1 **PREDICTING OFFERS FOR MOBILITY-ON-DEMAND SERVICES:**
2 **A MACHINE LEARNING-BASED APPROACH**
3
4
5
6 **Chenhao Ding**[*]
7 Chair of Traffic Engineering and Control
8 Technical University of Munich
9 Archisstraße 21
10 80333 Munich, Germany
11 Email: chenhao.ding@tum.de
12 ORCID: 0009-0004-0423-9614
13
14 **Florian Dandl**
15 Chair of Traffic Engineering and Control
16 Technical University of Munich
17 Archisstraße 21
18 80333 Munich, Germany
19 Email: florian.dandl@tum.de
20 ORCID: 0000-0003-3706-9725
21
22 **Klaus Bogenberger**
23 Chair of Traffic Engineering and Control
24 Technical University of Munich
25 Archisstraße 21
26 80333 Munich, Germany
27 Email: klaus.bogenberger@tum.de
28 ORCID: 0000-0003-3868-9571
29
30 [*] Corresponding Author
31
32
33 Word Count: 6997 words + 2 table(s) $\times$ 250 = 7497 words
34
35
36
37
38
39
40 Submission Date: November 12, 2024

1 **ABSTRACT**
2 The rapid development of mobile internet technology has driven the widespread adoption of Mobility-
3 on-Demand services like ride-hailing and ride-pooling, which are integral to daily travel. Ensuring
4 the quality of these services is vital for Transportation Network Companies such as Uber, Lyft, and
5 Didi, as it directly impacts user experience and revenue. One quality aspect is providing timely
6 and precise predictions of waiting and driving times for users, as this information affects their
7 decision-making. However, creating these predictions is conceptually and computationally chal-
8 lenging in high-demand scenarios due to the dynamic nature of the services and exhaustive vehicle
9 information processing. This paper introduces a machine learning-based model incorporating a
10 three-dimensional encoding method based on pick-up likelihood, vehicle capacity, and detour de-
11 gree to efficiently encode user requests and fleet status. A selection mechanism filters encoded
12 information in temporal and spatial domains, reducing computational burden. The model employs
13 an encoder-decoder architecture, with a Convolutional Neural Network encoding fleet status as an
14 image and Fully Connected Networks decoding it for various predictions. A case study using the
15 Manhattan taxi dataset demonstrates the model's effectiveness, showing that spatial domain infor-
16 mation significantly impacts predictive performance. The model achieves a True Negative Rate of
17 approximately 90% in identifying unserviceable requests and outperforms heuristic algorithms in
18 travel time prediction accuracy and efficiency, though it lags in waiting time prediction. Future
19 research should explore nonlinear encoding for vehicle pick-up likelihood and test various neural
20 network configurations to improve prediction performance and transferability.
21
22 *Keywords*: Mobility-on-Demand, Machine learning, Offer prediction

**INTRODUCTION**

With the rapid development of mobile internet technology, Mobility-on-Demand (MoD) services, such as ride-hailing (RH) and ride-pooling (RP), have gained widespread attention and become integral to people's daily travel routines. According to Uber Technologies, Inc. (*1*), approximately 9.4 billion trips were completed globally in 2023. In comparison, DiDi Global Inc. (*2*) reported serving 10.8 billion trips in China alone. RH offers users on-demand door-to-door transportation services, significantly enhancing travel flexibility and alleviating parking pressure from users in urban areas. Building on this, RP allows multiple passengers with identical or similar routes to share a vehicle, partially combining their journeys. This method increases vehicle occupancy, thereby improving transportation efficiency, reducing the number of vehicles on the road, alleviating traffic congestion, and cutting carbon emissions.

For Transportation Network Companies (TNCs) such as Uber, Lyft, and Didi, providing users with offers that are produced quickly and contain accurate forecasts for trip planning is crucial for ensuring a positive user experience and increasing revenue. When users access the MoD service app to input trip details such as origin, destination, and expected departure time, they expect the system to promptly provide essential information about available services, enabling effective travel planning. Based on real-time fleet status, such as the location and onboard passenger count of each vehicle, most TNCs provide users with estimated service details, including vehicle availability, estimated pick-up time, estimated travel time, and estimated cost, before users commit to the service. This information has to be provided quickly, as lengthy computation times result in slow app response times and can frustrate users, thereby decreasing their willingness to use the services. Additionally, the predictions must be as accurate as possible to reflect the actual services provided later. Inaccurate predictions can lead to users making poor decisions, which, over time, can erode their trust in the MoD services, reducing overall acceptance and loyalty (*3*).

This study contributes to the research on MoD offer prediction by:

1. Describing the MoD offer creation problem.
2. Proposing a machine learning (ML)-based offer prediction approach that employs an encoder-decoder architecture to enhance flexibility for different prediction tasks.
3. Developing a three-dimensional encoding method based on pick-up likelihood, vehicle capacity, and detour degree to achieve efficient encoding of vehicle status.
4. Introducing a selection mechanism that improves prediction efficiency by reducing the amount of information processed in both temporal and spatial domains.
5. Conducting a case study using the publicly available New York City taxi dataset (*4*) to validate the performance of the proposed approach in terms of prediction accuracy and efficiency.

The remainder of this paper is organized as follows. The **LITERATURE REVIEW** section summarizes and discusses the research progress of existing MoD simulation frameworks and the state-of-the-art ML applications in the MoD field. The **METHODOLOGY** section discusses the MoD offer creation problem and introduces the proposed approach, providing a detailed explanation of the design concepts and functionalities of the individual modules. The **CASE STUDY** section describes the experiments based on the Manhattan taxi dataset, presents the results, evaluates the prediction accuracy and efficiency of the proposed approach, and discusses the reasons behind the model's performance. Finally, the **CONCLUSION AND FUTURE WORK** section summarizes the research and outlines potential directions for future work.

1 **LITERATURE REVIEW**
2 **Modeling of MoD services**
3 Research on Mobility-on-Demand services has become a key focus in transportation studies. A
4 scientometric analysis by Guo and Li (*5*) reviews the development and research directions of MoD
5 systems. Similarly, Shaheen and Cohen (*6*) examine methodologies for studying MoD, evaluating
6 the impacts of shared modes and autonomous MoD (AMoD) on public transportation. For instance,
7 Paparella et al. (*7*) propose a time-invariant network flow model for two-person ride-sharing MoD
8 systems, utilizing a Poisson process to describe trip requests and analyze ride-sharing in the spa-
9 tiotemporal dimension. Pavone et al. (*8*) employ a fluid model and linear programming to devise
10 an optimal real-time rebalancing policy, validated through simulations and hardware experiments.
11      In addition to mathematical models, comprehensive simulation frameworks have been de-
12 veloped to simulate MoD fleet operations. Commercial solutions such as Aimsun (*9*) and PTV
13 (*10*) have been launched, while researchers have increasingly focused on open-source frameworks.
14 Ruch et al. (*11*) introduce AMoDeus, a software package within MATSim (*12*) for the quantitative
15 analysis of AMoD systems. Auld et al. (*13*) develop POLARIS, integrating dynamic simulation
16 of travel demand, network supply, and network operations. Kucharski and Cats (*14*) introduce
17 MaaSSim, an agent-based simulator for two-sided mobility platforms, modeling urban mobility
18 dynamics with agents such as travelers, drivers, and an intermediary platform matching demand
19 with supply. FleetPy, developed by Engelhardt et al. (*15*), focuses on user-operator interactions
20 and supports multiple operators within the transportation system, offering a modular structure for
21 customization and transferability.
22      Vehicle assignment is a core challenge for MoD services and is closely tied to offer cre-
23 ation. This problem can be addressed using rules, heuristics, or optimization algorithms. Fagnant
24 and Kockelman (*16*) tackle this problem by iteratively matching travelers with the nearest vehicles
25 within a 5-minute travel time, expanding the search zone if necessary. Zhang and Guhathakurta
26 (*17*) assign non-busy vehicles to travelers by selecting the vehicle that offers the lowest combined
27 cost of time and fare. Alonso-Mora et al. (*18*) present a scalable mathematical model that dy-
28 namically optimizes routes, starting with a greedy assignment and improving through constrained
29 optimization for quick, high-quality solutions. Kucharski and Cats (*19*) use a demand-driven al-
30 gorithm that matches trips into shared rides via a directed shareability multi-graph, narrowing the
31 search space for efficient graph searches. Engelhardt et al. (*20*) introduce a speed-up algorithm to
32 find the system's optimal assignments at each decision time step. Hyland and Mahmassani (*21*) de-
33 fine and model assignment strategies for RH, finding that optimization-based strategies and those
34 incorporating en-route pickups and drop-offs reduce fleet miles and wait times, particularly during
35 peak periods.
36      Another crucial aspect of offer creation is the interaction between users and operators (*22*).
37 To simplify this process, some models, like POLARIS (*13*), treat user on-demand requests as
38 service bookings, assuming that users will invariably accept the operator's offer. This approach re-
39 duces the complexity of the offer creation problem, transforming it into a straightforward vehicle
40 routing issue and allowing researchers to focus more specifically on vehicle assignment. Con-
41 versely, models like SimMobility (*23*) also simulate the operator's ability to reject requests. In
42 these models, operators decide whether to fulfill a user's request based on the current fleet sta-
43 tus. However, in reality, the decision-making process between users and operators is bidirectional.
44 Operators assess the feasibility of fulfilling user requests in real-time while users evaluate their
45 options among various travel modes (*24*). The offers provided by operators play a crucial role in

1  users' decisions to use MoD services. Models such as MaaSSim (*14*) and FleetPy (*15*) effectively
2  simulate this bidirectional decision-making process.
3        Nonetheless, most MoD simulation frameworks that offer mode choice functionality for
4  users rely on comprehensive and precise simulation results as the basis for decision-making rather
5  than quickly predicted offers. To generate these results, the model must evaluate all available vehi-
6  cle data and determine the optimal vehicle assignments. The speed of this process is influenced by
7  both fleet size and demand level, making rapid processing challenging in scenarios involving large
8  fleets and high demand. Although Engelhardt et al. (*15*) employed heuristic algorithms to predict
9  offers by integrating each passenger's pick-up and drop-off locations into fleet routing plans and
10 optimizing them based on the operator's objectives, this approach still does not entirely mitigate
11 the impact of fleet size and demand level. For such highly dynamic, data-intensive problems, ML
12 offers a potentially superior solution (*25*).

13 **ML Approaches for MoD services**
14 ML offers innovative solutions for MoD services from a big data perspective. Chen et al. (*26*)
15 propose UberNet, a deep learning convolutional neural network designed for short-term demand
16 prediction for RH. Chu et al. (*27*) introduce Multi-Scale Convolutional Long Short-Term Memory
17 (MultiConvLSTM), a deep learning model inspired by image and video processing techniques that
18 treats travel demand as image pixel values to capture both temporal and spatial correlations. Ke
19 et al. (*28*) propose the Fusion Convolutional Long Short-Term Memory Network (FCL-Net), in-
20 tegrating convolutional LSTM layers, standard LSTM layers, and convolutional layers to capture
21 spatiotemporal correlations while ranking variable importance using spatially aggregated random
22 forest. Wang et al. (*29*) propose a Deep Spatio-Temporal ConvLSTM framework for travel de-
23 mand forecasting, capturing both temporal and spatial dependencies with separate branches for
24 closeness, period, and trend components, and employing a linear fusion method for final predic-
25 tions. On the operational side, Lei et al. (*30*) combine optimization and ML to proactively relocate
26 vehicles and ensure fair income distribution. Guo and Xu (*31*) propose a ride-sharing vehicle dis-
27 patching and routing method, which includes vehicle routing decision-making through a Markov
28 decision process with deep reinforcement learning, and request-vehicle assignment based on rout-
29 ing learning values. However, offer prediction using ML remains an underexplored area.

30 **METHODOLOGY**
31 This paper uses simulation to analyze the performance of two different methods to address the offer
32 creation problem. This section formalizes the MoD offer creation problem, presents the simulation
33 framework, briefly summarizes the insertion heuristic approach, and discusses in detail the design
34 concepts and functionalities of the ML-based solution approach.

35 **Formulation of the MoD Offer Creation Problem**
36 The offer creation problem describes the task of an MoD operator to generate information about
37 a possible trip with the service after a user sends a request with their app. In this task, the MoD
38 operator controls a set of vehicles $\mathbf{V}$ in a street network $\mathbf{G} = (\mathbf{N}, \mathbf{E})$, where $\mathbf{N}$ are the nodes and
39 $\mathbf{E}$ are the edges containing travel distance and travel time costs. The operator has access to the
40 current fleet state $\mathbf{F_S}$, which includes the position and onboard passengers of each vehicle, along
41 with their assigned routing plans and any unprocessed requests. An app request $R \subset \mathbf{R}$ includes
42 information about earliest departure times $t_{ed}$, origin $p_o \in \mathbf{N}$, destination $p_d \in \mathbf{N}$, and the number

1  of passengers $n_{pax}$, i.e., $R = (t_{ed}, p_o, p_d, n_{pax})$. When a new app request $R$ is received, the MoD
2  Offer Creation Problem must be solved, requiring the operator to generate an offer $O_p$ containing
3  information on the expected pick-up time $t_{epu}$, expected drop-off time $t_{edo}$, and expected cost $c_e$,
4  i.e., $O_p = (t_{epu}, t_{edo}, c_e)$. This problem can be very challenging.

5     • For large-scale services, finding a reasonable solution for a new request based on the
6      current fleet state can be computationally expensive.
7     • MoD services receive new app requests dynamically, and future vehicle-user assignments
8      can affect the service levels of previously assigned requests. Figure 1 illustrates how the
9      dynamic mechanism works. At the simulation time step $t_0$, the MoD operator receives a
10      request $R_0$ and decides to assign it to the vehicle $v_0$, generating a routing plan $V_{rp\_0}^0$ and
11      a predicted offer $O_p^0 = (t_{epu\_0}^0, t_{edo\_0}^0)$ [1]. The user accepts $O_p^0$ and decides to use the MoD
12      service. At the simulation time step $t_1$, the MoD operator receives a new request $R_1$ and
13      decides to combine $R_1$ and $R_0$ to form an RP. However, as $R_1$ is added to $V_{rp\_0}^0$, $t_{epu\_0}^0$ and
14      $t_{edo\_0}^0$ of $O_p^0$ need to be updated as well, resulting in $t_{epu\_1}^0$ and $t_{edo\_1}^0$. As simulation time
15      advances, the user of $R_0$ is picked up at the moment $t_{pu}^0$. It is at this precise moment that
16      the realized pick-up time of this user can be truly determined. Similarly, when the user is
17      dropped off at $t_{do}^0$, the realized drop-off time becomes clear. Finally, it can be observed
18      that the pick-up and drop-off time in $O_p^0$ and the realized service $RS_0 = (t_{pu}^0, t_{do}^0)$ have
19      unavoidable deviations due to the addition of other requests.
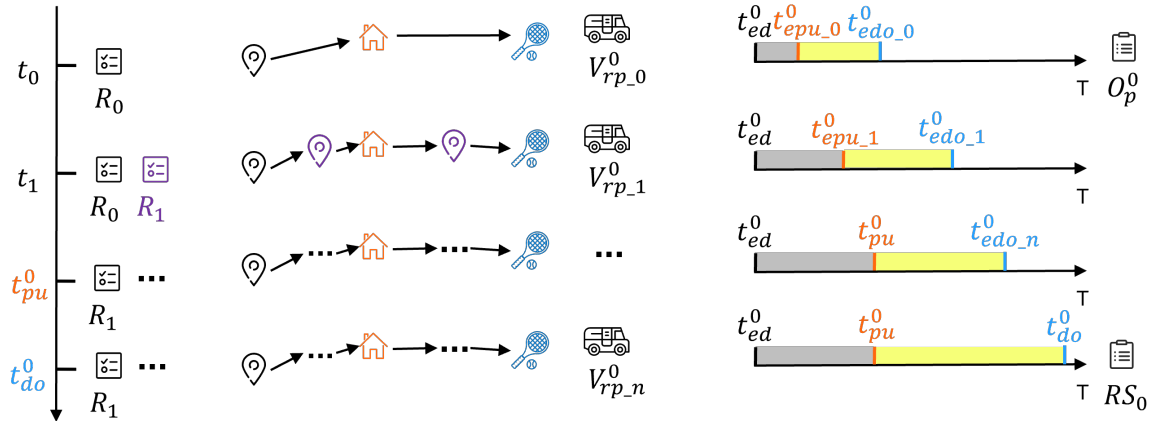


**FIGURE 1**: The dynamic mechanism of offer creation

20      This study investigates a constrained version of the general problem with the following
21  attributes:

22     • Users must be picked up within the maximum waiting time $t_{mw}$, i.e., $(t_{epu} - t_{ed}) \leq t_{mw}$.
23     • The detour time factor generated by the MoD service, $f_{dt} = (t_{edo} - t_{epu} - t_{od})/t_{od}$, where
24      $t_{od}$ denotes the direct travel time from origin to destination, cannot exceed the maximum
25      detour time factor $f_{mdt}$.
26     • The number of passengers $n_{pax}$ cannot exceed the number of vehicle's available seats.
27      If any of these conditions are not met by a vehicle, it will be marked as unable to provide
28  service. If no vehicles in the fleet can provide service, the request will be rejected.

---

[1]Pricing is not in the focus of this study and therefore omitted from now on.

1  **Simulation Framework**
2  In this study, FleetPy (*15*) is chosen to model the dynamics of realistic MoD services — including
3  the interaction of operators with users. During the simulation, offer creation problems appear and
4  have to be solved. FleetPy is an agent-based framework for MoD fleet simulation, as illustrated in
5  Figure 2. Users send travel requests $\mathbf{R}$ to an MoD operator, who then provides offer predictions
6  $\mathbf{O_P^{IR}}$ based on the current fleet state $\mathbf{F_S}$. By default, FleetPy employs an insertion heuristic in the
7  Immediate Response (IR) module, which integrates users' pick-up and drop-off locations into each
8  vehicle's routing plan. Offer attributes are derived from the optimal plan according to the operator's
9  objectives. At predetermined intervals, all accepted offers are batch-optimized and reassigned to
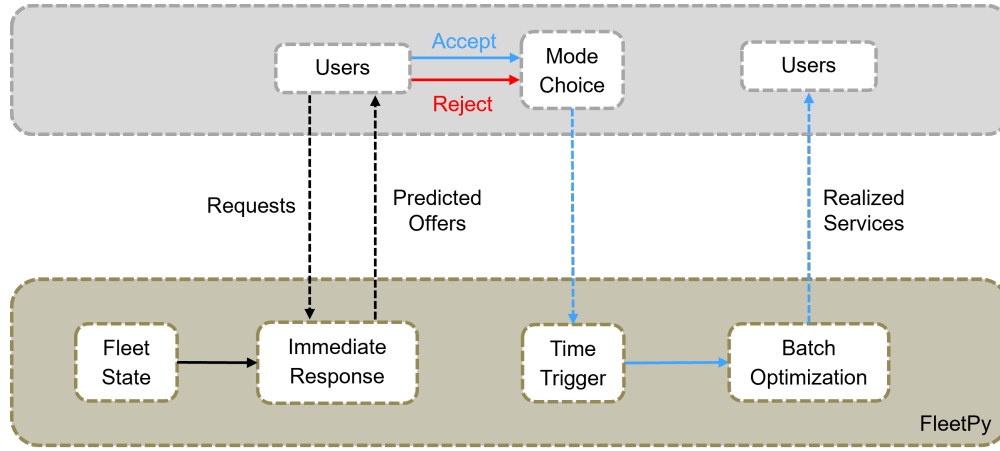10 vehicles to ensure a better system-wide solution.



**FIGURE 2**: The general workflow of FleetPy

11      The fleet state $\mathbf{F_S}$ is described by routing plans $V_{rp}$ for each vehicle. Besides the current
12 position and onboard requests, $V_{rp}$ consists of planned stops, denoted as $s_v^i = (p_{ps}^i, t_a^i, t_d^i, R_o^i, R_a^i)$.
13 Here, $p_{ps}^i$ denotes the position of planned stop $i$, $t_a^i$ denotes the arrival time at planned stop $i$, $t_d^i$
14 is the departure time at planned stop $i$, $R_o^i \subset \mathbf{R}$ is the users on board the vehicle at the planned
15 stop $i$, and $R_a^i \subset \mathbf{R}$ is the users alighting from the vehicle at the planned stop $i$. Collectively, these
16 individual routing plans form the overall fleet state, i.e., $V_{rp} \subset \mathbf{F_S}$.
17      Due to the dynamic nature of the underlying vehicle routing problem, there often is a
18 deviation between FleetPy's $\mathbf{O_p^{IR}}$ and the realized service (**RS**) received by the users.

19 **ML-Model Architecture**
20 This study proposes an ML-based model to predict offers for MoD services. During the simulation,
21 user requests $\mathbf{R}$ and the corresponding fleet status $\mathbf{F_S}$ are recorded as input data for the model. The
22 realized services **RS** provided for each request by fleetpy are documented as target values for the
23 ML module. The model then predicts whether a service can be provided. If it can, the model
24 outputs specific offer attribute values, $\mathbf{O_P^{ML}}$, where $O_p^{ml} = (t_{ew}^{ml}, t_{et}^{ml})$. By directly learning from
25 **RS**, the model can mitigate the effects of dynamic mechanisms. Figure 3 illustrates the overall
26 workflow of the proposed model.
27      The idea of this study's approach is to guide and construct the input features representing
28 vehicle availability for a user request along three dimensions: pick-up likelihood, vehicle capacity,
29 and detour degree. In the data preprocessing phase, the model first encodes the input $\mathbf{R}$ and $\mathbf{F_S}$ in

1  these three dimensions. Before being fed into the ML module, a selection mechanism is used to
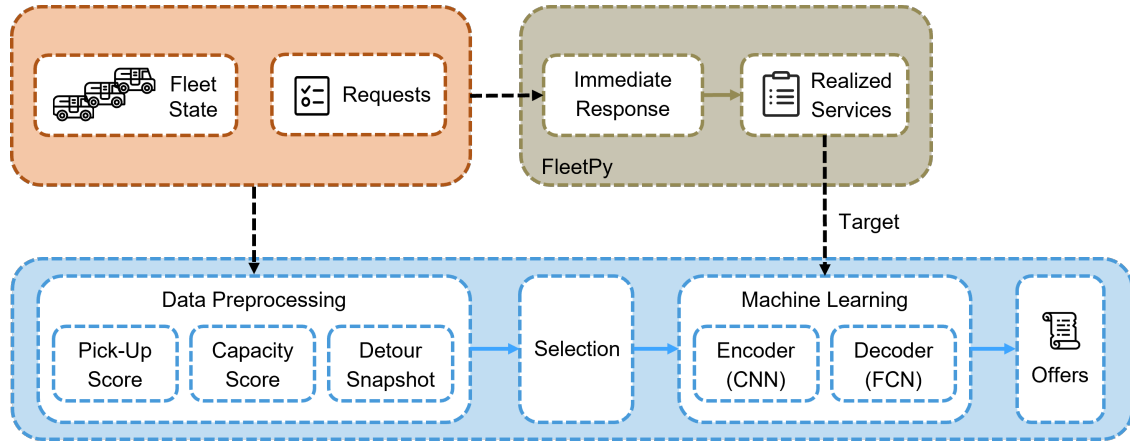2  streamline the encoded data, thereby improving the processing speed and prediction accuracy.



**FIGURE 3**: The overall workflow of the proposed model

3        Offer prediction involves two aspects: recognizing requests that cannot be served and pre-
4  dicting the offer attribute values for the requests that can be served. This approach employs an
5  encoder-decoder architecture to accommodate the specificity of the prediction task. The encoder
6  module comprehends and refines the encoded information to generate request embeddings (RE).
7  For different prediction needs, different decoder modules are used to decode the RE and output the
8  corresponding prediction results. In the following, each module is discussed in detail.

9  *Pick-Up Score*
10  An effective offer needs to ensure that a vehicle can reach the user's departure location within $t_{mw}$.
11  The pick-up score $s_{pu}$ is proposed in Eq. (1) to assess a vehicle's potential to fulfill this requirement.
12  For each planned stop of the vehicle in the assigned routing plan, the vehicle's pick-up ability is
13  scored. The time for a vehicle to travel from a planned stop to pick up a user, $t_{vpu}$, can be computed
14  according to Eq. (2); it consists of two parts: $t_{vo}$, the travel time from the stop to the user's origin,
15  and an additional component representing the time required for the vehicle to depart from the stop.
16  The longer it takes a vehicle to pick up a user, the lower the $s_{pu}$. If the travel time exceeds $t_{mw}$, the
17  vehicle's $s_{pu}$ at that planned stop becomes 0.
18        Calculating the $s_{pu}$ for each of the $m$ planned stops across all $n$ vehicles in the fleet yields
19  a pick-up score table $\mathbf{S_{pu}}$ (3). This table reflects the fleet's overall potential to reach passengers on
20  time.

21  $$s_{pu} = max(0, \frac{t_{mw} - t_{vpu}}{t_{mw}}) \tag{1}$$

22  $$t_{vpu} = t_{vo} + (t_d - t_{ed}) \tag{2}$$

23  $$\mathbf{S_{pu}} = \begin{bmatrix} s^i_{pu_1}, & \cdots, & s^i_{pu_m}, \\ \vdots & & \vdots \\ s^n_{pu_1}, & \cdots, & s^n_{pu_m} \end{bmatrix} \tag{3}$$

1 *Capacity Score*
2 To ensure that vehicles meet users' travel needs, they must provide sufficient seating capacity. The
3 status information of a vehicle at each planned stop includes details about passengers boarding and
4 alighting. Changes in the number of passengers on board directly impact the vehicle's ability to
5 accommodate new riders. This method simplifies the impact of these changes by focusing on the
6 number of available seats $n_{fs}$ as the vehicle departs each stop. The seating capacity is then scored
7 using Eq. (4), where $n_{mfs}$ represents the vehicle's maximum seating capacity. The more empty
8 seats a vehicle has, the higher its score. Similar to the calculation of $\mathbf{S_{pu}}$, the capacity scores for all
9 $n$ vehicles across all $m$ planned stops are computed, forming a capacity score table $\mathbf{S_{cap}}$ (5), which
10 reflects the overall seating capacity of the fleet in response to the request.

11
$$s_{cap} = \frac{n_{fs}}{n_{mfs} - n_{pax}} \tag{4}$$

12
$$\mathbf{S_{cap}} = \begin{bmatrix} s^i_{cap_1}, & \cdots, & s^i_{cap_m}, \\ \vdots & & \vdots \\ s^n_{cap_1}, & \cdots, & s^n_{cap_m} \end{bmatrix} \tag{5}$$

13 *Detour Snapshot*
14 Another crucial factor in evaluating a vehicle's service capability is detour time. This issue does
15 not need consideration for RH, where vehicles and users are uniquely paired. However, in RP,
16 where users share the vehicle with others, direct transportation from origin to destination is often
17 not feasible. Users may need to stop at other passengers' locations along the route, resulting in
18 detours. Although a certain amount of detour is acceptable, as the reduced trip cost from sharing
19 compensates for the increased travel time, excessively long detours can negatively impact the user
20 experience. Therefore, it's essential to factor in detour time for each vehicle and ensure it remains
21 within a reasonable limit to maintain service quality.
22     When constructing a new routing plan for each vehicle, the new planned stops $N_p$ for the
23 vehicle include two types of key points: the planned stops in the vehicle's current routing plan
24 and the origin and destination points of the new user. By connecting these points and using the
25 travel time between them as the edge lengths, we create a network that reflects the connectivity
26 relationships along with travel times. This network represents all potential routes for the vehicle
27 after incorporating the new user, forming what is known as a detour snapshot $S_{dt}$:
28
$$S_{dt} = \{t^{p_i}_{p_j} | p_i, p_j \in N_p, p_i \neq p_j\} \tag{6}$$
29 Here, $t^{p_i}_{p_j}$ denotes the travel time between nodes $p_i$ and $p_j$. By calculating detour snapshots for all
30 vehicles, we compile a detour snapshot table $\mathbf{S_{dt}}$, which reflects the fleet's overall detour degree.

31 *Selection Mechanism*
32 Offer prediction needs to be accomplished in a very short time, so it is necessary to improve
33 the computational efficiency of the model. This approach reduces the computational burden of
34 the model by designing a selection mechanism that streamlines information in both temporal and
35 spatial domains. In the temporal domain, the model reduces the feature dimensions of each vehicle,
36 while in the spatial domain, it filters the number of vehicles considered for offer prediction.
37     When assigning a vehicle to a user request, priority is given to nearby vehicles with suf-
38 ficient capacity. By combining the $\mathbf{S_{pu}}$ and $\mathbf{S_{cap}}$, we infer that a suitable vehicle has a non-zero

1  $s_{pu}$ and $s_{cap}$ at least at one of its planned stops. Therefore, vehicles with a $s_{pu}$ or $s_{cap}$ of 0 at all
2  planned stops are filtered out from subsequent calculations. Additionally, each vehicle's maximum
3  $s_{pu}$ is sorted, and only the top $n_{mv}$ vehicles are selected as candidates, retaining those closest to the
4  user's departure point.
5       Each vehicle may have multiple planned stops in its routing plan. Since these stops are
6  visited sequentially, the likelihood of a vehicle being able to arrive at these stops and then reach
7  the user's departure point within $t_{mw}$ diminishes for later stops. Therefore, this approach only
8  considers the first $n_{mps}$ stops of each vehicle, ignoring the subsequent stops to reduce the feature
9  dimensions of each vehicle.

10  *The Encoder-Decoder Architecture*
11  The encoder-decoder architecture is considered one of the most rational neural network designs
12  and is widely used in deep learning, particularly in models like the Transformer (*32*). The encoder
13  is responsible for understanding and processing the input information to create information em-
14  beddings. The decoder, on the other hand, provides the ML model with the flexibility to handle
15  various tasks. Depending on the task requirements, different decoder networks, tailored to the
16  specific task scenario, can be employed and combined with the encoder to form a complete neural
17  network prediction module.
18       The encoded information of the user's request and the corresponding fleet state is stored
19  in a matrix, where each row represents a vehicle, and the columns contain encoded information in
20  three dimensions: pick-up likelihood, vehicle capacity, and detour degree. This matrix serves as
21  a snapshot of the current fleet state, akin to an image that can be processed using machine vision
22  techniques. In this study, a Convolutional Neural Network (CNN) (*33*) is selected for extracting
23  and embedding the encoded information due to its effectiveness in handling such image-like data.
24       For each request, the offer prediction needs to be performed at most twice. First, it must be
25  determined whether the request can be adequately served with the current fleet state. If it can, the
26  specific offer attribute values need to be predicted. To meet these two-aspect prediction require-
27  ments, two different Fully Connected Networks (FCN) (*34*) are employed: one for identifying
28  rejected requests and the other for predicting the values of specific offer attributes.

29  **CASE STUDY**
30  This section outlines the experimental design of the study, presents the results, and concludes with
31  a brief discussion of the findings.

32  **Experimental Design**
33  *FleetPy*
34  To train the ML model, we configure FleetPy's simulation environment and construct the dataset
35  based on the results of FleetPy's runs. Table 1 shows all the key configurations of FleetPy. It is
36  important to note that $t_{mw}$ is set to 600 seconds, meaning that all users must be picked up within
37  10 minutes after their $t_{ed}$. Additionally, $f_{mdt}$ is set to 40%, indicating that the detour time incurred
38  by the RP services must be at most 40% of the $t_{od}$; otherwise, the user's request will be rejected.
39  All vehicles in the fleet have a maximum capacity $n_{mpax}$ of 4 passengers.
40       The operating area for the MOD service is confined to Manhattan, NYC. The network is
41  extracted from OpenStreetMap. To reflect real road traffic conditions, all link travel times are dy-
42  namically scaled based on the simulation time. The shortest paths between nodes are precomputed,

1  and a lookup table is created to represent the travel times between all origin and destination pairs.
2      The MoD demand in this study is derived from the New York City taxi dataset (*4*), a widely
3  used open-source resource. We utilize taxi trip data from November 11 to November 18, 2018,
4  selecting trips that both start and end in Manhattan. To reduce computational complexity, 10% of
5  these trips are randomly sampled and treated as RP demand, resulting in a total of 320,434 trip
6  requests. To ensure robustness, we generate three sets of trip data using three different random
7  seeds. All results presented in the following sections are the averaged outcomes from these three
8  datasets.
9      Additional details on the processing of network and trip data can be found in the work of
10 Engelhardt et al. (*35*).

**TABLE 1**: Configurations of FleetPy

| Parameter | Symbol | Value |
|---|---|---|
| Simulation duration | $t_{sim}$ | 24 h |
| Simulation step | $t_{step}$ | 1 s |
| Fleet size | $f_s$ | 200 |
| Maximum waiting time | $t_{mw}$ | 600 s |
| Maximum detour time factor | $f_{dt}$ | 40% |
| Demand | $D$ | Manhattan taxi data, from Nov. 11 to Nov. 18 2018 |
| Network | $G$ | Manhattan network |
| Vehicle capacity | $n_{mpax}$ | 4 |

11 *Dataset*
12 During the FleetPy simulation, the fleet state at $t_{ed}$ of each request is recorded and combined
13 with the request information to form the feature dataset. At the end of the simulation, acceptance
14 or rejection of requests, realized waiting times, and realized travel times are used to construct
15 the corresponding target dataset. Data from the entire week of November 12 to November 18 is
16 utilized to build the training and testing sets, with a split ratio of 70% for training and 30% for
17 testing. Additionally, data from the entire day of November 11 is used to construct the validation
18 set. Figure 4 shows the distribution of features and targets for each dataset.
19      It can be observed that the rejection rate of user requests in all the datasets is approximately
20 14%, indicating that the FleetPy configurations are sufficient to meet the travel demand used in this
21 study. The distributions of travel time and waiting time are consistent across all datasets, with the
22 mean travel time being around 644 seconds and the mean waiting time approximately 371 seconds.

23 *Scenario Specification*
24 The role of the selection mechanism is primarily defined by the parameters $n_{mv}$ and $n_{mps}$. In the
25 spatial domain, $n_{mv}$ controls the maximum number of vehicles the model considers. In the temporal
26 domain, $n_{mps}$ limits the scope of the model's perspective on the vehicle routing plan.
27      By analyzing the FleetPy simulation results, Figure 5 illustrates the distribution of the num-
28 ber of available vehicles corresponding to requests and the order of the planned stop at which newly
29 incoming requests are inserted into existing vehicle routes. It can be observed that while the fleet
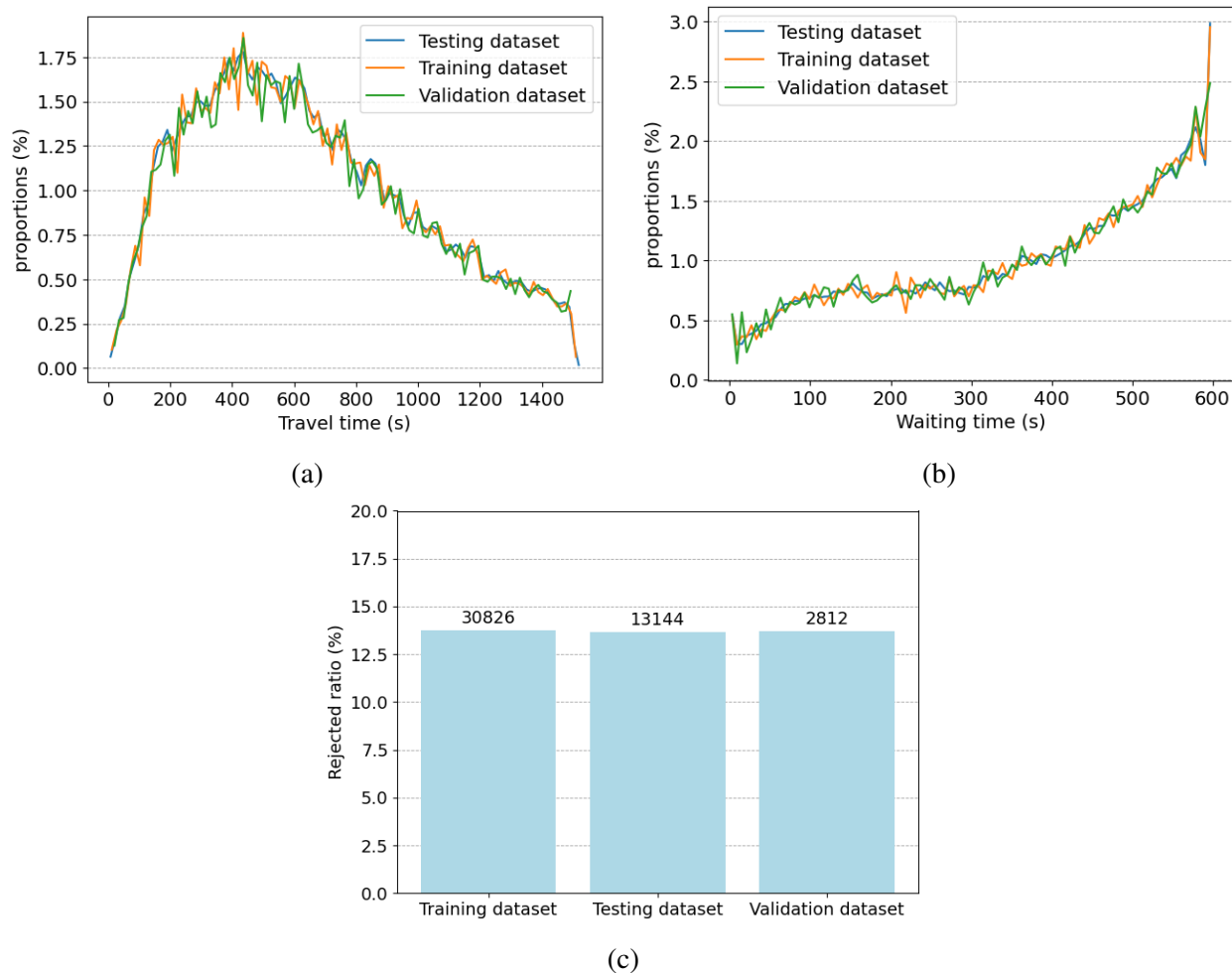
(a)

(b)

(c)

**FIGURE 4**: The distribution of features and targets for each dataset. (a) shows the distribution of travel time, (b) shows the distribution of waiting time, and (c) shows the ratio of requests rejected by FleetPy in each dataset, with the numbers above indicating the quantities.

1   size is 200 vehicles, over 95% of requests only have at most 60 vehicles potentially offering ser-
2   vice. Moreover, most requests are picked up before the vehicle's fifth planned stop.
3        This study designs various scenarios from these two aspects to explore the performance
4   of the proposed model and selection mechanism, respectively. Table 2 summarizes the model
5   configurations for these scenarios. In the base scenario (BS), $n_{mv}$ is set to 60 vehicles, and $n_{mps}$
6   is set to 5 planned stops. NV-20, NV-40, NV-100, and NV-200 represent different scales of $n_{mv}$,
7   ranging from considering only the 20 vehicles in the user's neighborhood to considering the entire
8   fleet to examine the role of the selection mechanism in the spatial domain. NPS-2, NPS-3, NPS-4,
9   NPS-6, and NPS-7 explore the role of selection mechanisms in the temporal domain, ranging from
10   considering only the 2 planned stops of each vehicle to considering 5 planned stops.
11        Each of the aforementioned scenarios includes two ML modules: one decoder is designed
12   as an FCN dealing with the classification problem (i.e. user acceptance or rejection), and the other
13   is an FCN capable of outputting two predicted offer attribute values (i.e. waiting and travel time).
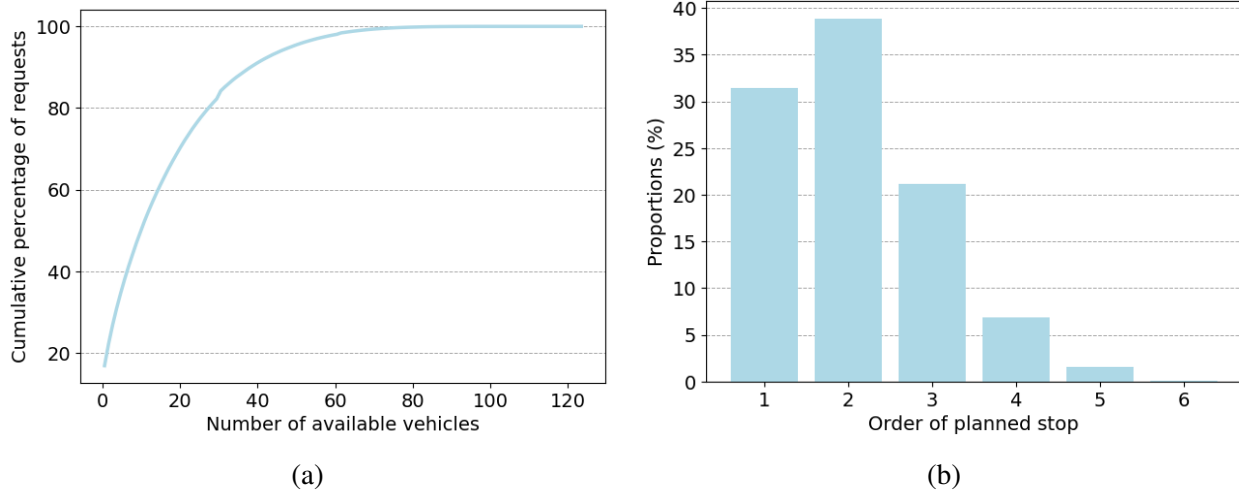14   The encoder of both modules is a CNN. For the classification problem, the ML module employs

(a)



(b)

**FIGURE 5**: The distribution of available vehicle numbers and valid planned stop order. (a) shows the number of available vehicles for each request, and (b) illustrates the order of the planned stop at which the newly incoming request is served.

**TABLE 2**: Scenario Configurations

| Scenario name | $n_{mv}$ | $n_{mps}$ |
|---|---|---|
| Base scenario (BS) | 60 | 5 |
| NV-200 | 200 | 5 |
| NV-100 | 100 | 5 |
| NV-40 | 40 | 5 |
| NV-20 | 20 | 5 |
| NPS-7 | 60 | 7 |
| NPS-6 | 60 | 6 |
| NPS-4 | 60 | 4 |
| NPS-3 | 60 | 3 |
| NPS-2 | 60 | 2 |

1 Binary Cross Entropy (BCE) as the loss function, while for the offer attribute values prediction, it
2 uses Mean Squared Error (MSE). The model is trained using Adaptive Moment Estimation (Adam)
3 (*36*) as the optimizer, with a learning rate of 0.001.

4 *Key Performance Indicators*
5 To compare the model's prediction performance across scenarios, Key Performance Indicators
6 (KPIs) need to be defined. The results are compared with both $\mathbf{O_P^{IR}}$ and **RS** to evaluate the model's
7 performance. When predicting request acceptance, it is crucial to avoid misclassifying unfeasible
8 requests as feasible, as this significantly impacts users' ride experiences. In the dataset, served
9 requests are labeled as 1 (positive), while rejected requests are labeled as 0 (negative). Therefore,

1  the model should prioritize accurately identifying rejected requests by increasing true negative
2  (TN) counts and decreasing false positive (FP) counts. Additionally, true positive (TP) and false
3  negative (FN) requests should also be considered to enhance the overall benefit for TNCs.
4          The True Negative Rate (TNR) (7) assesses the model's ability to recognize rejected re-
5  quests. Precision (8) evaluates the model's prediction accuracy for served requests, while Recall
6  (9) indicates the model's tendency to underreport. However, since approximately 86% of requests
7  are served, relying solely on a single metric does not fully capture the model's performance on pos-
8  itive cases. Therefore, the F1-Score (10), which combines Precision and Recall, is more suitable
9  for assessing the model's overall performance.
10         For predicting the attribute values of offers, two key features are considered: waiting time
11  and travel time. The Mean Absolute Error (MAE) of these two attributes is used to evaluate
12  the model's overall prediction performance. Additionally, heat maps of the joint distribution of
13  absolute errors for these two attributes are used to visualize the model's prediction accuracy at the
14  individual offer level.
15         Beyond accuracy, the efficiency of the prediction process is also crucial. The offer predic-
16  tion time of the proposed model, $t_{op}^{ml}$, is the combined duration of data preprocessing, data selection,
17  and ML processing. This approach provides a realistic and comprehensive view of the model's ef-
18  ficiency, making the evaluation practical and applicable to real-world applications. Based on the
19  processing time of the IR module, $t_{op}^{ir}$, the Efficiency Improvement Rate (EIR) (11) is used as a
20  KPI to evaluate the improvement in predictive efficiency of the proposed model.

21  $$\text{TNR} = \frac{TN}{TN + FP} \tag{7}$$

22  $$\text{Precision} = \frac{TP}{TP + FP} \tag{8}$$

23  $$\text{Recall} = \frac{TP}{TP + FN} \tag{9}$$

24  $$\text{F1} = 2 \times \frac{Precision \times Recall}{Precision + Recall} \tag{10}$$

25  $$\text{EIR} = \frac{t_{op}^{ir} - t_{op}^{ml}}{t_{op}^{ir}} \times 100 \tag{11}$$

26  **Results**
27  This section discusses the model's results across various scenarios, focusing on three key aspects
28  of offer prediction. The first part analyzes the model's ability to accurately recognize rejected
29  requests. The second part evaluates the model's accuracy in predicting offer attribute values. The
30  third part examines the model's operational efficiency.

31  *Recognition*
32  According to the FleetPy configuration used in this study, requests are determined to be serviceable
33  only by the IR module. Therefore, this study assumes the recognition results of the IR module
34  as ground truth and trains and validates the model based on these results. Figure 6 shows the
35  recognition performance of the model in all scenarios.
36         In BS, the model achieves a TNR of approximately 90%, indicating its effectiveness in
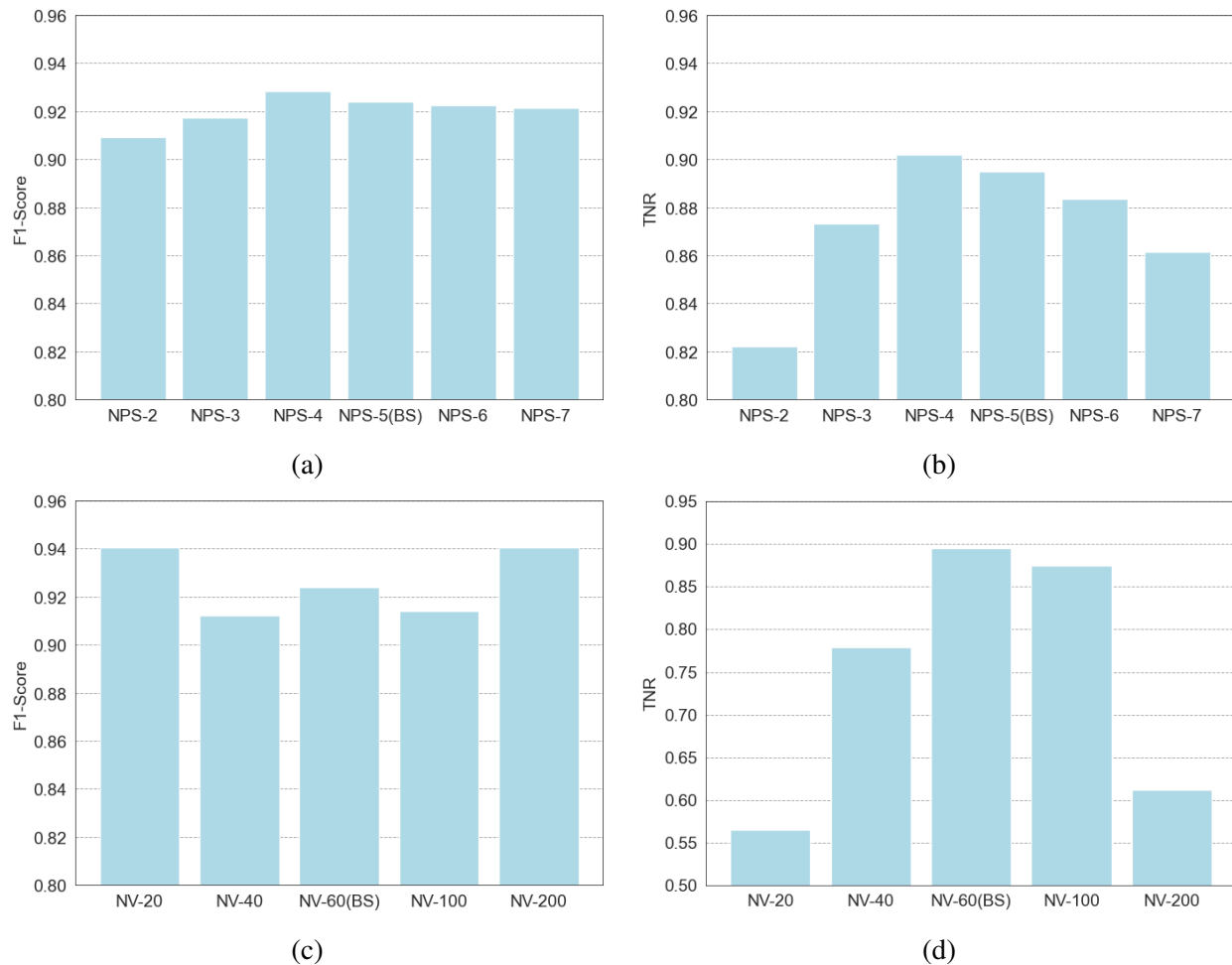
(a)



(b)



(c)



(d)

**FIGURE 6**: F1-Score and TNR for all scenarios. (a) and (b) illustrate the model's recognition performance for different $n_{mps}$, and (c) and (d) show the effect of $n_{mv}$ on the model's performance.

1  recognizing the majority of requests that need to be rejected, despite the rejection rate being only
2  around 14% in the dataset. NPS-4 slightly outperforms BS in temporal domain selection. As
3  shown in Figure 5, only about 2% of the requests are served at the fifth planned stop of the vehicle.
4  This suggests that ignoring this stop allows the model to focus on the first four stops, which contain
5  more relevant information. As the number of planned stops considered by the selection mechanism
6  decreases, the model's recognition rate for rejected requests declines. This reduction is attributed
7  to the loss of information in the shortened field of view, leaving the ML model with insufficient
8  features to make accurate judgments. Specifically, nearly 22% of the requests are served at the
9  vehicle's third planned stop, and ignoring this stop's feature information in NPS-2 results in a TNR
10 of about 82%. In scenarios where the selection mechanism considers additional planned stops,
11 such as NPS-6 and NPS-7, the model's TNR decreases. This decline occurs because all requests
12 in the dataset are serviced by or before the vehicle's fifth planned stop, making information about
13 subsequent stops a distraction that interferes with the model's judgment.
14       The model exhibits similar performance trends in selecting information in the spatial do-
15 main. Overly reducing or increasing the number of considered vehicles diminishes the model's

1  TNR. Notably, the model's performance exhibits different sensitivities to information selection
2  across the spatial and temporal domains, with spatial information selection having a greater im-
3  pact. When the amount of temporal information is scaled down to 70%, as in NPS-2, the model still
4  achieves a TNR higher than 82%. Conversely, when spatial information is scaled down to 70%, as
5  in NV-20, the model's accuracy in correctly identifying rejected requests drops to approximately
6  57%, which is only slightly better than random guessing.
7         This pattern is also observed in the model's ability to determine which requests can be
8  served. Compared to BS, NPS-4, which considers one less planned stop, achieves a higher F1-
9  Score of over 0.92. As $n_{mps}$ increases or decreases, the model's F1-Score decreases. The selection
10 mechanism in the spatial domain causes more significant changes in the F1-Score. Particularly in
11 NV-20 and NV-200, although the F1-Score exceeds 0.94 in both scenarios, higher than in all other
12 scenarios, this occurs because the model either receives too little valid information or too much
13 irrelevant information, impairing its ability to make accurate judgments. Consequently, the model
14 tends to classify more requests as serviceable, which is a more prevalent category in the dataset.

15 *Offer Attributes*
16 In FleetPy, $\mathbf{O_P^{IR}}$ is the prediction of **RS**. Therefore, in this study, **RS** is considered the ground
17 truth on which the model's training and validation are based. By comparing the mean error of the
18 model's predicted values with **RS**, we can evaluate the model's predictive performance. Addition-
19 ally, by comparing the predicted values of the model with $\mathbf{O_P^{IR}}$, we can determine how optimized
20 the proposed model is compared to the IR module. Figure 7 illustrates these two comparisons of
21 the model for all scenarios.
22         In the FleetPy configuration, $t_{mw}$ is set to 600 seconds. In BS, the model's mean prediction
23 error for this attribute is approximately 130 seconds. The best prediction result is observed in
24 NPS-4, where one less planned stop is considered in the temporal domain, reducing the error to
25 around 112 seconds. The model's prediction accuracy declines as the selection mechanism further
26 scales down the data in the temporal domain. It is notable that reducing the number of planned
27 stops considered has a more detrimental effect on the model's prediction accuracy than increasing
28 the temporal domain information.
29         In the spatial domain, both reducing or increasing the number of vehicles considered by the
30 model diminishes its predictive performance. Particularly in NV-20 and NV-200, insufficient fleet
31 information fails to provide adequate support for the model's predictions, while excessive infor-
32 mation dilutes the model's focus. This demonstrates that the selection mechanism can effectively
33 filter out key information to aid the model in predicting accurate offer attribute values. Addition-
34 ally, spatial information selection has a more significant impact on the model's prediction accuracy
35 compared to temporal information selection.
36         In the dataset, the average travel time for all requests is about 644 seconds. In BS, the
37 model's mean prediction error for this attribute is around 60 seconds. The best prediction result is
38 again observed in NPS-4, with an error of approximately 52 seconds. The impact of the selection
39 mechanism on the model's prediction accuracy for travel time, whether in the temporal or spatial
40 domain, follows the same trend as its impact on the prediction accuracy for waiting time. The
41 worst result is observed in NV-20, where the mean error reaches about 135 seconds.
42         The model's prediction results are compared with $\mathbf{O_P^{IR}}$. For travel time prediction, the IR's
43 mean prediction error is approximately 85 seconds, while the proposed model's prediction error
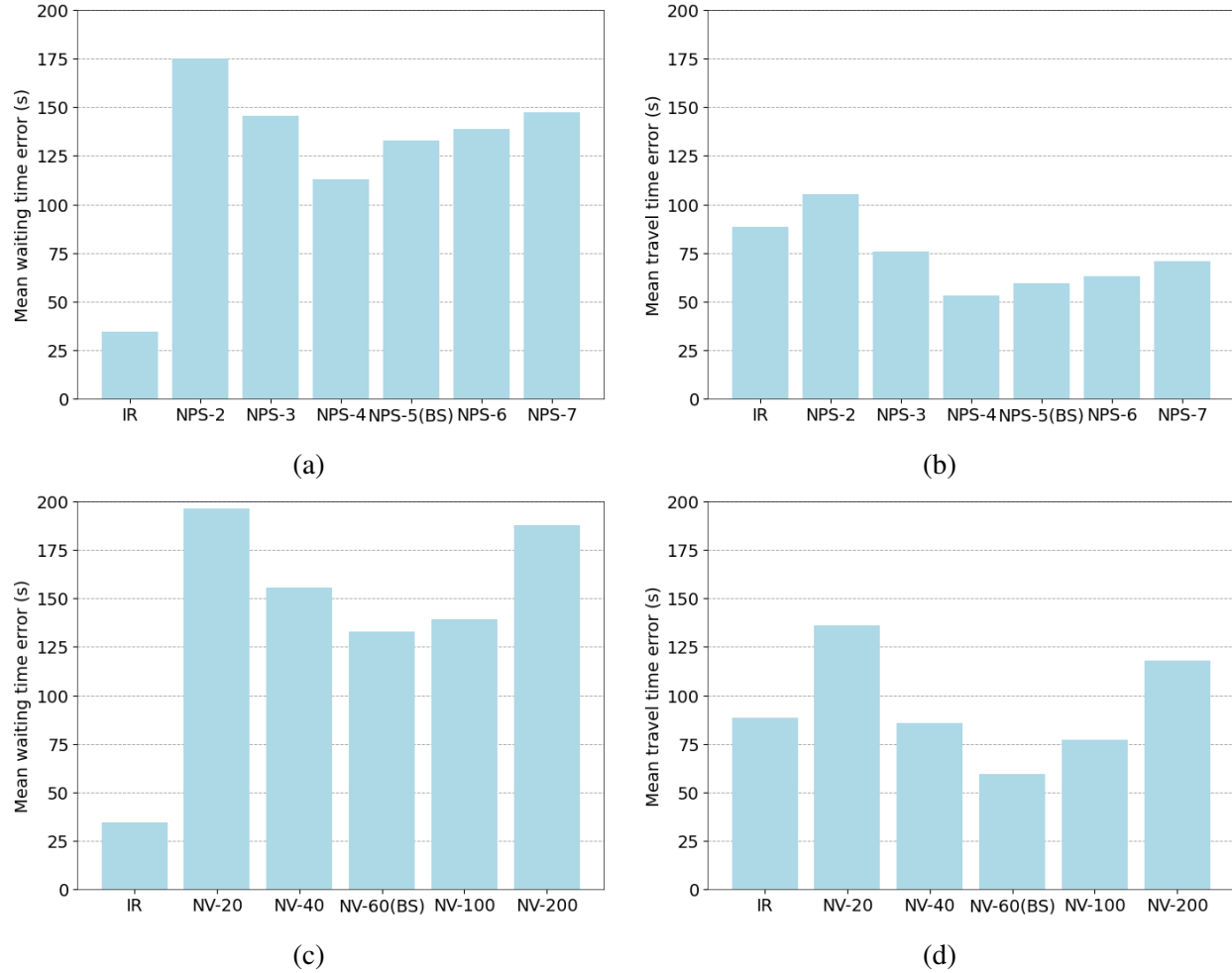44 is consistently lower across most scenarios, with the smallest error of about 52 seconds in NPS-4.

**FIGURE 7**: Comparison of mean errors in waiting time and travel time. Figures (a) and (b) show the mean errors of the model with different $n_{mps}$ values; Figures (c) and (d) show the mean errors of the model with different $n_{mv}$ values.

1   This indicates that the model outperforms the IR module in predicting travel time. For waiting time
2   prediction, the IR's mean prediction error is around 37 seconds, whereas the proposed model's
3   minimum prediction error is 112 seconds. This shows that the model's prediction accuracy for
4   waiting time is not as good as that of the IR. Theoretically, waiting time is primarily influenced
5   by $\mathbf{S_{pu}}$ in the fleet encoding information. However, in this study, the proposed model employs a
6   relatively simple linear equation to calculate $s_{pu}$. This linear mapping may not fully reflect the
7   complexities of real-world vehicle pick-up capabilities.
8           Figures 8 and 9 illustrate the joint distribution of absolute errors for waiting time and travel
9   time at the individual offer level. The horizontal axis represents the absolute prediction error for
10  waiting time, while the vertical axis represents the absolute prediction error for travel time for each
11  offer. From both figures, it can be observed that scenarios with lower mean errors also have more
12  concentrated joint absolute error distributions at the individual offer level. The heat map for NPS-
13  4, in particular, shows the most concentrated highlighted area in the lower left corner, indicating
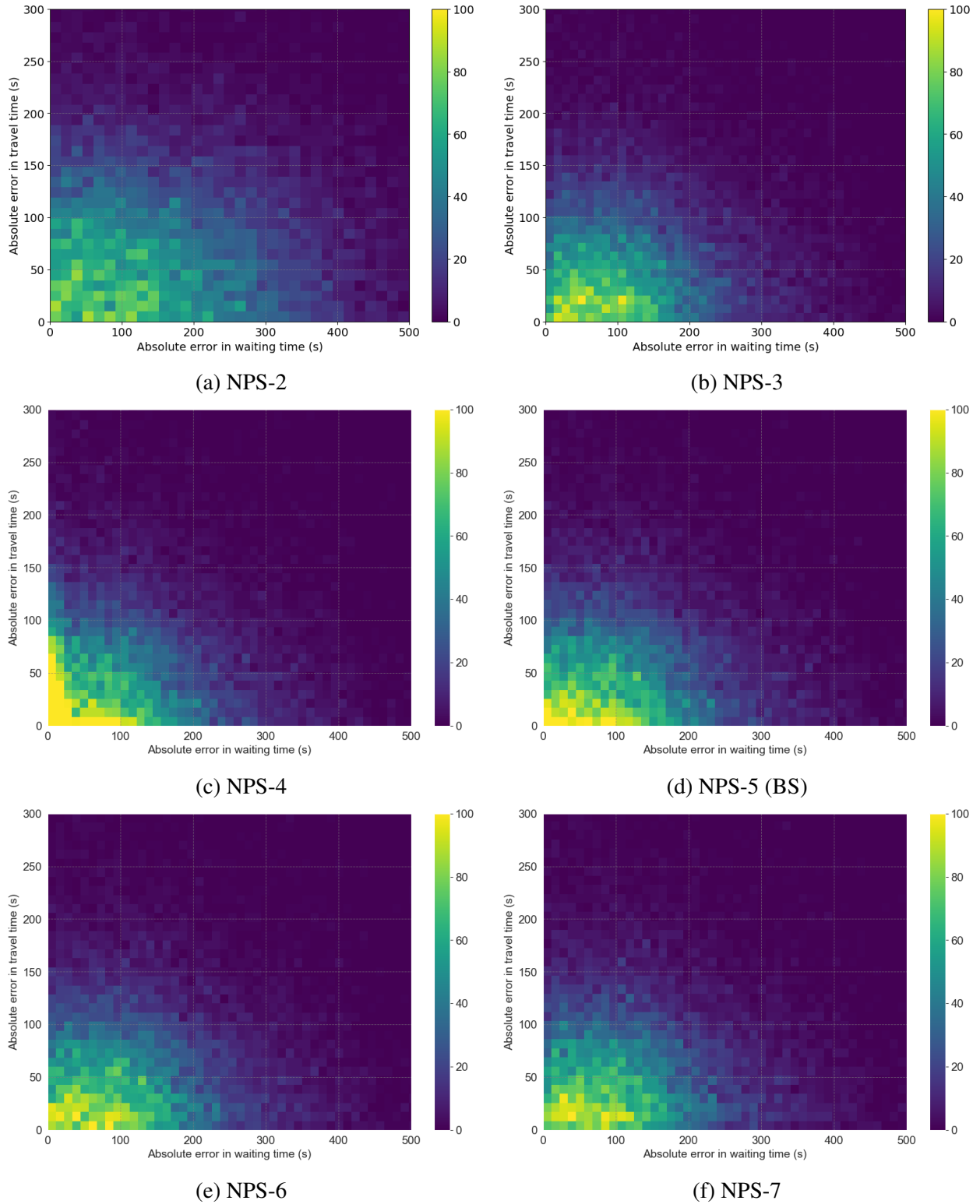14  that most offers have relatively low absolute errors for both attributes.

(a) NPS-2

(b) NPS-3

(c) NPS-4

(d) NPS-5 (BS)

(e) NPS-6

(f) NPS-7

**FIGURE 8**: Impact of $n_{mps}$ on model performance. Figures (a) to (f) present the joint distribution of absolute errors in the model's predictions across various scenarios characterized by different $n_{mps}$ values.
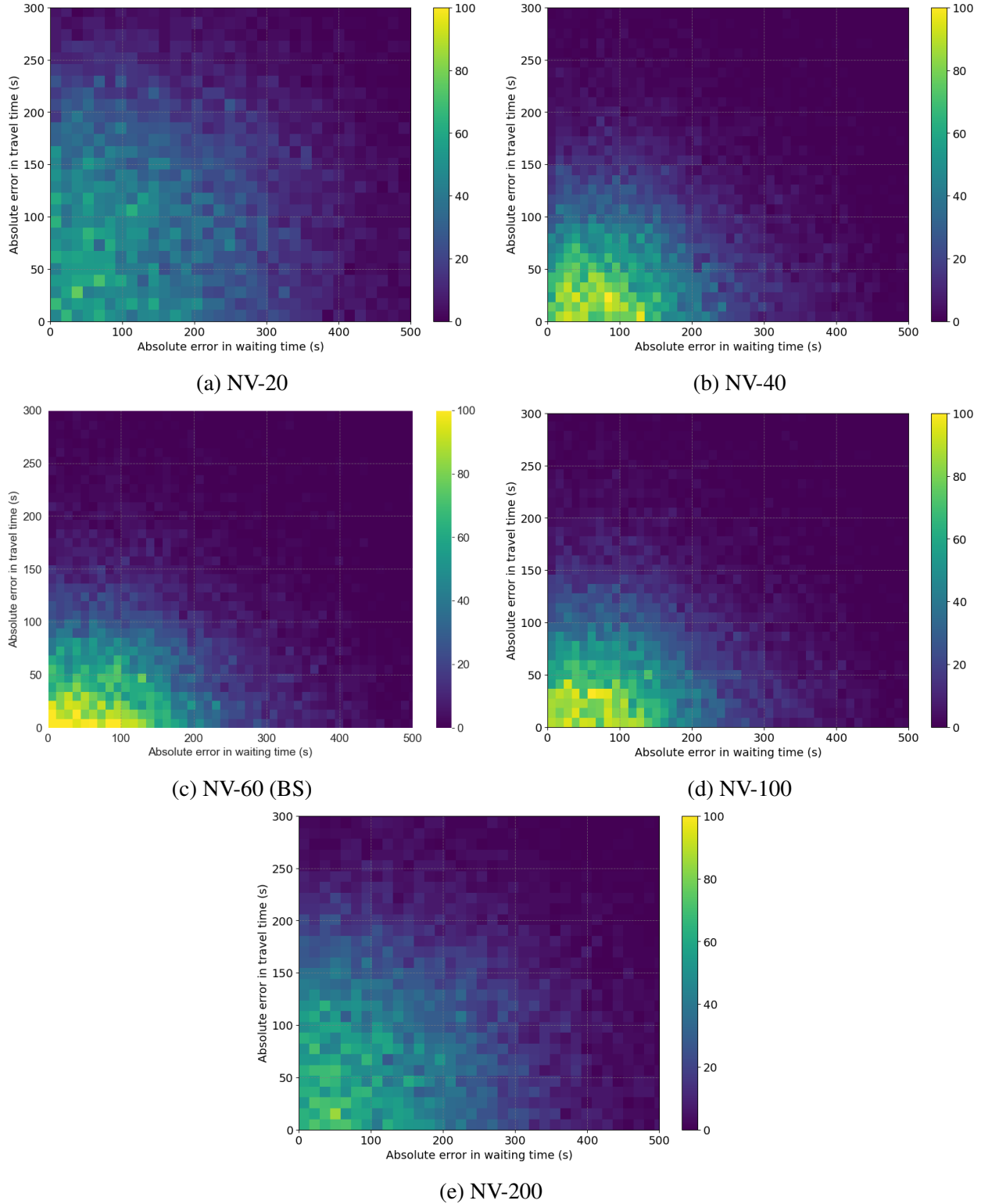
(a) NV-20

(b) NV-40

(c) NV-60 (BS)

(d) NV-100

(e) NV-200

**FIGURE 9**: Impact of $n_{mv}$ on model performance. Figures (a) to (f) present the joint distribution of absolute errors in the model's predictions across various scenarios characterized by different $n_{mv}$ values.

1  *Efficiency*
2  Another important indicator of the model's performance is efficiency. Using EIR, the predictive
3  efficiency of the proposed model is compared to that of the IR module. Figure 10 presents the
4  detailed comparison results.
5      Overall, by leveraging the selection mechanism to reduce the amount of information, the
6  proposed model improves prediction efficiency by approximately 47% to 78% in most scenarios
7  compared to the IR module. However, in NV-200, where the model needs to encode data from
8  all vehicles, prediction efficiency is significantly impacted, resulting in an average prediction time
9  of about 8% longer than that of the IR module. Notably, the model also does not achieve optimal
10  prediction performance under NV-200.
11      The trend in efficiency changes indicates that reducing the amount of information in both
12  the temporal and spatial domains generally enhances the model's prediction efficiency. However,
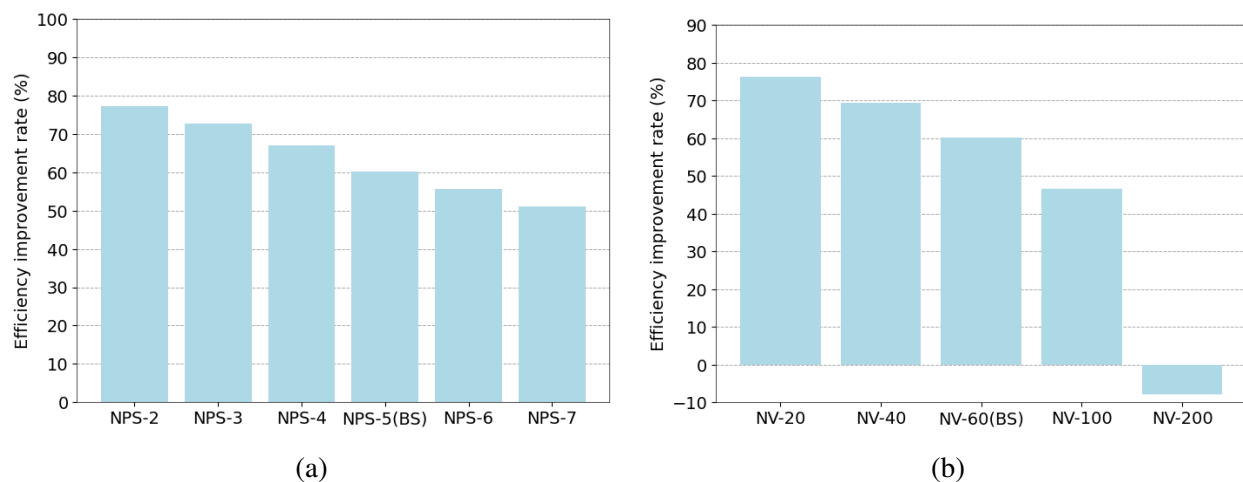13  considering prediction accuracy, indiscriminately reducing information is not recommended.



FIGURE 10: Comparison of EIR. (a) shows the EIR for the model across different values of $n_{mps}$;
(b) shows the EIR for the model across different values of $n_{mv}$.

14  **CONCLUSION AND FUTURE WORK**
15  Ensuring the quality of services is crucial for TNCs like Uber, Lyft, and Didi, as it directly impacts
16  user experience and revenue. A key challenge lies in providing fast and precise offer predictions
17  based on user requests and fleet conditions. Response strategies solving vehicle routing problems
18  for each incoming request struggle to deliver timely predictions in high-demand scenarios. The
19  exponential growth in MoD demand and fleet sizes further complicates this task, underscoring the
20  need for more efficient and scalable predictive methods.
21      This study proposes a machine learning-based model to effectively address these issues.
22  The model encodes fleet state combined with user request information from three dimensions:
23  pick-up likelihood, vehicle capacity, and detour degree, resulting in pick-up score, capacity score,
24  and detour snapshot. By designing a selection mechanism, the model filters the encoded informa-
25  tion in both temporal and spatial domains. This refined encoded information reduces the computa-
26  tional burden on the model, enhancing prediction efficiency and preventing the model's attention
27  from being dispersed, thereby improving prediction accuracy. Given that offer prediction involves

two aspects: identifying unserviceable requests and predicting specific offer attribute values, the model employs an encoder-decoder architecture. Fleet status is treated as an image, and a CNN is used in the encoder to embed image information. Different FCNs are then used to decode the embedded information, making various predictions.

Using FleetPy to simulate MoD services, a case study is conducted based on the Manhattan taxi dataset. Different scenarios are tested by controlling the selection mechanism's intensity in the temporal and spatial domains. Analysis and comparison of model results show that the selection mechanism critically impacts the model's predictive performance. Moreover, spatial domain information has a greater impact on the model's predictive performance than temporal domain information. When considering only 60 vehicles near the user and only the status of the first four stops for each vehicle, the model achieves a True Negative Rate of approximately 90% in identifying rejected requests, effectively recognizing requests the fleet cannot serve. Additionally, it outperforms the heuristic approach in both travel time prediction accuracy and prediction efficiency. However, the model's prediction error for waiting time is worse than the heuristic approach. This discrepancy is attributed to the model's use of a simple linear equation for encoding vehicle pick-up capability, which has limited expressive power.

Future research can improve the model's prediction performance for waiting time by using nonlinear encoding for vehicle pick-up capability. Additionally, different encoders and decoders can be tested to study the impact of various neural networks on the model's predictive ability. Furthermore, the model can be applied to different FleetPy configurations or tested under different networks and demands to examine its transferability. Finally, fleet control algorithms need to be adapted to work with a machine learning-based offer module; unlike the insertion heuristic-based offer module currently employed in FleetPy, machine learning-based approaches will accept requests, for which — according to the original constraints — no feasible vehicle-routing solution can be found.

## AUTHOR CONTRIBUTIONS
The authors confirm their contribution to the paper as follows: study conception and design: Chenhao Ding, Florian Dandl, Klaus Bogenberger; methodology and software: Chenhao Ding, Florian Dandl; data collection: Chenhao Ding; analysis and interpretation of results: Chenhao Ding, Florian Dandl; draft manuscript preparation: Chenhao Ding, Florian Dandl. All authors reviewed the results and approved the final version of the manuscript.

# REFERENCES

1. Uber Technologies, Inc., *Uber Announces Results for Fourth Quarter and Full Year 2023*. Uber Technologies, Inc., San Francisco, 2024.
2. DiDi Global Inc., *DiDi Announces Results for Fourth Quarter and Full Year 2023*. DiDi Global Inc., Beijing, 2024.
3. Kumar, A., A. Gupta, M. Parida, and V. Chauhan, Service quality assessment of ride-sourcing services: A distinction between ride-hailing and ride-sharing services. *Transport Policy*, Vol. 127, 2022, pp. 61–79.
4. NYC Taxi & Limousine Commission, *TLC Trip Record Data*. NYC Taxi & Limousine Commission, New York, 2018.
5. Guo, G. and X. Li, A scientometric review of mobility-on-demand car-sharing systems. *IEEE Intelligent Transportation Systems Magazine*, Vol. 15, No. 1, 2022, pp. 212–229.
6. Shaheen, S. and A. Cohen, Mobility on demand (MOD) and mobility as a service (MaaS): Early understanding of shared mobility impacts and public transit partnerships. In *Demand for emerging transportation systems*, Elsevier, 2020, pp. 37–59.
7. Paparella, F., L. Pedroso, T. Hofman, and M. Salazar, A time-invariant network flow model for two-person ride-pooling mobility-on-demand. In *2023 62nd IEEE Conference on Decision and Control (CDC)*, IEEE, 2023, pp. 4118–4123.
8. Pavone, M., S. L. Smith, E. Frazzoli, and D. Rus, Robotic load balancing for mobility-on-demand systems. *The International Journal of Robotics Research*, Vol. 31, No. 7, 2012, pp. 839–854.
9. Aimsun, S., *Aimsun Ride*, 2022, aimsun Ride Research Program.
10. Group, P. et al., *PTV MaaS Modeller*, 2022.
11. Ruch, C., S. Hörl, and E. Frazzoli, Amodeus, a simulation-based testbed for autonomous mobility-on-demand systems. In *2018 21st international conference on intelligent transportation systems (ITSC)*, IEEE, 2018, pp. 3639–3644.
12. Horni, A., K. Nagel, and K. W. Axhausen (eds.) *The Multi-Agent Transport Simulation MATSim*. Ubiquity Press, London, 2016.
13. Auld, J., M. Hope, H. Ley, V. Sokolov, B. Xu, and K. Zhang, POLARIS: Agent-based modeling framework development and implementation for integrated travel demand and network and operations simulations. *Transportation Research Part C: Emerging Technologies*, Vol. 64, 2016, pp. 101–116.
14. Kucharski, R. and O. Cats, MaaSSim–agent-based two-sided mobility platform simulator. *arXiv preprint arXiv:2011.12827*, 2020.
15. Engelhardt, R., F. Dandl, A.-A. Syed, Y. Zhang, F. Fehn, F. Wolf, and K. Bogenberger, *FleetPy: A Modular Open-Source Simulation Tool for Mobility On-Demand Services*, 2022.
16. Fagnant, D. J. and K. M. Kockelman, The travel and environmental implications of shared autonomous vehicles, using agent-based model scenarios. *Transportation Research Part C: Emerging Technologies*, Vol. 40, 2014, pp. 1–13.
17. Zhang, W. and S. Guhathakurta, Parking spaces in the age of shared autonomous vehicles: How much parking will we need and where? *Transportation Research Record*, Vol. 2651, No. 1, 2017, pp. 80–91.

1    18.    Alonso-Mora, J., S. Samaranayake, A. Wallar, E. Frazzoli, and D. Rus, On-demand high-
2           capacity ride-sharing via dynamic trip-vehicle assignment. *Proceedings of the National*
3           *Academy of Sciences*, Vol. 114, No. 3, 2017, pp. 462–467.
4    19.    Kucharski, R. and O. Cats, Exact matching of attractive shared rides (ExMAS) for system-
5           wide strategic evaluations. *Transportation Research Part B: Methodological*, Vol. 139,
6           2020, pp. 285–310.
7    20.    Engelhardt, R., F. Dandl, and K. Bogenberger, Speed-up heuristic for an on-demand ride-
8           pooling algorithm. *arXiv preprint arXiv:2007.14877*, 2020.
9    21.    Hyland, M. and H. S. Mahmassani, Dynamic autonomous vehicle fleet operations:
10          Optimization-based strategies to assign AVs to immediate traveler demand requests. *Trans-*
11          *portation Research Part C: Emerging Technologies*, Vol. 92, 2018, pp. 278–297.
12   22.    Fleckenstein, D., R. Klein, and C. Steinhardt, Recent advances in integrating demand man-
13          agement and vehicle routing: A methodological review. *European Journal of Operational*
14          *Research*, Vol. 306, No. 2, 2023, pp. 499–518.
15   23.    Adnan, M., F. C. Pereira, C. M. L. Azevedo, K. Basak, M. Lovric, S. Raveau, Y. Zhu,
16          J. Ferreira, C. Zegras, and M. Ben-Akiva, Simmobility: A multi-scale integrated agent-
17          based simulation platform. In *95th Annual Meeting of the Transportation Research Board*
18          *Forthcoming in Transportation Research Record*, The National Academies of Sciences,
19          Engineering, and Medicine Washington, DC, 2016, Vol. 2.
20   24.    Ben-Akiva, M. E. and S. R. Lerman, *Discrete choice analysis: theory and application to*
21          *travel demand*, Vol. 9. MIT press, 1985.
22   25.    Bengio, Y., A. Lodi, and A. Prouvost, Machine learning for combinatorial optimization:
23          a methodological tour d'horizon. *European Journal of Operational Research*, Vol. 290,
24          No. 2, 2021, pp. 405–421.
25   26.    Chen, L., P. Thakuriah, and K. Ampountolas, Short-term prediction of demand for ride-
26          hailing services: A deep learning approach. *Journal of Big Data Analytics in Transporta-*
27          *tion*, Vol. 3, 2021, pp. 175–195.
28   27.    Chu, K.-F., A. Y. Lam, and V. O. Li, Deep multi-scale convolutional LSTM network for
29          travel demand and origin-destination predictions. *IEEE Transactions on Intelligent Trans-*
30          *portation Systems*, Vol. 21, No. 8, 2019, pp. 3219–3232.
31   28.    Ke, J., H. Zheng, H. Yang, and X. M. Chen, Short-term forecasting of passenger demand
32          under on-demand ride services: A spatio-temporal deep learning approach. *Transportation*
33          *research part C: Emerging technologies*, Vol. 85, 2017, pp. 591–608.
34   29.    Wang, D., Y. Yang, and S. Ning, DeepSTCL: A deep spatio-temporal ConvLSTM for travel
35          demand prediction. In *2018 international joint conference on neural networks (IJCNN)*,
36          IEEE, 2018, pp. 1–8.
37   30.    Lei, Z., X. Qian, and S. V. Ukkusuri, Efficient proactive vehicle relocation for on-demand
38          mobility service with recurrent neural networks. *Transportation Research Part C: Emerg-*
39          *ing Technologies*, Vol. 117, 2020, p. 102678.
40   31.    Guo, G. and Y. Xu, A deep reinforcement learning approach to ride-sharing vehicle dis-
41          patching in autonomous mobility-on-demand systems. *IEEE Intelligent Transportation*
42          *Systems Magazine*, Vol. 14, No. 1, 2020, pp. 128–140.
43   32.    Vaswani, A., N. Shazeer, N. Parmar, J. Uszkoreit, L. Jones, A. N. Gomez, Ł. Kaiser,
44          and I. Polosukhin, Attention is all you need. *Advances in neural information processing*
45          *systems*, Vol. 30, 2017.

1    33.   He, K., X. Zhang, S. Ren, and J. Sun, Deep residual learning for image recognition. In
2          *Proceedings of the IEEE conference on computer vision and pattern recognition*, 2016,
3          pp. 770–778.
4    34.   Bengio, Y. et al., Learning deep architectures for AI. *Foundations and trends® in Machine*
5          *Learning*, Vol. 2, No. 1, 2009, pp. 1–127.
6    35.   Engelhardt, R., P. Malcolm, F. Dandl, and K. Bogenberger, Competition and Coopera-
7          tion of Autonomous Ridepooling Services: Game-Based Simulation of a Broker Concept.
8          *Frontiers in Future Transportation*, Vol. 3, 2022, p. 915219.
9    36.   Kingma, D. P. and J. Ba, Adam: A method for stochastic optimization. *arXiv preprint*
10         *arXiv:1412.6980*, 2014.