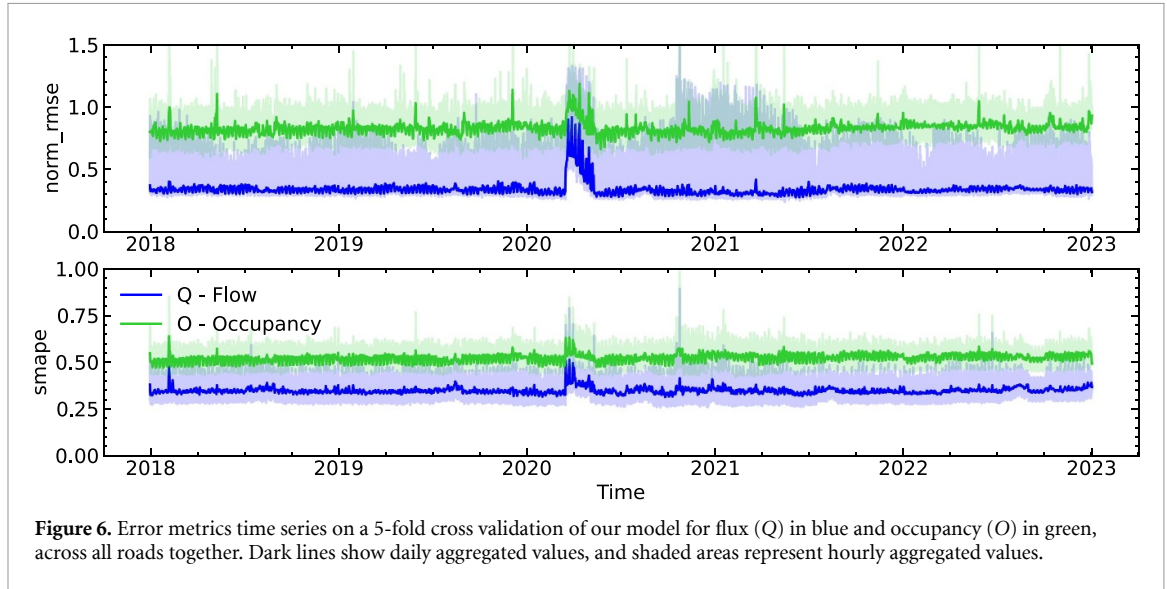


Figure 6. Error metrics time series on a 5-fold cross validation of our model for flux (*Q*) in blue and occupancy (*O*) in green, across all roads together. Dark lines show daily aggregated values, and shaded areas represent hourly aggregated values.

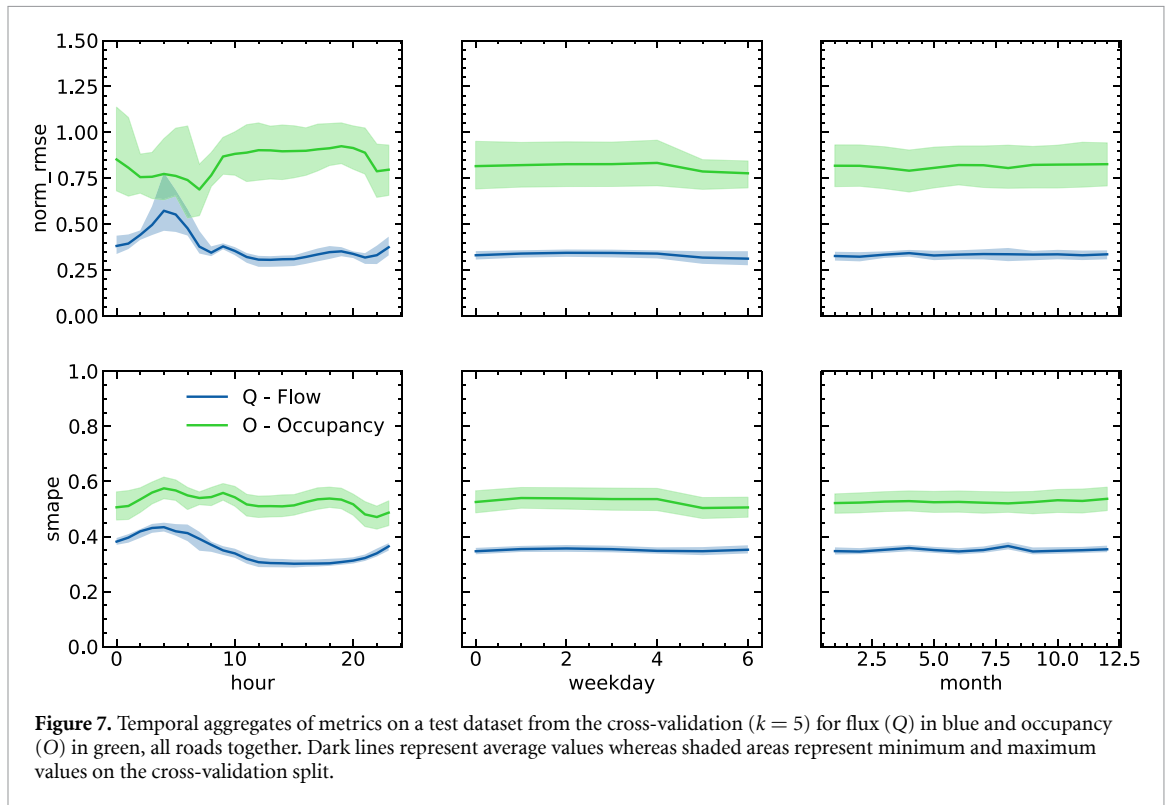


Figure 7. Temporal aggregates of metrics on a test dataset from the cross-validation ($k = 5$) for flux (*Q*) in blue and occupancy (*O*) in green, all roads together. Dark lines represent average values whereas shaded areas represent minimum and maximum values on the cross-validation split.

Looking at temporal patterns on figure 7, we can see that errors magnitude are mainly driven by original values magnitudes (slightly better scores during night hours, weekend, August). Busy hours occupancy values seem to be the hardest ones to predict accurately.

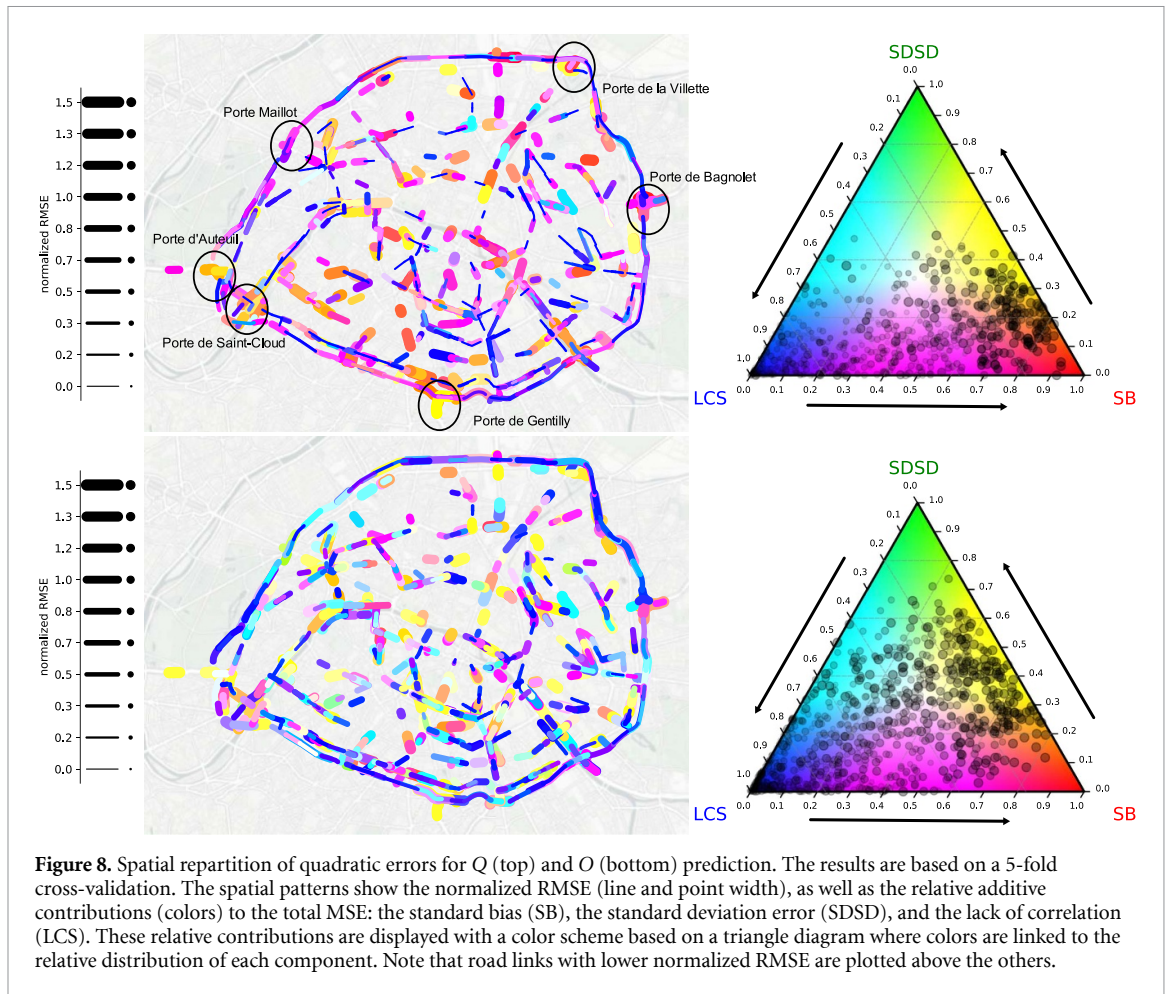


Figure 8 represents the spatial distribution of the normalized RMSE between model and the sensor data. The RMSE is proportional to the width of each road segment in the figure and the colors of each road refer to the relative contributions of each additive term explaining the model MSE error (equation (3)). The results show that the MSE for Q is dominated by SB (bias) and LCS (lack of correlation between predicted and observed time series). The highest errors are attributed to a large SB values (i.e. more than 50% of MSE), particularly concerning critical connections between the city and its outskirts (e.g. Porte de Gentilly, Porte de Bagnole, Porte de la Chapelle, near Porte d’Auteuil/Saint-Cloud, etc). These discrepancies can be attributed to the model’s limited ability to accurately recognize such road categories or behaviors. One can note that the significance of the contribution of SDSD (predicted variability magnitude difference from observed one) to the total MSE error is more pronounced in the context of O , where the highest errors are attributed to both a failure in predicting the mean value (SB) but also the variability (SDSD) of the occupancy on these road links.

Comparing the performances, it is evident that our model performs better on peripheral roads than on intra-city road links. The behavior of peripheral roads appears to be more predictable, contributing to more favorable prediction outcomes. By comparison, the variability of Q and O illustrated in figure 12 is of greater importance for intra-city roads. Additionally, the quality of the OSM data used to define road features may be called into question, as it may lead to incorrect map matching results and subsequent conclusions. This is because intra-city roads undergo more frequent design modifications (e.g. lane closures, cycle paths, bus paths), making the OSM data less reliable for such road links than for the peripheral. The presence of outdated information in OSM for intra-city roads might thus partly explain their less accurate predictions, as shown in figure 8.

Our model is predicting the variable O with a lesser accuracy than Q . This result can be attributed to the fact that information related to O is more susceptible to perturbations, such as a car parking on the sensor or misplacement of the sensor on the lane (e.g. questions regarding proximity to traffic signals and the choice of monitored lane). Unlike Q , which benefits from multiple sensors ($2 \times N_{lanes} - 1$) to detect vehicles straddling two lanes, O relies on only one sensor, making it more vulnerable to disruptions. The increasing deployment

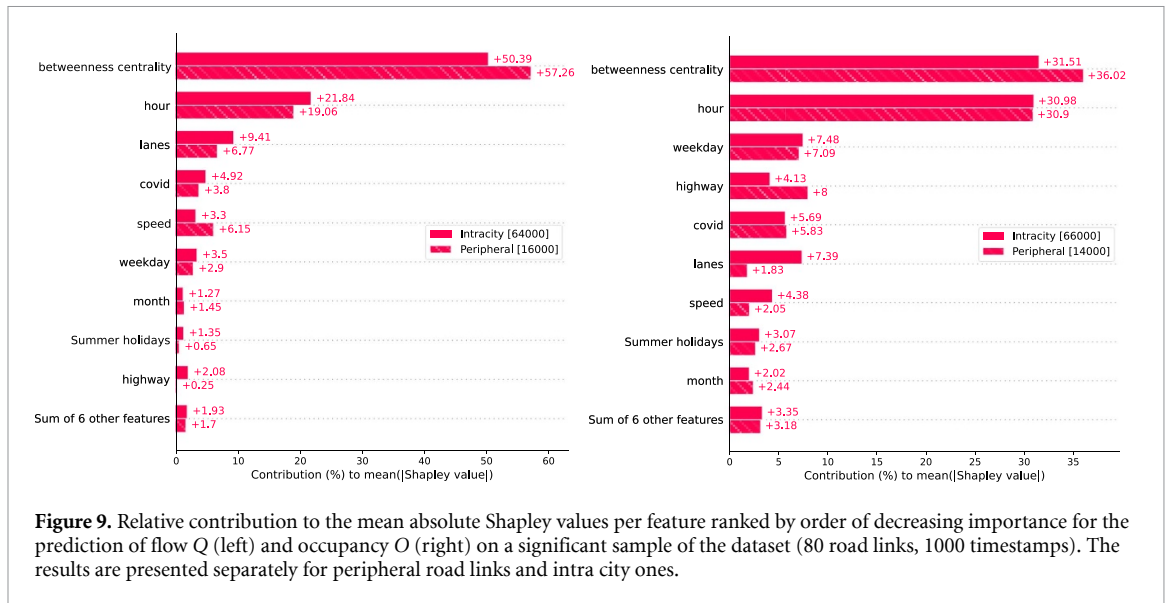


Figure 9. Relative contribution to the mean absolute Shapley values per feature ranked by order of decreasing importance for the prediction of flow Q (left) and occupancy O (right) on a significant sample of the dataset (80 road links, 1000 timestamps). The results are presented separately for peripheral road links and intra city ones.

of cameras as monitoring devices holds promise for improving data quality and, consequently, input data accuracy.

Additionally, the relative variance in the distribution of O is larger than that of Q for a similar hour (figure 12). This higher variability in O data poses additional complexities for the model, as it must account for larger fluctuations and uncertainties in the observations, making accurate predictions more challenging.

Our model lacks consideration for the spatial interdependencies among distinct road segments, wherein information at the onset of a road influences conditions at its terminus. Algorithms such as Graph Neural Networks could potentially mitigate this issue, yielding superior predictive capabilities for variables Q and O , as demonstrated in previous applications like the prediction of arrival times in Google Maps [48].

Exploring temporal factors, alternative approaches were considered, including the incorporation of strike data, construction activities leading to lane closures, and adverse weather conditions. However, integrating these data sources were challenging and failed to enhance model performance commensurate with the complexity of the task, resulting in their exclusion from the scope of this study.

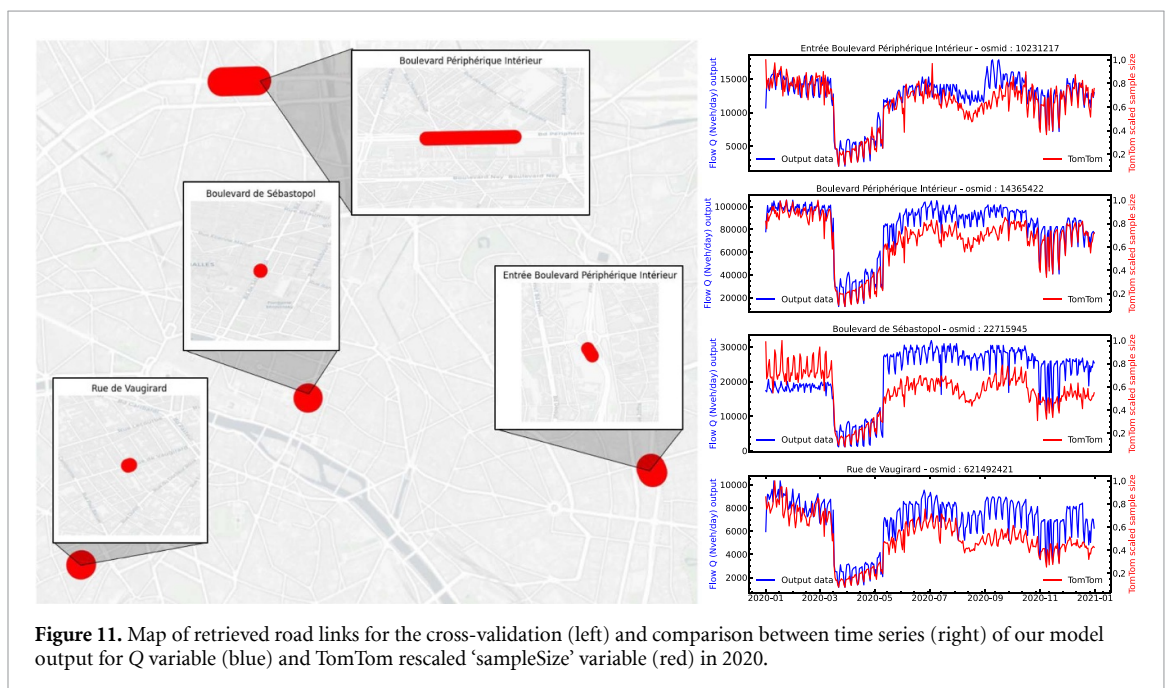
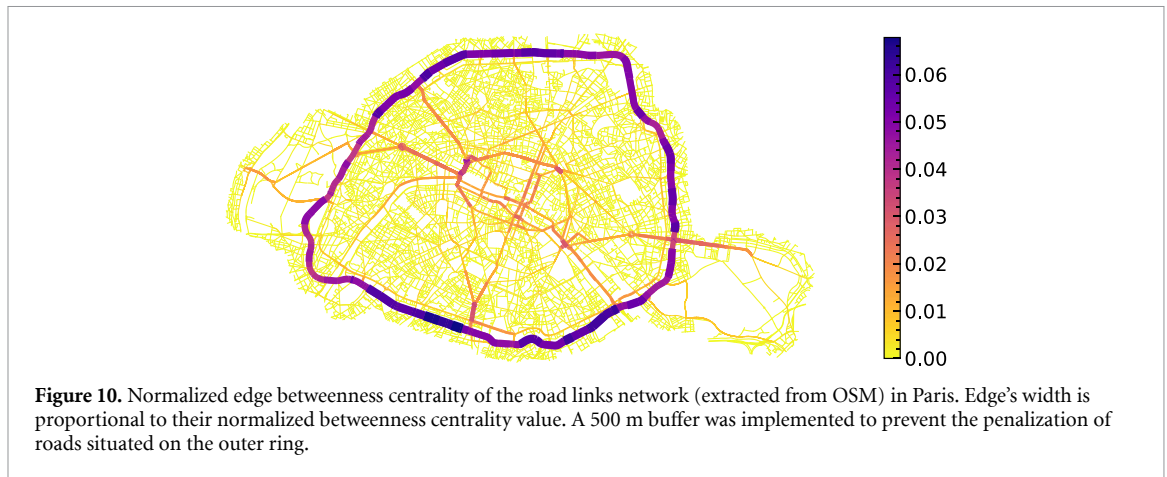
5.5. Explaining traffic variations

To answer our second research question about the main features that contribute significantly to the predictions of Q and O patterns generated by our model, we calculate Shapley values [49] and compare their mean absolute values across features, as depicted in figure 9. Sampling is made by randomly selecting specific road links and timestamps, because the model size together with the number of observations were too heavy to compute Shapley values on all the dataset. Thus, we computed the global feature importance multiple times with a random sampling to first compute the model (150 road links, 6000 timestamps) and then another random sampling (80 road links, 1000 timestamps) to compute the Shapley values. This approach showed an absence of major changes between the different random samplings, allowing us to propose a representative result for the whole dataset with a reasonable computation time (figure 9).

Figure 9 shows two distinct spatial clusters of roads, namely peripheral road links and intracity ones. Notably, the ‘lanes’ variable (stating for the number of lanes) obviously emerges as a potent explanatory factor for both Q and O predictions on intra city road links, being respectively the third and fourth feature in terms of relative importance for such segments. We could attribute that to the greater diversity of infrastructure and importance on this cluster, characterized by its road width and sizes.

We found that spatial features taken together (betweenness-centrality, lanes, speed, highway) carry greater importance for Q than temporal features (hour, COVID-19 stringency, weekday) while it is nearly the case for O . The hour feature yet emerges as of key importance, given that human behavior is predominantly influenced by diurnal cycles. The hour of the day and weekday number show a greater importance for predicting O (linked to traffic congestion) than Q . Obviously, the ‘lanes’ and ‘speed’ features show greater importance for Q than for O . In a nutshell, figure 9 indicates that while both Q and O values are mainly related to infrastructure-related features, temporal proxies play a greater role in explaining congestion than the magnitude of the traffic.

Of notable interest is the result that Betweenness centrality seems to surpass all other features in terms of importance (roughly 30%–55% of mean ($|\text{Shapley value}|$)) within our dataset, providing richer insights than



the classifications derived from OpenStreetMap. Figure 10 is depicting the values of normalized betweenness centrality for the whole OSM network of the city of Paris, where main axes for road traffic can easily be observed. Consequently, one could argue that this particular feature holds significant explanatory power, making it relevant for tasks such as urban planning modelization.

5.6. Cross-validation on non-monitored road links

To validate the results of our model using independent data, we conducted an additional analysis acquiring daily data from the TomTom traffic stats API [50] for four distinct road links that were not part of our initial dataset. These selected cross-validation road segments encompassed a peripheral route, a ramp leading to the peripheral road, a bustling boulevard, and a typical city street. In figure 11, we compared the rescaled 'sampleSize' variable (unique vehicles traveling on segments) to the model's predictions for the year 2020 on the flow variable Q . The results of this comparison indicate a successful replication of temporal correlations, capturing essential aspects such as lockdown periods and seasonal traffic patterns.

However, the proportionality between the two signals exhibit inconsistencies in some cases. For instance, Boulevard Sebastopol exhibits higher Tomtom signal before the pandemic than after compared to the model (figure 11). This could potentially be attributed to variations in data providers that TomTom relies upon for the computation of their statistics. Additionally, our results for this evaluation reveal that the coefficient linking the two signals differed between the various street types. In this particular example, the share of road users relying on TomTom data ranged roughly between 20% and 50%. These heterogeneous shares are consistent with other comparisons with input data (SI section 6) that indicate potential challenges pertaining to spatial homogeneity in the context of this type of floating car data.

6. Conclusion

We presented a comprehensive transportation model for each non-residential road segment in the city of Paris, based on the upscaling of sensor's point scale measurements, and operating at an hourly time scale throughout the period from 2018 to 2022. This transportation model captures the flow and the occupancy variables, indicating respectively the magnitude of traffic and the level of congestion on road segments. We analyzed the performances and the limitations of the approach used here, and discussed the relative relevancy of the employed features. In future work, the next logical step is to develop a comprehensive CO₂ and pollutant inventory for the entire city with a similar high resolution. Thus, our model will be used to gain insight on traffic variables for each road link in the city. A new emission model should incorporate an emission factor that considers the associated mean speed of the vehicle fleet, derived from both flow and occupancy data. By considering both flow and occupancy information, we expect to capture speed-dependent emission factors, thus providing a more robust estimation of CO₂ and pollutant emissions. We already observed that calculating mean speed from the flow-to-occupancy ratio can introduce significant errors, primarily arising from uncertainties linked to these two variables. Therefore, implementing a capping mechanism may be necessary to mitigate the risk of substantial errors in speed, which could consequently result in elevated errors in emission factor estimations.

Moreover, the methodology employed in this study is not limited to Paris alone. With access to comparable data from other urban centers (for instance Madrid [51], Los Angeles [52], Berlin [53] or Auckland [54]) this methodology can be replicated and applied to assess transportation dynamics as well as pollutant inventories. We recommend using a model trained with these local datasets in order to prevent biases, as related road segments in different cities may exhibit distinct traffic patterns. This cross-applicability underscores its potential as a versatile tool for sustainable monitoring on an urban scale. Nonetheless, it is important to acknowledge that this approach may require methodological adjustments or calibration efforts, such as rescaling with census data, to ensure the accurate modeling of these traffic variables in different urban contexts.

By integrating various predictor features with hourly road link data, we established a robust framework that can be harnessed to evaluate the potential effects of proposed changes to the road infrastructure, such as the introduction of new road types or access gates. The model's accuracy in capturing transportation patterns can guide decision-makers in crafting policies that align with the goals of reducing environmental impact, enhancing public health, and improving the overall quality of life in urban areas. Furthermore, the approach offers a dynamic framework that can adapt to changing circumstances, making it well-suited for evaluating scenarios in response to unforeseen events like pandemics, climate change, or shifts in transportation behavior. The versatility of this model positions it as a valuable asset for shaping the future of Paris and other cities committed to sustainable urban development.

Data availability statement

The data cannot be made publicly available upon publication because the cost of preparing, depositing and hosting the data would be prohibitive within the terms of this research project. The data that support the findings of this study are available upon reasonable request from the authors.

Disclaimer

The authors declare that they have no known competing financial interests or personal relationships that could have appeared to influence the work reported in this paper.

During the preparation of this work the authors used the service ChatGPT 3.5 in order to improve readability and language. After using this service, the authors reviewed and edited the content as needed and take full responsibility for the content of the publication.

ORCID iDs

Xavier Bonnemaizon  <https://orcid.org/0009-0007-5262-5175>

Philippe Ciais  <https://orcid.org/0000-0001-8560-4943>

References

- [1] Ministère de la Transition Ecologique 2022 Chiffres clés des transports *Chiffres clés transports 2022* (available at: www.statistiques.developpement-durable.gouv.fr/edition-numerique/chiffres-cles-transports-2022/19-emissions-de-gaz-a-effet-de-serre-du-transport.php) (Accessed 17 April 2023)

- [2] European Parliament EU ban on sale of new petrol and diesel cars from 2035 explained (News European Parliament) (available at: <https://www.europarl.europa.eu/topics/en/article/20221019STO44572/eu-ban-on-sale-of-new-petrol-and-diesel-cars-from-2035-explained>) (Accessed 17 April 2023)
- [3] Pinto J A, Kumar P, Alonso M F, Andreão W L, Pedruzzi R, Dos Santos F S, Moreira D M and Albuquerque T T D A 2020 Traffic data in air quality modeling: a review of key variables, improvements in results, open problems and challenges in current research *Atmos. Pollut. Res.* **11** 454–68
- [4] Pant P and Harrison R M 2013 Estimation of the contribution of road traffic emissions to particulate matter concentrations from field measurements: a review *Atmos. Environ.* **77** 78–97
- [5] Ville de Paris Le Bilan Carbone® de Paris 2018 (available at: <https://cdn.paris.fr/paris/2020/02/06/dc2edb10d13ae050815850f721f5a837.pdf>) (Accessed 16 January 2024)
- [6] IPCC AR6 Synthesis Report 2022 (available at: https://report.ipcc.ch/ar6syr/pdf/IPCC_AR6_SYR_LongerReport.pdf) (Accessed 22 April 2023)
- [7] EMISIA COPERT documentation (available at: www.emisia.com/wp-content/uploads/2023/09/1.A.3.b.i-iv-Road-transport-2023_Sep.pdf) (Accessed 15 November 2022)
- [8] Direction de la Voirie et des Déplacements Comptage routier—Historique—Données trafic issues des capteurs permanents (available at: <https://parisdata.opendatasoft.com/explore/dataset/comptages-routiers-permanents-historique/information/>) (Accessed 17 April 2023)
- [9] Ayaz S, Khattak K S, Khan Z H, Minallah N, Khan M A and Khan A N 2022 Sensing technologies for traffic flow characterization: from heterogeneous traffic perspective *J. Appl. Eng. Sci.* **20** 29–40
- [10] Lenormand M, Picornell M, Cantú-Ros O G, Tugores A, Louail T, Herranz R, Barthelemy M, Frías-Martínez E and Ramasco J J 2014 Cross-checking different sources of mobility information *PLoS One* **9** e105184
- [11] Xu Y, Clemente R D and González M C 2021 Understanding vehicular routing behavior with location-based service data *EPJ Data Sci.* **10** 12
- [12] Hörl S and Balac M 2021 Synthetic population and travel demand for Paris and Île-de-France based on open and publicly available data *Transp. Res. C* **130** 103291
- [13] Sadowski A, Galar Z, Walasek R, Zimon G and Engelseth P 2021 Big data insight on global mobility during the Covid-19 pandemic lockdown *J. Big Data* **8** 78
- [14] Apronti D, Ksaibati K, Gerow K and Hepner J J 2016 Estimating traffic volume on Wyoming low volume roads using linear and logistic regression methods *J. Traffic Transp. Eng.* **3** 493–506
- [15] Tarunesh I and Chung E 2020 Predicting traffic volume and occupancy at failed detectors *Transp. Res. Proc.* **48** 1072–83
- [16] Xing J, Wu W, Cheng Q and Liu R 2022 Traffic state estimation of urban road networks by multi-source data fusion: review and new insights *Physica A* **595** 127079
- [17] Oda T, Haga C, Hosomi K, Matsui T and Bun R 2021 Errors and uncertainties associated with the use of unconventional activity data for estimating CO₂ emissions: the case for traffic emissions in Japan *Environ. Res. Lett.* **16** 084058
- [18] Guevara M et al 2022 European primary emissions of criteria pollutants and greenhouse gases in 2020 modulated by the COVID-19 pandemic disruptions *Earth Syst. Sci. Data* **14** 2521–52
- [19] Google COVID-19 community mobility report *COVID-19 Community Mobility Report* (available at: www.google.com/covid19/mobility?hl=fr) (Accessed 19 April 2023)
- [20] Huo D et al 2022 Carbon Monitor Cities near-real-time daily estimates of CO₂ emissions from 1500 cities worldwide *Sci. Data* **9** 1
- [21] TomTom Traffic Index—Live traffic statistics and historical data *TomTom Traffic Index—Live traffic statistics and historical data* (available at: www.tomtom.com/traffic-index/) (Accessed 19 April 2023)
- [22] Biswal A, Singh V, Malik L, Tiwari G, Ravindra K and Mor S 2023 Spatially resolved hourly traffic emission over megacity Delhi using advanced traffic flow data *Earth Syst. Sci. Data* **15** 661–80
- [23] EMISIA SA COPERT | the industry standard emissions calculator (available at: www.emisia.com/utilities/copert/) (Accessed 4 December 2023)
- [24] Li Y, Lv C, Yang N, Liu H and Liu Z 2020 A study of high temporal-spatial resolution greenhouse gas emissions inventory for on-road vehicles based on traffic speed-flow model: a case of Beijing *J. Clean. Prod.* **277** 122419
- [25] US EPA MOVES and Mobile Source Emissions Research (available at: www.epa.gov/moves) (Accessed 6 December 2023)
- [26] Observatoire Parisien des Mobilités *Barometre Trimestriel—Paris Data* (available at: <https://opendata.paris.fr/pages/barometre/>) (Accessed 17 April 2023)
- [27] Ville de Paris La Zone à faibles émissions (ZFE) (available at: www.paris.fr/pages/la-zone-a-faibles-emissions-zfe-pour-lutter-contre-la-pollution-de-l-air-16799) (Accessed 22 April 2023)
- [28] Nicolini G et al 2022 Direct observations of CO₂ emission reductions due to COVID-19 lockdown across European urban districts *Sci. Total Environ.* **830** 154662
- [29] OpenStreetMap *OpenStreetMap* (available at: www.openstreetmap.org/) (Accessed 19 March 2023)
- [30] Parif A 2019 Emissions de polluants atmosphériques et de gaz à effet de serre (available at: www.airparif.asso.fr/sites/default/files/pdf/Bilan2019.pdf)
- [31] Uber Movement let's find smarter ways forward, together (available at: www.uber.com/fr/fr/business/movement-decommissioning/) (Accessed 19 April 2023)
- [32] Mahajan V, Cantelmo G, Rothfeld R and Antoniou C 2022 *IET Intell. Transp. Syst.* **17** 804–24
- [33] Uber Movement is no longer active *Uber* (available at: www.uber.com/fr/fr/business/movement-decommissioning/) (Accessed 3 July 2024)
- [34] Buisson C and Lesort J-B 2010 *Comprendre le Trafic Routier* (Certu)
- [35] Gumbel E J 1935 Les valeurs extrêmes des distributions statistiques *Ann. Inst. Henri Poincaré* **5** 115–58 (available at: http://www.numdam.org/item/AIHP_1935__5_2_115_0/)
- [36] Hale T et al 2021 A global panel database of pandemic policies (Oxford COVID-19 Government Response Tracker) *Nat. Hum. Behav.* **5** 529–38
- [37] Le calendrier scolaire—data.gouv.fr (available at: www.data.gouv.fr/fr/datasets/le-calendrier-scolaire/) (Accessed 19 April 2023)
- [38] Barthelemy M 2004 Betweenness centrality in large complex networks *Eur. Phys. J. B* **38** 163–8
- [39] Tortora M, Cordelli E and Soda P 2022 PyTrack: a map-matching-based python toolbox for vehicle trajectory reconstruction *IEEE Access* **10** 112713–20
- [40] SharedStreets · GitHub (available at: <https://github.com/sharedstreets>) (Accessed 3 January 2023)

- [41] Newson P and Krumm J 2009 Hidden Markov map matching through noise and sparseness *Proc. 17th ACM SIGSPATIAL Int. Conf. on Advances in Geographic Information Systems—GIS'09* (ACM Press) p 336
- [42] Barthelemy M 2015 From paths to blocks: new measures for street patterns *Environ. Plan. B* **44**
- [43] Shwartz-Ziv R and Armon A 2021 *Tabular Data: Deep Learning is Not All You Need* (<https://doi.org/10.48550/arXiv.2106.03253>)
- [44] Kobayashi K and Salam M U 2000 Comparing simulated and measured values using mean squared deviation and its components *Agron. J.* **92** 345–52
- [45] Daganzo C F 1994 The cell transmission model: a dynamic representation of highway traffic consistent with the hydrodynamic theory *Transp. Res. B* **28** 269–87
- [46] Daganzo C F and Geroliminis N 2008 An analytical approximation for the macroscopic fundamental diagram of urban traffic *Transp. Res. B* **42** 771–81
- [47] Nande F, Weber M-L, Bouchet S and Loup P 2022 Learning from the crisis: a study of the conditions promoting remote workers' well-being @GRH **44** 13–41
- [48] Derrow-Pinion A et al 2021 ETA prediction with graph neural networks in Google Maps *Proc. 30th ACM Int. Conf. on Information & Knowledge Management* pp 3767–76
- [49] Shapley L S 1953 *Contributions to the Theory of Games (AM-28), Volume II* ed H W Kuhn and A W Tucker (Princeton University Press) pp 307–18
- [50] TomTom *Traffic Stats API* (available at: <https://developer.tomtom.com/traffic-stats/documentation/product-information/introduction>) (Accessed 7 November 2023)
- [51] Datos Abiertos Madrid *Tráfico. Histórico de datos del tráfico desde 2013—Portal de datos abiertos del Ayuntamiento de Madrid* (available at: <https://datos.madrid.es/portal/site/egob/menuitem.c05c1f754a33a9f8e4b2e4b284f1a5a0/?vgnextoid=33cb30c367e78410VgnVCM1000000b205a0aRCRD&vgnnextchannel=374512b9ace9f310VgnVCM100000171f5a0aRCRD&vgnnextfmt=default>) (Accessed 30 October 2023)
- [52] LADOT *Traffic Counts Summary | Los Angeles—Open Data Portal* (available at: <https://data.lacity.org/Transportation/LADOT-Traffic-Counts-Summary/94wu-3ps3>) (Accessed 30 October 2023)
- [53] Open data Berlin *Verkehrsdetektion Berlin | offene Daten Berlin* (available at: <https://daten.berlin.de/datensaetze/verkehrsdetektion-berlin>) (Accessed 30 October 2023)
- [54] Transport Auckland *Traffic counts Auckland Transport* (available at: <https://at.govt.nzhttps://at.govt.nz/about-us/reports-publications/traffic-counts>) (Accessed 13 November 2023)