

MIT Open Access Articles

Robust matching-integrated vehicle rebalancing in ride-hailing system with uncertain demand

The MIT Faculty has made this article openly available. *Please share* how this access benefits you. Your story matters.

Citation: Guo, Xiaotong, Caros, Nicholas S and Zhao, Jinhua. 2021. "Robust matching-integrated vehicle rebalancing in ride-hailing system with uncertain demand." Transportation Research Part B: Methodological, 150.

As Published: 10.1016/J.TRB.2021.05.015

Publisher: Elsevier BV

Persistent URL: https://hdl.handle.net/1721.1/140212

Version: Original manuscript: author's manuscript prior to formal peer review

Terms of use: Creative Commons Attribution-NonCommercial-NoDerivs License



Robust Matching-Integrated Vehicle Rebalancing in Ride-Hailing System with Uncertain Demand

Xiaotong Guo^a, Nicholas S. Caros^a, Jinhua Zhao^{b,*}

^aDepartment of Civil and Environmental Engineering, Massachusetts Institute of Technology, Cambridge, MA 02139, USA

^bDepartment of Urban Studies and Planning, Massachusetts Institute of Technology, Cambridge, MA 02139, USA

Abstract

With the rapid growth of the mobility-on-demand (MoD) market in recent years, ride-hailing companies have become an important element of the urban mobility system. There are two critical components in the operations of ride-hailing companies: driver-customer matching and vehicle rebalancing. In most previous literature, each component is considered separately, and performances of vehicle rebalancing models rely on the accuracy of future demand predictions. To better immunize rebalancing decisions against demand uncertainty, a novel approach, the matching-integrated vehicle rebalancing (MIVR) model, is proposed in this paper to incorporate driver-customer matching into vehicle rebalancing problems to produce better rebalancing strategies. The MIVR model treats the driver-customer matching component at an aggregate level and minimizes a generalized cost including the total vehicle miles traveled (VMT) and the number of unsatisfied requests. For further protection against uncertainty, robust optimization (RO) techniques are introduced to construct a robust version of the MIVR model. Problem-specific uncertainty sets are designed for the robust MIVR model. The proposed MIVR model is tested against two benchmark vehicle rebalancing models using real ride-hailing demand and travel time data from New York City (NYC). The MIVR model is shown to have better performances by reducing customer wait times compared to benchmark models under most scenarios. In addition, the robust MIVR model produces better solutions by planning for demand uncertainty compared to the non-robust (nominal) MIVR model.

Keywords: Ride-hailing, Vehicle Rebalancing, Robust Optimization, Demand Uncertainty.

1. Introduction

Advanced wireless communication and cloud computing technologies coupled with the growing popularity of shared mobility have led to a fast-growing Mobility-on-Demand (MoD) market in recent years [1]. Ride-hailing companies, also known as Transportation Network Companies (TNCs), such as Uber and Lyft have become ubiquitous forms of MoD in most

^{*}Corresponding author

Email addresses: xtguo@mit.edu (Xiaotong Guo), caros@mit.edu (Nicholas S. Caros), jinhua@mit.edu (Jinhua Zhao)

cities over the past decade. The number of worldwide active drivers for Uber grew from 6 almost zero in 2010 to over 3 million in 2017, while Lyft, a relative latecomer to the market. 7 had 1.4 million active drivers in the US and Toronto in 2017 [2]. Two of the primary 8 innovations that allowed them to capture a significant market share from their established 9 competitors, the taxi industry, were: 1) matching trip requests with drivers using a mobile 10 app rather than curbside hailing or an in-advance booking system, and 2) responding to 11 changes in demand by incentivizing or actively dispatching drivers to high-demand areas. 12 These innovations have been identified as two important ride-hailing operations problems in 13 the literature: the driver-customer matching problem and the vehicle rebalancing problem [3]. 14

One of the key technological competence requirements for efficient operation of ridehailing platforms is the algorithmic approaches for optimally matching drivers and customers in real-time [4]. Given a list of available vehicles and trips requested by customers, the matching algorithm pairs drivers and customers according to specific objectives and feasibility constraints. Moreover, matching decisions need to be made quickly, typically within seconds. Researchers have been seeking solutions to improve the operational and computational performance of the on-demand driver-customer matching problem.

Because the spatial distributions of supply and demand in the ride-hailing system are 22 often unbalanced, platforms can improve the operational performance by actively rebalancing 23 idle vehicles to areas where the demand is expected to exceed supply based on estimates of 24 future demand. Algorithms for rebalancing idle vehicles have been proposed for ride-hailing 25 platforms to reduce wait times for customers [5, 6, 7, 8]. However, the performance of vehicle 26 rebalancing algorithms depends on the accurate future demand estimations. Rebalancing 27 decisions generated with inaccurate demand forecasts could have negative impacts on the 28 system performance. Incorporating robustness into the vehicle rebalancing algorithm is 29 one approach to protect solutions against demand uncertainty that arise from inaccurate 30 estimates of future demand [7]. 31

While rebalancing and matching are often treated as separate operations in the literature [3], both problems relate to dispatching idle vehicles, either to pick up customers or to increase supply in areas with high expected demand. A common objective for the drivercustomer matching problem is minimizing the vehicle miles traveled (VMT) and unsatisfied requests [9, 10] while the primary objective for the vehicle rebalancing problem is minimizing the VMT and a functional term measuring the system-wide service availability for future demand [5, 6, 7].

In the vehicle rebalancing problem, the overall goal for improving the system-wide service 39 availability for incoming customers is to minimize the number of unsatisfied requests, which 40 coincides with the objective of the driver-customer matching problem. The functional term in 41 the objective of the vehicle rebalancing problem can therefore be treated as an approximation 42 to represent the number of unsatisfied requests. However, the maximum system-wide service 43 availability for incoming customers does not necessarily lead to the minimum number of 44 unsatisfied requests if there are inaccurate future demand estimates. To immunize vehicle 45 rebalancing decisions against the inherent demand uncertainty, we introduce the driver-46 customer matching component into the vehicle rebalancing problem in order to explicitly 47 model the number of unsatisfied requests. 48

⁴⁹ Nonetheless, there is a methodological difference between driver-customer matching prob ⁵⁰ lems and vehicle rebalancing problems. The driver-customer matching problem is typically

⁵¹ solved by an agent-based model, where each driver and customer are considered individually.
⁵² For the vehicle rebalancing problem, most methods divide the study area into several sub⁵³ regions and the vehicle rebalancing problem is solved at an aggregate level, where vehicles
⁵⁴ are rebalanced between sub-regions.

To resolve this methodological difference, we propose the matching-integrated vehicle rebalancing (MIVR) model where the area partitioning method is retained and the matching component is modeled at an aggregate level. The objective of the MIVR model is to minimize the total VMT and the number of unsatisfied requests. The aggregate matching component of the MIVR model provides a satisfying approximation of the vehicle pickup distance and the number of unsatisfied requests when using small regions.

- Figure 1 provides a toy example to illustrate the benefits of the MIVR model compared to an independent vehicle rebalancing (VR) model, where the service availability is represented by the absolute difference between estimated future demand and supply. Compared to the
- ⁶⁴ independent rebalancing scenario, the matching-integrated rebalancing scenario dispatches
- $_{\rm 65}$ $\,$ the idle vehicle to a location near sub-regions with estimated future demand. This "smart" $\,$
- ⁶⁶ rebalancing decision compensates for inaccurate future demand estimation by harmonizing
- ⁶⁷ vehicle pickup distance across different demand profiles.



Figure 1: Example scenarios comparing regular VR decisions and the MIVR decisions. There are 16 unit squares (sub-regions) and a trip request is equally likely to appear in any of the orange sub-regions in the next time interval. For the independent rebalancing scenario, the rebalancing distance is 2 and the expected pick-up distance is 3 (four possible pick-up distances 0, 3, 4, 5 with $\frac{1}{4}$ probability on each case). For the matching-integrated rebalancing scenario, the rebalancing distance is 2 and the expected pick-up distance is 2.0 (four possible pick-up distances 1, 2, 2, 3 with $\frac{1}{4}$ probability on each case).

To further protect the vehicle rebalancing decisions against demand uncertainty, we introduce robust optimization (RO) techniques to construct a robust MIVR model. Problemspecific uncertainty sets are established to better reflect the uncertainty within ride-hailing demand. In short, the ride-hailing matching process and RO techniques can be incorporated into the rebalancing procedure to produce better vehicle rebalancing decisions for platforms when facing demand uncertainty. The contributions of this paper can be summarized as follows:

- Proposing the MIVR model to incorporate driver-customer matching information to improve vehicle rebalancing problems with explicit modeling of unsatisfied requests for the first time, to the best of authors' knowledge.
- Proposing the robust MIVR model to consider demand uncertainty and designing
 problem-specific uncertainty sets to better reflect the inherent demand uncertainty in
 the ride-hailing system.
- Using simulations to show performance improvements of the MIVR model compared to an independent VR model and a state-of-the-art empty-car routing policy with real demand data and travel times from New York City (NYC). In high supply scenarios, a *Pareto* improvement can be found for the MIVR model when compared to the VR model at aggregate level regarding the overall VMT, the average customer wait time and the number of unsatisfied requests.

 Comparing the nominal MIVR and the robust MIVR under multiple uncertain scenarios by solving a driver-customer matching problem with realized demand and vehicle distributions after rebalancing. The robust MIVR model is shown to perform better under demand uncertainty, especially in conditions of high supply relative to demand.

The remainder of the paper is organized as follows. Section 2 reviews the relevant literature. Section 3 describes the nominal and robust MIVR models and the robust counterpart. Section 4 includes the empirical study design and descriptions for data used in this paper. Benchmark comparisons, scenario testing results and robust solution performances are described in Section 5. Finally, Section 6 recaps the main contributions of this work, outlines the limitations and provides future research directions.

97 2. Existing Literature

98 2.1. Ride-hailing Matching and Rebalancing

Ride-hailing matching is a variant of the classical Dial a Ride Problem, where customer 99 trips are matched with vehicles such that generalized costs are minimized. These costs can 100 include VMT, customer wait time, and penalties for poor service quality. Development of 101 new algorithms for this problem is a very active field of research and the methods have been 102 used by platform operators in practice [11]. Agatz et al. [4] provided a comprehensive survey 103 of literature related to optimization of driver-to-passenger for dynamic ride sharing between 104 travelers with similar itineraries. In a more recent survey, Mourad et al. [12] reviews research 105 related to optimization of shared mobility systems more broadly, which includes ride-hailing. 106 The authors identify demand uncertainty as a critical issue in modeling shared mobility sys-107 tems, and identify stochastic programming and multi-scenario optimization as two possible 108 modeling techniques. Finally, Ho et al. [13] presents an overview of recent research relating to 109

the general Dial a Ride Problem. While this survey is focused on applications such as paratransit and demand-responsive transit, the taxonomy and solution techniques are applicable to ride-hailing problems. Like Mourad et al. [12], the authors find that the development of models and solution methods that include stochastic demand is an important research direction.

Ride-hailing is one type of on-demand service platform, which is characterized by the 115 waiting time sensitivity of customers and service providers without fixed work schedules. 116 Other on-demand service platforms include food and goods delivery services such as Door-117 Dash and Uber Eats, and ride-pooling platforms. Several recent papers have examined the 118 dynamics of on-demand service platforms. It has been shown that customers' sensitivity 119 to delay has a significant impact on optimial pricing and wage setting [14]. Another paper 120 determines optimal prices and wages under different levels of demand, and calibrates param-121 eters using actual ride-hailing data [15]. Cachon et al. [16] develops a model for dynamic 122 pricing in on-demand platforms, demonstrating that such policies benefit stakeholders by 123 expanding access to service during periods of peak demand. Theoretical relationships be-124 tween pricing, demand and detour policies within ride-pooling platforms, which are similar 125 to ride-hailing but with the possibility of shared trips, have also been investigated [17]. 126

Given the size and dynamic nature of the ride-hailing matching problem in large cities. 127 many approaches involve metaheuristic methods to generate sub-optimal solutions [18, 19]. 128 Recently, researchers have investigated the role of matching radii and matching time periods 129 on the optimal solution [20]. Lyu et al. [21] develops an online matching algorithm that 130 considers multiple objectives, and provides a theoretical optimality guarantee for the online 131 solution. Xu et al. [11] proposed a dynamic programming approach to matching that seeks to 132 optimize matching decisions over a long time horizon. Their method, which did not consider 133 demand uncertainty, has been adopted by a leading ride-hailing platform. 134

Optimal rebalancing of idle ride-hailing vehicles has shown to substantially improve sys-135 tem performance. Typical considerations in designing a rebalancing algorithm are the dura-136 tion of the decision period and the costs included in the objective function. Chen and Levin 137 [22] proposed a simple linear programming (LP) model to select vehicle rebalancing flows 138 that minimize travel cost for five minute periods. Zhang et al. [23] showed that a stable 139 predictive control algorithm could be used for dispatching and rebalancing an autonomous 140 ride-hailing fleet in a discrete time system. At each decision period, a mixed-integer linear 141 programming (MILP) is solved to minimize rebalancing travel time. Their method pro-142 duced significant reductions in peak wait times compared to the no rebalancing scenario. 143 Similarly, Iglesias et al. [24] proposed a model predictive control algorithm for operating the 144 ride-hailing system in real-time by leveraging short-term demand forecasts. They utilized 145 the Long Short-Term Memory (LSTM) neural networks to forecast future customer demand 146 for each origin and destination pair and their proposed algorithm outperformed a state-of-147 the-art rebalancing strategy by reducing up to 89.6% of the average customer wait time. 148 Wallar et al. [5] developed an online vehicle rebalancing algorithm that discretized an area 149 into optimal rebalancing sub-regions, resulting in an average wait time reduction of 37%150 compared to the scenario without rebalancing idle vehicles. Braverman et al. [25] formulates 151 a fluid-based optimization model for idle vehicle rebalancing in ride-hailing systems. The 152 authors use a nine-region network and real-life ride-hailing data to show how the fluid-based 153 model results in a higher fraction of passengers served compared to benchmark models. We 154

include the Braverman et al. [25] model as a benchmark to test the results of our own model. 155 Al-Kanj et al. [26] combined the matching and vehicle rebalancing into a single dynamic 156 programming method for autonomous electric vehicles. Their approach employs incentives 157 rather than centralized control to rebalance vehicles, meaning that rebalancing decisions 158 made by the platform are subject to some amount of non-responsiveness by the passenger 159 or vehicle. Dandl et al. [27] also solves for matching and rebalancing decisions in a single 160 optimization model to inform a simulation that tests how demand forecasting accuracy affects 161 the system performance. The authors use an agent-based model, where the objective function 162 is a combination of penalties and rewards for matching and for reducing demand-supply 163 imbalances. Their simulation assumes that all requests are served and customers will wait 164 indefinitely for pickup. In contrast, our method includes matching information and explicitly 165 models customer wait time and unsatisfied requests in order to make rebalancing decisions. 166 In addition to optimization methods, machine learning (ML) approaches have been pro-167 posed to predict demand in rebalancing vehicles [6, 28]. There has also been considerable 168

work on other practical methods, beyond explicit vehicle rebalancing, to achieve greater balance between supply and demand in ride-hailing systems. These methods include dynamic
pricing [16, 29], providing more information to drivers [30], reward schemes [31], alternative
market structures [32] and carpooling incentives [33].

173 2.2. Robust Optimization

RO is a common approach to handle data uncertainty in optimization problems. The 174 general approach is to specify a range for an uncertain parameter (the "uncertainty set"), and 175 optimize over the worst-case realizations within the bounded uncertainty set. The method 176 is therefore well suited to applications where there is considerable uncertainty related to the 177 model input parameters, and when data uncertainties can lead to significant penalties or 178 infeasibility in practice. The solution method for robust optimization problems involves gen-179 erating a deterministic equivalent, called the robust counterpart. Computational tractability 180 of the robust counterpart has been a major practical difficulty [34]. A variety of uncertainty 181 sets have been identified for which the robust counterpart to a robust optimization problem 182 is reasonably tractable [35]. 183

The RO field has grown substantially over the past two decades. Seminal papers in the late 1990s [36, 37] and early 2000s [38] established the field. Comprehensive surveys on the early literature were done by Ben-Tal et al. [34] and Bertsimas et al. [35]. The development of the robust optimization technique has allowed researchers to tackle problems with data uncertainty in a range of fields. Examples can be found for renewable energy network design [39], supply chain operations [40] and health care logistics [41].

190 2.3. Applications of RO in Ride-hailing Operations

In recent years, robust optimization applications in transportation, and ride-hailing rebalancing more specifically, have attracted considerable research attention. Liu et al. [42] considered uncertain local demand in their matching algorithm for ridesharing operations. Miao et al. [43] proposed an RO model for the taxi dispatching problem and tested it using NYC taxi data. They also proposed a data-driven approach to construct the uncertainty set based on historical demand data with a probability guarantee, building on previous data-driven RO theory proposed by Bertsimas et al. [44]. He et al. [45] tackled the robust ride-hailing

rebalancing problem using linear decision rules (LDR) to create a multi-period adaptive RO 198 (ARO) model. Their ARO-based approach is heavily based upon theory developed by Bert-199 simas et al. [46]. To the best of authors' knowledge, no existing papers have incorporated 200 matching component and robust optimization techniques into vehicle rebalancing problems. 201 This research gap is important to address given the prominent role of ride-hailing in urban 202 transportation. Demonstrating how robustness and matching-integrated rebalancing can be 203 combined in ride-hailing operations, and evaluating whether this combination of methods is 204 advantageous, can help to improve future ride-hailing operations. 205

206 3. Methodology

207 3.1. Problem Description

Given an operation period \mathcal{T} , we first divide it into Ω identical time intervals indexed by 208 $k = 1, 2, ..., \Omega$, where the length of each time interval is Δ^1 . Figure 2 displays the framework 209 of the MIVR model. The MIVR model is solved in a rolling-horizon manner, where decision 210 variables are determined repeatedly at the beginning of each time interval. At the beginning 211 of time interval k, κ future time intervals are incorporated in the MIVR model, and only 212 the vehicle rebalancing decisions of the current time interval k are implemented. When 213 proceeding to the next time interval, vehicle locations are observed and updated as the 214 input for the MIVR model. Let $(k, k+1, ..., k+\kappa-1)$ represent time intervals considered 215 at time k, to simplify the notation, these time intervals are indexed by $k = 1, 2, ..., \kappa$. The 216 study region is partitioned into n sub-regions, each sub-region i has an estimated demand 217 $r_i^k \geq 0$ at time k. We define the following two sets: $N = \{1, ..., n\}$ representing the set 218 of sub-regions and $K = \{1, ..., \kappa\}$ representing the set of time intervals considered in the 219 problem. 220



Figure 2: MIVR model framework. Each time interval has length Δ . Grey intervals indicate past time intervals that have been optimized. The green interval represents the current decision time interval and red intervals stand for look-ahead time within the MIVR model.

The MIVR model introduces the driver-customer matching component into the vehicle rebalancing problem by considering interzonal matchings based on estimated demand. Within a time interval k, the vehicle rebalancing phase happens at the beginning of the interval and the driver-customer matching phase is conducted at the end of the interval. In the vehicle rebalancing phase, decision variables are represented by $x_{ij}^k \in \mathbb{N}$ denoting the number of idle

¹The choice of Δ should depend on the size of sub-regions.

vehicles rebalanced from sub-region i to sub-region j at time k. Let $S_i^k \in \mathbb{N}$ indicate the 226 number of available vehicles in sub-region i at time k for the matching phase. Let d_{ij}^k, w_{ij}^k 227 denote the travel distance and time from sub-region i to sub-region j at time k, respectively, 228 which can be approximated by the distance and travel time between the centroids of two 220 sub-regions. We define a parameter $a_{ij}^k \in \{0,1\}$ denoting whether an idle vehicle can be 230 rebalanced from sub-region i to sub-region j at time k, where $a_{ij}^k = 0$ if rebalancing between 231 sub-regions i, j is feasible at time k. The vehicle rebalancing from sub-region i to sub-region 232 j at time k is feasible if $w_{ij}^k \leq \Delta$, stipulating that the vehicle can be rebalanced to the des-233 tination sub-region j within time interval k. Then the feasibility constraint of rebalancing 234 between sub-regions is given by: 235

$$a_{ij}^k \cdot x_{ij}^k = 0 \quad \forall i, j \in N, \ \forall k \in K.$$

$$\tag{1}$$

This constraint does not prevent long-distance rebalancing decisions that occur over several time periods, but rather limits the movement of rebalancing vehicles within a single time period to zones that are reachable within that time period.

In the driver-customer matching phase, matching is considered between sub-regions with-239 out considering actual demand and detailed locations of customers and vehicles. Let $y_{ij}^k \in \mathbb{N}$ 240 denote the number of customers in sub-region i matched with vehicles in sub-region j at 241 time k. It is worth mentioning that decision variables y_{ij}^k of the matching component only 242 serve as auxiliary variables in the MIVR model, which focuses on computing the rebalanc-243 ing decisions. When vehicle are rebalanced and requests are collected, the driver-customer 244 matching problem can then be solved by a separate driver-customer matching problem given 245 the realized demand. Let $T_i^k \in \mathbb{N}$ denote the number of unsatisfied requests in sub-region i 246 at time k. Then constraints related to the matching phase are: 247

$$\sum_{j=1}^{n} y_{ji}^{k} \le S_{i}^{k} \quad \forall i \in N, \ \forall k \in K$$
(2a)

$$\sum_{j=1}^{n} y_{ij}^{k} \le r_{i}^{k} \quad \forall i \in N, \ \forall k \in K$$
(2b)

$$T_i^k = r_i^k - \sum_{j=1}^n y_{ij}^k \quad \forall i \in N, \ \forall k \in K$$
(2c)

Constraints (2a) and (2b) restrict the interzonal matching decisions by the number of 248 available vehicles S_i^k and estimated demand r_i^k . Constraints (2c) define the number of un-249 satisfied requests, which is equivalent to the number of customers who have not been as-250 signed drivers within the current matching phase. When matching customers and drivers, 251 a maximum pickup time constraint is imposed to guarantee that customers do not experi-252 ence excessive wait times. Let \bar{w} denote customers' maximum pickup time and parameter 253 $b_{ij}^k \in \{0,1\}$ denote whether customers in sub-region i can be matched with drivers in sub-254 region j at time k, where $b_{ij}^k = 0$ indicates a feasible interzonal matching. The matching between customers in sub-region i and drivers in sub-region j at time k is feasible if $w_{ji}^k \leq \bar{w}$, 255 256 which enforces the maximum pickup time constraint. The matching feasibility constraint is 257

 $_{258}$ then

$$b_{ij}^k \cdot y_{ij}^k = 0 \quad \forall i, j \in N, \ \forall k \in K.$$
(3)

Next, we establish the connection between the two phases. Let $V_i^k, O_i^k \in \mathbb{N}$ represent the 259 number of vacant and occupied vehicles for sub-region i at the beginning of time interval 260 k, respectively. The initial vehicle locations, $V_i^1, O_i^1, \forall i \in N$, are inputs for the MIVR 261 model. Other inputs to the model are regional transition matrices P^k, Q^k , which describe 262 the dynamics of occupied vehicles. The entry (i, j) for P^k , P^k_{ij} , denotes the probability that 263 an occupied vehicle located in sub-region i at time k will be in sub-region j and stay occupied 264 at time k+1. The entry (i, j) for Q^k , Q^k_{ij} , indicates the probability that an occupied vehicle 265 starting in sub-region i at time k will be in sub-region j and become vacant at time k + 1. 266 In reality, the regional transition matrices depend on the spatio-temporal demand flows 267 as well as the operator's dispatching and rebalancing strategies. The matching and rebalanc-268 ing decisions in the MIVR model are defined at interzonal level, and the regional transition 269 matrices formulated with interzonal level decision variables are approximations to the real 270 matrices. To reduce the model complexity, we further approximate the real regional transi-271

tion matrices with static matrices estimated from the historical data. The impact of utilizing
static transition matrices will be elaborated in the results section. These matrices must
satisfy the following constraints:

$$\sum_{j=1}^{n} (P_{ij}^{k} + Q_{ij}^{k}) = 1, \quad \forall i \in N, \ \forall k \in K.$$

Then, we specify the following relationships between S_i^k, V_i^k, O_i^k and decision variables x_{ij}^k, y_{ij}^k :

$$\sum_{j=1}^{n} x_{ij}^{k} \le V_{i}^{k} \quad \forall i \in N, \ \forall k \in K$$
(4a)

$$S_{i}^{k} = V_{i}^{k} + \sum_{j=1}^{n} x_{ji}^{k} - \sum_{j=1}^{n} x_{ij}^{k} \quad \forall i \in N, \ \forall k \in K$$
(4b)

$$V_{i}^{k+1} = S_{i}^{k} - \sum_{j=1}^{n} y_{ji}^{k} + \sum_{j=1}^{n} Q_{ji}^{k} O_{j}^{k} \quad \forall i \in N, \ \forall k \in K \setminus \{\kappa\}$$
(4c)

$$O_i^{k+1} = \sum_{j=1}^n y_{ji}^k + \sum_{j=1}^n P_{ji}^k O_j^k \quad \forall i \in N, \ \forall k \in K \setminus \{\kappa\}$$

$$(4d)$$

Where constraints (4a) ensure that the number of vehicles in sub-region i that can be rebalanced to other sub-regions is bounded by the number of vacant vehicles. Constraints (4b) show that the available vehicles in sub-region i at time k consist of vacant and rebalanced vehicles. Similarly, constraints (4c) indicate that the set of vacant vehicles in sub-region iat time k + 1 is comprised of currently vacant vehicles at time k and currently occupied vehicles that become vacant in the next time interval. The number of unmatched vehicles at

time k, denoted by $S_i^k - \sum_{j=1}^n y_{ji}^k$, is equal to the difference between the number of available vehicles and the number of vehicles dispatched for interzonal matching. The number of 283 284 occupied vehicles at time k that become vacant at time k+1 in sub-region i is represented 285 by $\sum_{j=1}^{n} Q_{ji}^{k} O_{j}^{k}$. Constraints (4d) state that occupied vehicles in sub-region *i* at time k+1 are comprised of currently vacant vehicles that become occupied in the next interval as 286 287 well as currently occupied vehicles at time k. The number of vacant vehicles that become 288 occupied in sub-region i at time k + 1 because of interzonal matching at time k is indicated 289 by $\sum_{j=1}^{n} y_{ji}^{k}$. The number of occupied vehicles at time k that stay occupied at time k+1 in 290 sub-region *i* is enforced by $\sum_{j=1}^{n} P_{ji}^{k} O_{j}^{k}$. 291

The objective for the MIVR model is minimizing the number of unsatisfied requests 292 and the total vehicle travel distance, which consists of vehicle rebalancing distance and 293 vehicle pickup distance. To construct the objective function as the generalized VMT for 294 ride-hailing operations, we assume γ to be a parameter indicating the penalty VMT induced 295 by each unsatisfied request. Let β represent a parameter that defines the relative weighting 296 of rebalancing distance and pickup distance. The parameter β controls the trade-off between 297 the total non-occupied VMT (from the system perspective) and the service quality (from 298 the customer perspective). A larger β indicates a higher priority on minimizing the vehicle 299 pickup distance, which leads to better service quality with a smaller customer wait time. 300 When $\beta = 1$, the MIVR model purely minimizes the total VMT and the number of unsatisfied 301 requests without explicitly putting any weight on the customer wait times². 302

$$(MIVR) \quad \min \quad Z = \sum_{k=1}^{\kappa} \sum_{i=1}^{n} \sum_{j=1}^{n} x_{ij}^{k} d_{ij}^{k} + \beta \cdot \sum_{k=1}^{\kappa} \sum_{i=1}^{n} \sum_{j=1}^{n} y_{ij}^{k} d_{ji}^{k} + \gamma \cdot \sum_{k=1}^{\kappa} \sum_{i=1}^{n} T_{i}^{k}$$
(5a)

s.t. Constraints
$$(1), (2a) - (2c), (3), (4a) - (4d)$$

$$x_{ij}^k, y_{ij}^k \in \mathbb{N} \quad \forall i, j \in N, \ \forall k \in K$$
(5b)

$$S_i^k, V_i^k, O_i^k, T_i^k \in \mathbb{N} \quad \forall i \in N, \ \forall k \in K$$

$$(5c)$$

The MIVR model is an integer linear programming (ILP) problem with integer variables 303 $x_{ij}^k, y_{ij}^k, S_i^k, V_i^k, O_i^k$ and T_i^k . ILP problems of this size and complexity can be difficult to 304 solve in a reasonable time frame. To improve the computational performance of our model 305 while producing satisfying results, we relax all integer variables in the problem to positive 306 real numbers \mathbb{R}^+ . The rebalancing decisions used for implementations can be generated by 307 rounding down the solutions generated by the relaxed model. The approximated rebalancing 308 decisions are guaranteed to be feasible regarding to constraints (4a), which impose an upper-309 bound on the number of vehicles that can be rebalanced. 310

By incorporating matching decisions within vehicle rebalancing problem, the model also considers future matching distances in addition to the rebalancing distance, leading to "smarter" rebalancing decisions. Essentially, the MIVR reduces the cost of inaccurate demand estimation when rebalancing idle vehicles. Meanwhile, the MIVR model is a forward-

²The MIVR model implicitly weights the customer wait times because of the correlation between the vehicle pickup distance and wait times.

looking model by incorporating κ future time intervals into the model.

316 3.2. Robust Optimization Model Formulation

The estimation of the future demand r_i^k is crucial for vehicle rebalancing problems in 317 ride-hailing systems. Previous studies have assumed the number of customers in any sub-318 region followed a Poisson distribution [5, 6]. However, in most applications we have limited 319 knowledge about the "true" distribution for the future demand. The assumption that com-320 plex customer behaviour can be described by a simple probability distribution might be too 321 strong. Instead of imposing a probability distribution on the future demand, we introduce 322 the robust optimization technique where the uncertain demand parameters are described by 323 uncertainty sets rather than specific probability distributions. The uncertainty sets specify 324 a range for the uncertain demand where the demand can lie anywhere in the range. 325

First, we define the uncertainty set for the robust MIVR model. For the uncertainty in the demand originating in sub-region *i* within time interval *k*, we construct an uncertainty set \mathcal{U} from the intersection of two different sets: a box uncertainty set $\tilde{\mathcal{U}}_i^k$ and a polyhedral uncertainty set $\bar{\mathcal{U}}^k$ which constrains the total variation in demand across all sub-regions. The uncertainty set \mathcal{U} was selected to reflect the actual range of demand variability across different sub-regions without producing solutions that are too conservative in practice.

The box uncertainty set imposes upper and lower bounds of ρ standard deviations between estimated regional demand and the mean regional demand at each time interval k. The parameter ρ is set according to the operator's level of risk tolerance, with a higher ρ representing a lower tolerance for risk. The mean μ_i^k and standard deviation σ_i^k of the demand in sub-region *i* during time *k* are estimated with the historical data. The box uncertainty set for estimated demand r_i^k is then

$$\tilde{\mathcal{U}}_{i}^{k}(\rho) = \left\{ r_{i}^{k} : \left| \frac{r_{i}^{k} - \mu_{i}^{k}}{\sigma_{i}^{k}} \right| \le \rho \right\} \quad \forall i \in N, \ \forall k \in K.$$

The polyhedral uncertainty set limits the total offset in the sum of the demand during a time interval across all sub-regions. This second restriction is intuitive; within a given time interval, demand may be above or below the mean in one region, but the total demand across the entire service area could be expected to remain at a similar level compared to previous days under most scenarios. Sub-regions with unusually high demand should be offset by other nearby sub-regions of low demand. The polyhedral uncertainty set for estimated demand r_i^k is

$$\bar{\mathcal{U}}^k(\Gamma) = \left\{ (r_1^k, ..., r_n^k) : \left| \sum_{i=1}^n (r_i^k - \mu_i^k) \right| \le \Gamma \right\} \quad \forall k \in K,$$

where Γ is the parameter to control the level of uncertainty for the polyhedral uncertainty set. It is worth noting that the construction of the uncertainty set indicates how much uncertainty the operator would like to tolerate in the operation. In reality, there exists scenarios where the total demand at certain time intervals exceed the historical mean by far, for instance ride-hailing demand after concerts or large events. It is wise for the ride-hailing operator to not take such unusual demand scenarios into consideration.

The combined uncertainty set \mathcal{U} for the estimated demand r_i^k is:

$$\mathcal{U} = \left[\bigcap_{i=1}^{n} \bigcap_{k=1}^{\kappa} \tilde{\mathcal{U}}_{i}^{k}(\rho)\right] \cap \left[\bigcap_{k=1}^{\kappa} \mathcal{U}^{k}(\Gamma)\right]$$

By defining an uncertain parameter $\zeta \in \mathbb{R}^{n\kappa}$ and letting $r_i^k = \mu_i^k + \zeta_i^k \sigma_i^k$, we can write \mathcal{U} as follows:

$$\mathcal{U} = \left\{ \boldsymbol{\zeta} : \left\| \boldsymbol{\zeta} \right\|_{\infty} \le \rho; \left\| e^T (\boldsymbol{\zeta}^k \circ \sigma^k) \right\| \le \Gamma, \forall k \in K \right\},$$
(6)

where $\boldsymbol{\zeta}^{k}, \sigma^{k} \in \mathbb{R}^{n}$ are vectors for a specific time interval $k, e \in \mathbb{R}^{n}$ is a vector with all entries equal to one, and $\zeta^{k} \circ \sigma^{k}$ indicates the element-wise product for vectors ζ^{k} and σ^{k} . The parameters ρ and Γ control the size of the uncertainty set for estimated demand, and can be adjusted based on the operators' risk tolerance or desired probability guarantee for constraints involving uncertain parameters. Increasing the value of ρ and Γ leads to more conservative rebalancing decisions for the robust model.

Combining the MIVR model with the uncertainty set described above, we propose a robust MIVR model:

$$(P) \quad \min_{x_{ij}^k, y_{ij}^k} \quad Z = \sum_{k=1}^{\kappa} \sum_{i=1}^n \sum_{j=1}^n x_{ij}^k d_{ij}^k + \beta \cdot \sum_{k=1}^{\kappa} \sum_{i=1}^n \sum_{j=1}^n y_{ij}^k d_{ji}^k + \gamma \cdot \sum_{k=1}^{\kappa} \sum_{i=1}^n T_i^k$$
(7a)

s.t.
$$S_i^k = V_i^k + \sum_{j=1}^n x_{ji}^k - \sum_{j=1}^n x_{ij}^k \quad \forall i \in N, \ \forall k \in K$$
 (7b)

$$V_{i}^{k+1} = S_{i}^{k} - \sum_{j=1}^{n} y_{ji}^{k} + \sum_{j=1}^{n} Q_{ji}^{k} O_{j}^{k} \quad \forall i \in N, \ \forall k \in K \setminus \{\kappa\}$$
(7c)

$$O_i^{k+1} = \sum_{j=1}^n y_{ji}^k + \sum_{j=1}^n P_{ji}^k O_j^k \quad \forall i \in N, \forall k \in K \setminus \{\kappa\}$$
(7d)

$$\sum_{j=1}^{n} x_{ij}^{k} \le V_{i}^{k} \quad \forall i \in N, \ \forall k \in K$$
(7e)

$$\sum_{j=1}^{n} y_{ji}^{k} \le S_{i}^{k} \quad \forall i \in N, \ \forall k \in K$$
(7f)

$$\sum_{j=1}^{n} y_{ij}^{k} \le \mu_{i}^{k} + \zeta_{i}^{k} \sigma_{i}^{k} \quad \forall i \in N, \ \forall k \in K, \ \forall \boldsymbol{\zeta} \in \mathcal{U}$$
(7g)

$$T_i^k = \mu_i^k + \zeta_i^k \sigma_i^k - \sum_{j=1}^n y_{ij}^k \quad \forall i \in N, \ \forall k \in K, \ \forall \boldsymbol{\zeta} \in \mathcal{U}$$
(7h)

$$b_{ij}^k \cdot y_{ij}^k = 0 \quad \forall i \in N, \ \forall k \in K$$
(7i)

- $a_{ij}^k \cdot x_{ij}^k = 0 \quad \forall i \in N, \ \forall k \in K$ (7j)
- $x_{ij}^k, y_{ij}^k \ge 0 \quad \forall i, j \in N, \; \forall k \in K$ (7k)
- $S_i^k, V_i^k, O_i^k, T_i^k \ge 0 \quad \forall i \in N, \ \forall k \in K$ (71)

The problem (P) becomes infeasible even with a small value of ρ if the coefficient of 362 variation³ for uncertain demand is large for some sub-regions during certain time intervals. 363 Particularly, the problem (P) is infeasible if $\exists i \in N, \exists k \in K \text{ and } \rho \geq \frac{\mu_i^k}{\sigma_i^k}$. Because when 364 inequality $\rho \geq \frac{\mu_i^k}{\sigma_i^k}$ holds, the box uncertainty set $\tilde{\mathcal{U}}_i^k(\rho)$ allows ζ_i^k to take values smaller 365 than $-\frac{\mu_i^k}{\sigma_i^k}$, which leads to a negative uncertain demand, i.e., $r_i^k = \mu_i^k + \zeta_i^k \sigma_i^k < 0$. The 366 constraint (7g) is infeasible when the right-hand side is negative since the decision variable 367 y_{ij}^k is non-negative. To prevent infeasibility that can result from demand uncertainty, we add 368 restrictions on the uncertainty set in the problem (P) to guarantee that estimated demand 360 is non-negative: 370

$$\mu_i^k + \zeta_i^k \sigma_i^k \ge 0 \quad \forall i \in N, \ \forall k \in K, \ \forall \boldsymbol{\zeta} \in \mathcal{U}$$
(8)

When modeling robust optimization problems, equality constraints with uncertain parameters should be avoided as much as possible since they dramatically shrink the feasible region and often lead to infeasibility [47]. For the problem (P) with uncertain parameter ζ , we must therefore reformulate equality constraints (7h). Equality constraints (7h) can be avoided by eliminating variable T_i^k through substitution. After this variable elimination step, objective function of problem (7a) becomes:

$$\min_{x_{ij}^k, y_{ij}^k} \left\{ \sum_{k=1}^{\kappa} \sum_{i=1}^n \sum_{j=1}^n x_{ij}^k d_{ij}^k + \beta \cdot \sum_{k=1}^{\kappa} \sum_{i=1}^n \sum_{j=1}^n y_{ij}^k d_{ji}^k + \max_{\zeta \in \mathcal{U}} \left[\gamma \cdot \sum_{k=1}^{\kappa} \sum_{i=1}^n (\mu_i^k + \zeta_i^k \sigma_i^k - \sum_{j=1}^n y_{ij}^k) \right] \right\}. \quad (9)$$

The objective function (9) with min-max formulation can be reformulated by introducing an auxiliary variable ω :

min
$$Z = \omega$$
 (10a)

s.t.
$$\sum_{k=1}^{\kappa} \sum_{i=1}^{n} \sum_{j=1}^{n} x_{ij}^{k} d_{ij}^{k} + \sum_{k=1}^{\kappa} \sum_{i=1}^{n} \sum_{j=1}^{n} (\beta \cdot d_{ji}^{k} - \gamma) y_{ij}^{k} + \gamma \cdot \sum_{k=1}^{\kappa} \sum_{i=1}^{n} (\mu_{i}^{k} + \zeta_{i}^{k} \sigma_{i}^{k}) \le \omega \quad \forall \zeta \in \mathcal{U}$$
(10b)

However, robust counterparts for equivalent formulations of the same problem are not 379 necessarily equivalent [47]. To reformulate the problem while maintaining an identical robust 380 counterpart, we make variables T_i^k adaptive, meaning that both variables are "wait-and-see"⁴ 381 variables relating to uncertain parameters ζ , i.e., $T_i^k = T_i^k(\zeta)$. Introducing adaptive vari-382 ables turns the initial RO problem into an Adaptive Robust Optimization (ARO) problem. 383 A commonly-used approximation method for solving ARO problems is the application of 384 Linear Decision Rules (LDRs), which has been shown to perform well in practice [34, 48]. 385 Also, if the coefficients for the variables to be eliminated in the equality constraint do not in-386 clude uncertain parameters and the constraint is linear in the uncertain parameters, making 387 such variables adaptive and applying LDRs is equivalent to directly eliminating them [47]. 388

³Ratio of the standard deviation to the mean.

⁴The value of "wait-and-see" variables are determined only after the future demand is revealed.

Substitutions with equality constraint (7h) satisfies both conditions, therefore we eliminate variables T_i^k in the problem (P) to ensure no uncertain parameters appear in equality constraints. The reformulation (P') is equivalent to an approximation for the original robust formulation (P) together with restriction (8) on the uncertainty set by applying LDRs:

$$(P') \quad \min \quad Z = \omega \tag{11a}$$
s.t.
$$\sum_{k=1}^{\kappa} \sum_{i=1}^{n} \sum_{j=1}^{n} x_{ij}^{k} d_{ij}^{k} + \sum_{k=1}^{\kappa} \sum_{i=1}^{n} \sum_{j=1}^{n} (\beta \cdot d_{ji}^{k} - \gamma) y_{ij}^{k} + \gamma \cdot \sum_{k=1}^{\kappa} \sum_{i=1}^{n} (\mu_{i}^{k} + \zeta_{i}^{k} \sigma_{i}^{k}) \le \omega \quad \forall \zeta \in \mathcal{U} \tag{11b}}$$

Constraints (7b) - (7g), (7i) - (7l), (8)

After the reformulation, uncertain parameters only appear in the constraints. The next step is to derive the robust counterpart for the robust MIVR model. Constraints (7g), (11b) and Equation (8) with uncertain parameter ζ can be written as the following generic formulation:

$$L(\cdot) + v^T \zeta \le c \quad \forall \zeta \in \mathcal{U}, \tag{12}$$

where $L(\cdot)$ indicates a function of decision variables in problem (P'), v is a vector in dimension $n\kappa$ and c is a scalar. The robust counterpart for the generic constraint (12) is

$$\begin{cases} L(\cdot) + \rho \|\theta_0\|_1 + \Gamma \sum_{k=1}^{\kappa} (\eta_1^k + \eta_2^k) \le c \\ (\eta_1^{k'} - \eta_2^{k'}) \sigma_i^{k'} = \theta_{k'}^{i,k} & \forall i \in N, \ \forall k = k' \in K \\ \theta_{k'}^{i,k} = 0 & \forall i \in N, \ \forall k \neq k' \in K \\ \eta_1^k, \eta_2^k \ge 0 & \forall k \in K \\ \sum_{k=0}^{\kappa} \theta_k = v \end{cases}$$
(13)

Where $\theta_k \in \mathbb{R}^{n\kappa}$ and $\theta_{k'}^{i,k}$ represents (ik)-th entry of vector $\theta_{k'}$, $\forall k' \in K$. The full derivation of the generic robust counterpart of (12) can be found in Appendix A. Then we derive the robust counterpart for problem (P'):

$$(RC) \quad \min \quad Z = \omega$$
(14a)
s.t. Constraints $(7b) - (7f), (7i) - (7l)$
$$\sum_{k=1}^{\kappa} \sum_{i=1}^{n} \sum_{j=1}^{n} x_{ij}^{k} d_{ij}^{k} + \sum_{k=1}^{\kappa} \sum_{i=1}^{n} \sum_{j=1}^{n} (\beta \cdot d_{ji}^{k} - \gamma) y_{ij}^{k} + \gamma \cdot \sum_{k=1}^{\kappa} \sum_{i=1}^{n} \mu_{i}^{k} + \rho \cdot \sum_{k=1}^{\kappa} \sum_{i=1}^{n} \tilde{\theta}_{0}^{i,k}$$
$$+ \Gamma \cdot \sum_{k=1}^{\kappa} (\eta_{1}^{k} + \eta_{2}^{k}) \leq \omega$$
(14b)

$$(\eta_1^{k'} - \eta_2^{k'})\sigma_i^{k'} = \theta_{k'}^{i,k} \quad \forall i \in N, \ \forall k = k' \in K$$

$$(14c)$$

$$\theta_{k'}^{i,k} = 0 \quad \forall i \in N, \ \forall k \neq k' \in K$$
(14d)

$$\sum_{k'=0}^{\kappa} \theta_{k'}^{i,k} = \gamma \cdot \sigma_i^k \quad \forall i \in N, \ \forall k \in K$$
(14e)

$$-\tilde{\theta}_0^{i,k} \le \theta_0^{i,k} \le \tilde{\theta}_0^{i,k} \quad \forall i \in N, \ \forall k \in K$$
(14f)

$$\eta_1^k, \eta_2^k \ge 0 \quad \forall k \in K \tag{14g}$$

$$\sum_{j=1}^{n} y_{ij}^{k} + \rho \sum_{k'=1}^{\kappa} \sum_{i'=1}^{n} (\tau_{1,i,k}^{i',k'} + \tau_{2,i,k}^{i',k'}) + \Gamma \sum_{k'=1}^{\kappa} (\tau_{3,i,k}^{k'} + \tau_{4,i,k}^{k'}) \le \mu_{i}^{k} \quad \forall i \in N, \ \forall k \in K$$

$$(14h)$$

$$\tau_{1,i,k}^{i',k'} - \tau_{2,i,k}^{i',k'} + \sigma_{i'}^{k'}(\tau_{3,i,k}^{k'} - \tau_{4,i,k}^{k'}) = 0 \quad \forall i', i \in N, \ \forall k', k \in K, \ (i',k') \neq (i,k)$$
(14i)

$$\tau_{1,i,k'}^{i',k'} - \tau_{2,i,k}^{i',k'} + \sigma_{i'}^{k'}(\tau_{3,i,k}^{k'} - \tau_{4,i,k}^{k'}) = -\sigma_i^k \quad \forall i' = i \in N, \; \forall k' = k \in K$$
(14j)

$$\tau_{1,i,k}^{i',k'}, \tau_{2,i,k}^{i',k'} \ge 0 \quad \forall i', i \in N, \ \forall k', k \in K$$
(14k)

$$\tau_{3,i,k}^{k'}, \tau_{4,i,k}^{k'} \ge 0 \quad \forall i \in N, \ \forall k', k \in K$$

$$(141)$$

$$\rho \sum_{\substack{k'=1\\i\neq i}}^{\kappa} \sum_{\substack{i'=1\\i\neq i}}^{n} (\nu_{1,i,k}^{i',k'} + \nu_{2,i,k}^{i',k'}) + \Gamma \sum_{\substack{k'=1\\k'=1}}^{\kappa} (\nu_{3,i,k}^{k'} + \nu_{4,i,k}^{k'}) \le \mu_i^k \quad \forall i \in N, \ \forall k \in K$$
(14m)

$$\nu_{1,i,k}^{i',k'} - \nu_{2,i,k}^{i',k'} + \sigma_{i'}^{k'}(\nu_{3,i,k}^{k'} - \nu_{4,i,k}^{k'}) = 0 \quad \forall i', i \in N, \ \forall k', k \in K, \ (i',k') \neq (i,k) \quad (14n)$$

$$\nu_{1,i,k}^{i',k'} - \nu_{2,i,k}^{i',k'} + \sigma_{i'}^{k'}(\nu_{3,i,k}^{k'} - \nu_{4,i,k}^{k'}) = -\sigma_i^k \quad \forall i' = i \in N, \ \forall k' = k \in K$$
(140)

$$\nu_{1,i,k}^{i',\kappa}, \nu_{2,i,k}^{i',\kappa} \ge 0 \quad \forall i', i \in N, \ \forall k', k \in K$$

$$(14p)$$

$$\nu_{3,i,k}^{k'}, \nu_{4,i,k}^{k'} \ge 0 \quad \forall i \in N, \ \forall k', k \in K$$

$$(14q)$$

The constraints (14b) - (14g) represent the robust counterpart of constraints (11b). Constraints (14h) - (14l) are the robust counterpart of constraints (7g) while constraints (14m) - (14q) are the robust counterpart of Equation (8). Compared to problem (P'), the robust counterpart (RC) introduces $(4n^2\kappa^2 + 5n\kappa^2 + 2n\kappa + 2\kappa)$ new auxiliary continuous variables. Although the number of decision variables increases considerably in the robust counterpart, this LP problem can be solved efficiently even for large-scale instances.

408 4. Empirical Study Design

In this section, we describe a real-time ride-hailing simulator used to compare the MIVR model with an independent VR model. To justify the benefit of introducing the robust optimization technique into the vehicle rebalancing problem, a separate matching problem is solved over multiple demand scenarios to evaluate robust solutions and compare the nominal MIVR model with the robust MIVR model. We also describe the data used in the experiments.

415 4.1. Ride-hailing Simulator

The ride-hailing simulator is used to compare the nominal MIVR model with a benchmark VR model described in Appendix B. The results produced by this simulator allow us to evaluate the impact of the MIVR model independent of the robust optimization component. The simulation framework is shown in Figure 3.

⁴²⁰ Data Input. Data input for the ride-hailing simulator including the road network for ⁴²¹ the studied region with a shortest path distance matrix and a predecessor matrix, the set ⁴²² of n sub-regions N, a distance matrix d_{ij}^k and travel time matrix w_{ij}^k between centroids of



Figure 3: Ride-hailing simulation framework.

sub-regions, the set of Ω time intervals with length Δ , a mean μ_i^k of demand for each subregion during each time interval, a full day of ride-hailing demand, and regional probability transition matrices for occupied vehicles P, Q and vacant vehicles P_v, Q_v . Details of the P, Qmatrix estimation methods are provided in Appendix D. Data sources are described in detail in Section 4.3.

Simulation Parameters. Table 1 presents and explains the simulation parameters. Rebalancing decisions are solved with a model considering κ look-ahead time intervals.

Simulation Parameter	Explanation	Base Case Value
α	Cost parameter for regular rebalancing model	10^{2}
eta	Weight parameter for pickup distance	1
γ	Cost parameter for unsatisfied requests	10^{2}
T_{start}	Start time of simulation	00:00
T_{end}	End time of simulation	24:00
Δ	Decision time interval length	$300 \ (seconds)$
δ	Matching batch size	$30 \ (seconds)$
κ	Number of time intervals considered in model	6
$ar{w}$	Maximum pickup time	$300 \ (seconds)$
ilde w	Maximum wait time	$300 \ (seconds)$
N_v	Number of vehicles	3000
$ar{v}$	Average vehicle speed	20 (mph)

Table 1: Simulation parameters and base case value.

430 Demand Generation. Due to privacy concerns, historical TNC trip datasets typically

do not provide exact addresses or coordinates for trip origins and destinations. Given the
demand data at sub-regional level, we randomly assign road nodes within sub-regions as
origins and destinations.

Vehicle Initialization. At the start of the simulation period, the N_v vehicles are equally likely to be in any sub-region *i*. The initial location for a vehicle within a sub-region *i* is randomly assigned to a road node within *i*. All vehicles are considered to be available at the beginning of the simulation.

Simulator. There are two main components contained in the simulator: the vehicle 438 rebalancing engine and driver-customer matching engine. Vehicle locations are updated at 439 the beginning of each simulation iteration. The simulator works as follows: at the beginning 440 of current simulation iteration, vacant and occupied vehicle locations are updated and used as 441 the input for vehicle rebalancing engine; vacant vehicles are rebalanced based on rebalancing 442 decision variables for the current iteration; within each simulation iteration, the driver-443 customer matching engine can be run multiple times depend on the matching batch size 444 (e.g., 30 seconds); vehicles with assigned customers become occupied and start to pick up 445 customers and finish their trips. 446

⁴⁴⁷ Driver-customer Matching Engine. The optimal assignment problem for matching drivers ⁴⁴⁸ with customers in the simulator can be found in Appendix C. The objective of the optimal ⁴⁴⁹ assignment problem is minimizing the number of unsatisfied requests while minimizing the ⁴⁵⁰ pickup distance. The batch size of driver-customer matching engine is δ and customers will ⁴⁵¹ leave the ride-hailing system if they wait longer than the maximum wait time \tilde{w} .

Simulation Results. We evaluated the simulation with the following vehicle-related indicators: number of served customers, non-occupied VMT and number of rebalancing trips.
Customer wait time is used as the customer-related indicator to evaluate the simulation.
The customer wait time includes two components: the time for the vehicle to be assigned to
the customer, and the time for the assigned vehicle to travel to the pickup location.

457 4.2. Robust Solution Evaluation

Evaluating the solutions from the robust model requires multiple different demand scenarios due to the stochastic inputs. We compare the average performance of the model across all demand scenarios in the study period for different uncertainty set sizes.

To evaluate the model performance under each demand scenario, we solve a separate driver-customer matching problem after the demand is realized and the (nominal or robust) rebalancing decision x_{ij}^k (generated with estimated demand) is executed. The drivercustomer matching problem solved here is identical to the one solved in the simulator. The overall pickup time and the number of unsatisfied customers are used as outputs to evaluate robust solutions.

467 4.3. Data Description

The study area used in the experiments is the island of Manhattan in NYC. We used the high-volume ride-hailing trip data collected by the NYC Taxi and Limousine Commission [49] as the demand data. The sub-regions used in the experiments are "taxi zones" defined within the high-volume ride-hailing trip dataset. There are 63 taxi zones on the island of Manhattan (N = 63). For benchmark comparisons of the nominal MIVR model, weekdays in June 2019 were chosen as the analysis period. Only trips that began and ended on the island of Manhattan were included. The mean and standard deviation of daily trip count by zone are shown in Figure 4 to illustrate the overall demand pattern. Demand is generally concentrated around dense residential areas on the eastern and western sides of Manhattan. There was an average of 294,422 high-volume ride-hailing trips per weekday during the sample period.



Figure 4: Average daily demand by zone (trips).

The full day of ride-hailing demand used in the simulation is from June 10, 2019. We chose a non-holiday Wednesday as it represents a typical day of demand from the study period. Figure 5 shows the comparison between the real demand and estimated demand⁵ aggregated into 5-minute time intervals. Based on the relationship between total real demand and total estimated demand, we identify four discrete demand scenarios over which the model can be tested:

I Low demand with accurate estimation (0 - 6): overall demand is relatively low and consistent with the historical average for this period.

- ⁴⁸⁷ II High demand with accurate estimation (6 10): overall demand is high and consistent ⁴⁸⁸ with the historical average for this period.
- III Demand underestimation (11 17): the total demand exceeds the historical average for
 this period.
- ⁴⁹¹ IV Demand overestimation (20 24): the total demand is lower than the historical average ⁴⁹² for this period.

It is worth mentioning that an accurate prediction of the total demand does not lead to accurate sub-regional demand predictions. Demand uncertainties exist in every demand scenario and the overall level of uncertainty is higher in scenarios where demand is underestimated or overestimated. Simulation results for each demand scenario are shown in the

⁵Mean demand μ_i^k is used as estimated demand.



Figure 5: Estimated and real demand with four different types of demand scenarios. The demand is aggregated into 5-minute time intervals.

⁴⁹⁷ next section in order to illustrate the difference in model performance across two dimensions:
⁴⁹⁸ demand level and prediction accuracy.

For evaluations of the robust MIVR model, we utilized the actual demand data for the 499 65 week days from April to June 2019 to reflect the real demand uncertainty. Mean μ_i^k and 500 standard deviation σ_i^k used in the robust MIVR model are generated from the same period. 501 The interzonal travel times for each time interval, w_{ij}^k , were collected from real travel 502 speed data provided by the Uber Movement database for the study period of June 2019 [50]. 503 Hourly link-level travel speed is available for every link with at least five unique trips during 504 the hour. First, the average hourly speed across all days in the study period was determined. 505 The average hourly link travel speed was then used as an input to find the shortest path travel 506 time between each zone pair for each hour in the day. Dijkstra's algorithm [51] was used 507 to determine the shortest path between zone centroids. The regional transition probability 508 matrices for occupied and vacant vehicles, P, Q, P_v and Q_v , are generated based on the real 500 travel time and demand data, and details are shown in Appendix D. 510

511 5. Results

All experiments in this paper are conducted on a 3.0 GHz AMD Threadripper 2970WX Processor with 128 GB Memory. The integer linear program and linear program in the experiments are solved with Gurobi 9.0 [52].

Presentation and discussion of the results is organized into three subsections. Section 5.1 compares the MIVR model to two benchmark models: the VR model described earlier, and a recent state-of-the-art rebalancing model [25]. Section 5.2 explores the sensitivity of the MIVR results to variation in the model inputs. Finally, Section 5.4 provides the results for the robust MIVR model.

520 5.1. Benchmark Comparison

First, we compare the MIVR model with the benchmark VR model described in Appendix B and a fluid-based empty-car routing policy (FERP) proposed by Braverman et al. [25]. The performance of each model is assessed with the ride-hailing simulator described in Section 4. To ensure a fair comparison, each vehicle rebalancing model uses the same demand profile and initial vehicle locations for each scenario.

526 5.1.1. Benchmark VR Model Comparison

The base case scenario (full-day simulation) is tested with the simulation parameters shown in Table 1. The base case considers a scenario with 3000 vehicles, i.e., $N_v = 3000$, and 6 future time intervals in the vehicle rebalancing model, i.e., $\kappa = 6$. The base case scenario purely minimizes the number of unsatisfied requests and the total non-occupied VMT, i.e., $\beta = 1$. Both vehicle- and customer-related metrics are presented in Figure 6, where each figure shows the distributions for vehicles or customers for both MIVR and VR model results.

As shown in Figure 6a, the MIVR model reduces the non-occupied travel distance on 534 average when compared to the VR model. Also, the number of vehicles with extremely 535 long travel distance is reduced when utilizing the MIVR model. Figure 6b displays the 536 rebalancing trip distributions, indicating that the MIVR dispatches fewer vacant vehicles 537 for rebalancing purposes. The distribution of the number of served customers per vehicle 538 is shown in Figure 6c. Although the average number of customers served by each vehicle 539 is identical for two models, vehicles utilization is more evenly distributed under the MIVR 540 model compared to the VR model. Figure 6d compares the wait time between the MIVR and 541



(c) Number of customers served (per vehicle) distribution

(d) Customer wait time distribution

Figure 6: Vehicle- and customer-related metrics in the simulation for the base case.

VR models. The average wait times are 65.5 and 68.5 seconds for each model, respectively. This occurs because the MIVR model reduces the number of customers with longer wait times. The fraction of unsatisfied requests for both models is less than 0.1%. Under the base case scenario, the MIVR model reduces customer wait time by 4.4% on average and total non-occupied VMT by 8.5%.

To better understand the model performance relative to the magnitude of demand and 547 the level of prediction accuracy, we compared the MIVR model with the VR model over the 548 four demand scenarios described in Section 4.3. Figure 7 displays the non-occupied vehicle 549 travel distance distributions and Figure 8 shows the customer wait time distribution over 550 the four demand scenarios. For the low demand with accurate estimation (I) and demand 551 underestimation (III) scenarios, the MIVR model outperforms the VR model by significantly 552 reducing customer wait time while also reducing the average vehicle non-occupied travel 553 distance. In the high demand with accurate estimation scenario (II), the MIVR model 554 reduces customer wait time by proactively rebalancing vehicles more frequently than the VR 555 model. In the demand overestimation scenario (IV), the MIVR model is outperformed by the 556 VR model as the VR model leads to lower average customer wait time and average vehicle 557 non-occupied travel distance. The detailed simulation results for each demand scenario can 558 be found in Appendix E. 559



(c) Demand underestimation (11 - 17)

Distance (mile) (d) Demand overestimation (20 - 24)

Figure 7: Vehicle non-occupied travel distance distributions for different demand scenarios.

⁵⁶⁰ To summarize, the MIVR model dispatches more vacant vehicles than the VR model



(a) Low demand with accurate estimation (0 - 6)

(b) High demand with accurate estimation (6 - 10)



Count

Figure 8: Customer wait time distributions for different demand scenarios.

when the level of estimated demand is high (given a specific fleet size N_v). On the other hand, fewer vehicles are dispatched by the MIVR model compared to the VR model when the level of estimated demand is low. This conclusion is further substantiated in Section 5.2.1, which discusses the results under different fleet sizes. We observe that the MIVR model is less proactive on dispatching vacant vehicles compared to the VR model when the fleet size is large relative to the level of demand.

In this section, we have shown that the performance of rebalancing models, as measured 567 by the average customer wait time, depends on the accuracy of demand prediction and the 568 level of demand. When the error in demand prediction is low, the MIVR model reduces 569 the average customer wait time compared to the VR model. Model performance is penal-570 ized when the error in demand prediction is high (the total demand is underestimated or 571 overestimated). Additionally, a rebalancing model which dispatches more vacant vehicles 572 suffers higher penalties due to inaccurate demand estimation. In the demand scenario III, 573 the level of predicted demand is low and the MIVR model dispatches fewer vacant vehicles 574 than the VR model. Therefore, the MIVR model performs better than the VR model by 575 reacting less often to inaccurate demand estimation. In the demand scenario IV, the level 576 of predicted demand is high and the MIVR model dispatches more vacant vehicles than the 577 VR model. The MIVR model experiences a higher penalty due to inaccurate demand esti-578 mations because of a proactive rebalancing strategy; hence, it performs worse than the VR 579

model under these conditions. The demand scenario IV implies that the demand prediction serves a critical role in the performance MIVR model. These results therefore demonstrate the value of a *robust* MIVR model that explicitly considers demand uncertainty.



583 5.1.2. Benchmark FERP Comparison

Figure 9: Benchmark comparison results between MIVR and FERP models.

To further evaluate the performance of proposed MIVR model, we compared our approach 584 with a state-of-the-art method for solving the vehicle rebalancing problem [25]. Braverman 585 et al. [25] formulated a fluid-based optimization problem to generate a static empty-car 586 routing policy. To guarantee a fair comparison, we chose a two-hour time period (7AM -587 9AM) with historical demand and travel time data from June 2019 and 3000 vehicles to 588 compute a static empty-car routing policy. We implemented the static routing policy in the 589 simulator to dispatch vacant vehicles at each time interval instead of solving an optimization 590 problem. Comparison results are shown in Figure 9. 591

Figure 9a displays the distributions of non-occupied vehicle travel distance and Figure 9b shows the vehicle rebalancing trip distributions. The MIVR model dispatches vacant vehicles more proactively than the FERP. The distributions of number of customers served per vehicle are presented in Figure 9c, where vehicles are utilized slightly more evenly by the MIVR model than the FERP. Figure 9d displays the customer wait time distributions. The MIVR model reduces the average customer wait time by 18% while increasing total nonoccupied VMT by 24%. The proportion of unsatisfied requests for both approaches is less than 0.1%, which is a result of the adequate supply of vehicles. The MIVR model optimizes rebalancing decisions during each time interval and the FERP maintains the same vehicle rebalancing policy throughout the simulation period. In general, the MIVR model provides better service quality for customers by producing a more proactive rebalancing strategy, but it also results in a somewhat higher non-occupied VMT.

604 5.2. Scenario Testing

Second, we test the sensitivity of the results when changing input parameters of the MIVR model, including the fleet size, N_v , the length of decision time interval, κ , the weight parameter for pickup distance, β , and the size of the sub-regions. To avoid the effect of inaccurate demand estimation when testing different scenarios⁶, we tested different scenarios with N_v , κ and β over a four-hour time period (6 AM - 10 AM) assuming perfect future demand predictions. Alternative scenarios are generated by changing the simulation parameters for the base case.

612 5.2.1. Fleet size N_v

Results for scenarios with varying fleet sizes, represented by N_v in the simulation pa-613 rameters, are shown in Figure 10. When there is a limited number of vehicles $(N_v \leq 4000)$ 614 in the system, the MIVR model generates more rebalancing trips per vehicle compared to 615 the VR model. When there are sufficient vehicles in the system ($N_v = 5000$ or 6000), the 616 MIVR model dispatches fewer vacant vehicles and reduces the total non-occupied VMT com-617 pared to the VR model. This is intuitive; for the MIVR model, less rebalancing is needed 618 when there is a higher concentration of idle vehicles since more passengers can be picked up 619 (within the maximum wait time constraint) without significant rebalancing. Therefore, the 620 MIVR model reduces the total non-occupied VMT. The MIVR model decreases the average 621 customer wait time under all scenarios with different fleet sizes compared to the VR model. 622 Customer wait time decreases significantly for the MIVR model when a larger fleet is avail-623 able. Even though rebalancing is not as critical for a large fleet, the MIVR model continues 624 to minimize pickup distance and therefore customer wait time. The proportion of unsatisfied 625 requests is marginally decreased for the MIVR model compared to the VR model, regardless 626 of fleet size. 627

The scenario testing with different fleet sizes implies the existence of the *Pareto* improve-628 ment at aggregate level for the MIVR model compared to the VR model. When a sufficient 629 number of vehicles is available, the MIVR model reduces the total non-occupied VMT, aver-630 age vehicle rebalancing trips and average customer wait time while satisfying more requests 631 compared to the VR model. For instance, when there are 6000 vehicles in the system (with 632 $\kappa = 6$), the MIVR model reduces the total non-occupied VMT by 33%, average vehicle 633 rebalancing trips by 22% and average customer wait time by 36% when compared to the VR 634 model. Under this scenario, the MIVR model clearly outperforms the VR model, indicating 635 that the *Pareto* improvement exists. 636

⁶The effect of input parameters on the simulation results can be overshadowed by the effect induced by inaccurate demand estimations when two have contradictory effects on certain performance metrics.



(a) Average number of rebalancing trips per vehicle and to- (b) Average wait time per customer and proportion of untal non-occupied VMT. Bars indicate the average rebalancing satisfied requests. Bars indicate the average wait time per trip number per vehicle and dashed lines show the total non- customer and dashed lines show the proportion of unsatisfied occupied VMT. requests.

Figure 10: Scenario testing results for different fleet size N_v .

637 5.2.2. Decision time interval length κ



Figure 11: Scenario testing results for different decision time interval length κ .

Figure 11 shows the results under scenarios with varying decision time intervals κ . Both 638 models dispatch more vehicles when considering additional future time intervals (i.e. κ 639 becomes large), and similar amount of vacant vehicles are dispatched by both models. Also, 640 the total non-occupied VMT increases when considering more future time intervals for both 641 models, and the MIVR model leads to less non-occupied VMT compared to the VR model 642 for all scenarios. With respect to customer wait time, considering additional time intervals 643 benefits both models and the MIVR model reduces wait times for all scenarios compared 644 to the VR model. The MIVR outperforms the VR model on the proportion of unsatisfied 645 requests for all scenarios. 646

Note that selecting number of time intervals presents a trade-off between system performance and computation time. Increasing κ linearly increases the size of the problem, which may result in a solution time that is too long to use in practice. The average computation time for solving the MIVR model with $\kappa = 6$ is 3.8 seconds and the average computation time for the MIVR model with $\kappa = 12$ is 7.5 seconds. Platform operators must therefore choose a look-ahead window that is suited to their system size and computational capacity.



Figure 12: Sensitivity testing results for the weight parameter β in the MIVR model. Solid line indicates the average customer wait time and dashed line represents the total non-occupied VMT.

The weight parameter β in the MIVR model controls the trade-off between the total non-occupied VMT and the service quality. In previous experiments, $\beta = 1$ was used as a base case, leading to a MIVR model which purely minimized the total non-occupied VMT and the number of unsatisfied requests. In this section, different values of β are tested based on the base case simulation setting assuming perfect future demand predictions, and the total non-occupied VMT and the average customer wait time are shown in Figure 12.

When β becomes larger, the MIVR model puts more weight on the service quality (cus-660 tomer wait times), and the total non-occupied VMT gets larger. The average customer wait 661 time monotonically decreases when β increases. By increasing the value of β to 3, the average 662 wait time is reduced by 3% while increasing the total non-occupied VMT by 4%. However, 663 the MIVR model becomes more vulnerable to the demand uncertainty when the value of β 664 is large. This is because more vacant vehicles are rebalanced when β is large, where a larger 665 penalty is induced by the inaccurate demand estimations. Therefore, the service quality can 666 be diminished if β is too large. 667

⁶⁶⁸ On the other hand, a negative weight is put on the service quality when $\beta < 1$, meaning ⁶⁶⁹ that the service quality is sacrificed to reduce the total non-occupied VMT. For the scenario ⁶⁷⁰ with $\beta = 0.5$, the total non-occupied VMT is reduced by 0.5% and the average wait time ⁶⁷¹ is increased by 4% compared to the base case. Since the vehicle rebalancing distance is ⁶⁷² highly correlated with customer wait time, reducing β does not significantly decrease the ⁶⁷³ total non-occupied VMT.

674 5.2.4. Sub-regional size

The MIVR model performance relies on the size of sub-regions. Smaller sub-regions leads to more rebalancing options (decision variables) and a better overall model performance. However, the model complexity increases when considering smaller sub-regions. To quantify the effect of changing the size of sub-regions, we combined 63 taxi zones into 13 larger zones and ran simulations for the 13 large sub-regions. Comparison results are shown in Figure 13.

Figure 13a and 13b show the distributions of non-occupied vehicle travel distance and vehicle rebalancing trips. Fewer sub-regions with larger size reduces the opportunities for





Figure 13: Results comparison between simulations with 63 regular sub-regions and 13 large sub-regions.

rebalancing vacant vehicles between sub-regions. Therefore, both the average vehicle non-683 occupied VMT and rebalancing trips are significantly decreased. The distribution of number 684 of customers served per vehicle is shown in Figure 13c, where vehicles are more evenly utilized 685 by the MIVR model under a smaller sub-region size. Figure 13d displays the customer wait 686 time distributions for both scenarios. Compared to the scenario with larger sub-regions, the 687 scenario with 63 sub-regions leads to 20% and 23% reductions on the average customer wait 688 time for the VR and the MIVR, respectively. Differences between the MIVR and VR models 689 hold regardless of the size of the sub-regions. 690

As for the computation complexity, the average running time for producing rebalancing decisions during each iteration by the MIVR model under a regular sub-region size is 3.95 seconds. The average running time for the MIVR model under a larger sub-region size is 0.18 seconds. Reducing the number of sub-regions from 63 to 13 saves approximately 95% of the computation time on generating rebalancing decisions. In general, the size of sub-regions should be chosen to balance computation complexity and model performance.

⁶⁹⁷ 5.3. Impact of Regional Transition Matrices

In the MIVR model, we utilized static regional transition matrices P and Q, which are estimated from the historical data, to reflect the movement of occupied vehicles. However, the true regional transition matrices depend spatio-temporal demand flows and operators' dispatching and rebalancing strategies. In this section, we will quantify the impact of approximating true regional transition matrices with the historical data.

To incorporate the true regional transition matrices in the model, we modified the sim-703 ulator by estimating regional transition matrices for occupied vehicles based on preceding 704 matching decisions at the beginning of each simulation period. By using the previous match-705 ing decisions in the simulation, only regional transition matrices between the current time 706 period k and the next time period k+1 can be evaluated accurately. Therefore, we imple-707 mented a MIVR model with $\kappa = 2$ in the simulation, indicating that two time intervals were 708 considered when making rebalancing decisions. Other simulation parameters are identical to 709 the base case scenario. Such a modified simulator is able to produce rebalancing decisions 710 based on the true regional transition matrices at each time interval. 711

To quantify the impact of approximating regional transition matrices with the historical data, we compared results from the modified simulator to results from a standard simulator described in section 4.1 with $\kappa = 2$, which guarantees identical look-ahead windows in the MIVR model. Results are compared within a four-hour time period (8AM - 12PM) and detailed comparison results are shown in Figure 14.



3000

2500

2000

1000

500 0

50

100

00 1500

(a) Vehicle non-occupied travel distance distribution

(b) Vehicle rebalancing trip distribution

Dynamic: mean = 88.3; stdev = 84.9

Static; mean = 90.5; stdev = 88.0



(c) Number of customers served (per vehicle) distribution



150

200 nds) 250

300

350

Figure 14: Comparison results between simulators with dynamic and static regional transition matrices. *Dynamic* indicates that regional transition matrices are estimated at the beginning of every simulation time interval. *Static* implies that regional transition matrices estimated by the historical data are utilized.

Figure 14a shows the distributions of vehicle non-occupied travel distance. Utilizing true

regional demand matrices reduces the total non-occupied VMT by 1.9%. Distributions of
vehicle rebalancing trips and number of customers served are displayed in Figure 14b and
14c, where two simulators have identical performance on average. Figure 14d presents the
distributions of customer wait time. Using true regional transition matrices reduces the
average customer wait time by by 2.4%.

The comparison results imply that approximating the true regional transition matrices 723 with static matrices estimated from the historical data has a marginal impact on model 724 performance. This is intuitive; the regional transition matrices are used for constructing a 725 forward-looking vehicle rebalancing model. In the simulation, only the rebalancing decisions 726 for the first time interval will be implemented, although rebalancing decisions for κ time 727 periods are generated. When moving to the next time period, real-time information (e.g., 728 vehicle locations) is updated and a separate MIVR model considering κ time intervals is 729 solved. Therefore, regional transition matrices have a limited impact on rebalancing decisions 730 at the first time interval, which subsequently has a marginal impact on model performance. 731 732

733 5.4. Robust Model Results

To evaluate the robust optimization model, we tested multiple scenarios with different 734 levels of uncertainty as defined by the uncertainty set size parameters ρ and Γ . Each robust 735 solution was generated for the robust MIVR model considering 6 future time intervals, i.e., 736 $\kappa = 6$. The model parameters were set as $\beta = 1$, $\gamma = 10^2$ and $\bar{w} = \tilde{w} = 300$. For the 737 number of vehicles N_v , we considered the scenario with 3000 vehicles, indicating a sufficient 738 supply (almost all customers can be served) given the demand profile, and 2000 vehicles, 739 representing an insufficient supply. The initial vehicle distributions V_1 and O_1 are generated 740 using the following process: each vehicle in the fleet with size N_v is either vacant or occupied 741 with equal probability and is randomly assigned to a sub-region. To test the performance 742 for different solutions, we utilized the real demand data from 9 AM - 9:30 AM for 65 work 743 days from April to June 2019 and solved a driver-customer matching problem with realized 744 demand and vehicle distributions after rebalancing. The performance of each solution was 745 evaluated based on the average values of the total pickup time and the number of unsatisfied 746 requests over the 65 demand scenarios. The solution generated by the nominal MIVR model 747 was used as the benchmark for evaluating robust solutions. The performance of each robust 748 solution is displayed as the percentage reduction in performance measurements compared to 749 the nominal solution. 750

For the scenario with insufficient supply $(N_v = 2000 \text{ and a proportion of customers can})$ 751 not be served), Table 2 shows the results about the total pickup time and Table 3 752 displays the percentage reduction for the number of unsatisfied requests over the nominal 753 MIVR model. Introducing uncertainty into the model generates solutions that outperform 754 the nominal solution for all values of ρ and Γ . The uncertain parameter ρ significantly affects 755 the total pickup time and the number of unsatisfied requests while the uncertain parameter 756 Γ has limited impact on them. When a high level of uncertainty is considered in development 757 of the robust MIVR model, more customers can be served with less total pickup time. 758

For the scenario with sufficient number of vehicles ($N_v = 3000$ and almost all customers can be served), the percentage reduction of the total pickup time is shown in Table 4. The robust MIVR model benefits more when having a large fleet of vehicles in the system. The

Γ	0	1	2	3	4	5	6	7	8	9	10
0.0	0.0	0.0	0.0	0.0	0.0	0.0	0.0	0.0	0.0	0.0	0.0
0.1	0.21	0.21	0.21	0.21	0.21	0.21	0.21	0.21	0.21	0.21	0.21
0.2	0.51	1.38	1.18	0.66	1.18	0.66	1.18	1.18	0.66	0.79	1.38
0.3	2.45	2.45	4.52	2.45	2.45	2.45	2.45	4.52	4.52	2.45	4.52
0.4	3.47	4.15	4.15	4.15	4.15	4.15	5.67	4.15	4.15	5.67	5.6
0.5	5.62	5.62	5.62	5.62	5.62	5.62	5.62	5.62	7.68	7.68	5.62
0.6	7.89	7.89	7.89	7.89	7.89	7.89	7.89	7.89	7.89	7.89	7.89
0.7	8.78	10.32	10.32	10.32	10.32	10.32	10.32	10.32	10.32	10.32	10.32
0.8	13.24	13.24	13.24	13.24	13.24	13.24	13.24	13.24	13.24	13.24	13.24
0.9	17.17	17.17	17.17	18.59	17.17	17.17	17.17	19.78	18.59	18.59	17.17
1.0	21.19	19.92	21.23	21.23	21.23	21.23	21.23	21.23	21.23	21.23	21.23

Table 2: Percentage reduction in the total pickup time compared to the nominal MIVR solution with insufficient supply ($N_v = 2000$), for different values of ρ and Γ . Gray cells indicate uncertain scenarios with the largest reduction in pickup time.

Γρ	0	1	2	3	4	5	6	7	8	9	10
0.0	0.0	0.0	0.0	0.0	0.0	0.0	0.0	0.0	0.0	0.0	0.0
0.1	0.2	0.2	0.2	0.2	0.2	0.2	0.2	0.2	0.2	0.2	0.2
0.2	0.2	0.41	0.21	0.2	0.21	0.2	0.21	0.21	0.2	0.2	0.41
0.3	0.17	0.17	0.61	0.17	0.17	0.17	0.17	0.61	0.61	0.17	0.61
0.4	0.14	0.14	0.14	0.14	0.14	0.14	0.22	0.14	0.14	0.22	0.3
0.5	0.08	0.08	0.08	0.08	0.08	0.08	0.08	0.08	0.3	0.3	0.08
0.6	0.15	0.15	0.15	0.15	0.15	0.15	0.15	0.15	0.15	0.15	0.15
0.7	0.2	0.22	0.22	0.22	0.22	0.22	0.22	0.22	0.22	0.22	0.22
0.8	0.46	0.46	0.46	0.46	0.46	0.46	0.46	0.46	0.46	0.46	0.46
0.9	0.56	0.56	0.56	0.56	0.56	0.56	0.56	0.56	0.56	0.56	0.56
1.0	0.54	0.54	0.54	0.54	0.54	0.54	0.54	0.54	0.54	0.54	0.54

Table 3: Percentage reduction in the number of unsatisfied requests compared to the nominal MIVR solution with insufficient supply ($N_v = 2000$). Gray cells indicate uncertain scenarios with the largest reduction in unsatisfied requests.

⁷⁶² largest total pickup time reduction for the robust MIVR model with sufficient supply is ⁷⁶³ 41.03% compared to 21.23% for the scenario with insufficient supply. Under the scenario ⁷⁶⁴ with sufficient supply, all customers can be served and introducing uncertainty into the model ⁷⁶⁵ generates solutions that outperform the nominal solution for all values of ρ and Γ .

The robust MIVR model protects the rebalancing decisions against demand uncertainty by restricting the number of rebalancing trips compared to the nominal MIVR model, which is shown in Figure 15. When dispatching fewer vacant vehicles compared to the nominal case, the penalty incurred due to inaccurate demand estimations is decreased and the system becomes more robust against the demand uncertainty, hence has less total pickup time.

Γ	0	1	2	3	4	5	6	7	8	9	10
0.0	0.0	0.0	0.0	0.0	0.0	0.0	0.0	0.0	0.0	0.0	0.0
0.1	4.2	4.2	4.2	4.2	4.2	4.2	4.2	4.2	4.2	5.23	4.19
0.2	6.4	6.4	6.4	7.44	7.01	7.27	6.67	6.45	6.4	6.4	6.4
0.3	12.51	12.51	12.51	12.0	12.51	12.0	12.0	12.0	12.0	12.51	12.51
0.4	16.23	15.33	15.33	15.33	15.33	15.33	15.33	15.33	15.33	15.33	15.33
0.5	18.19	18.32	18.32	17.57	18.32	18.32	18.32	18.32	18.32	18.32	18.32
0.6	24.14	22.96	22.96	22.98	22.98	22.96	22.98	22.96	22.96	22.96	22.98
0.7	25.62	25.18	25.18	25.18	25.18	25.18	25.18	25.18	25.18	25.18	25.18
0.8	30.89	29.39	29.22	29.13	31.44	29.39	30.98	29.82	29.4	29.73	29.39
0.9	39.82	36.55	38.02	38.92	36.6	36.44	37.16	36.44	38.1	36.44	36.44
1.0	38.22	39.1	39.49	39.01	39.01	40.41	39.49	39.41	41.03	40.47	40.93

Table 4: Percentage reduction in the total pickup time compared to the nominal MIVR solution with sufficient supply $(N_v = 3000)$, for different values of ρ and Γ .

The number of rebalancing trips is significantly restricted (less than 50% compared to the nominal MIVR model) when introducing a high level of uncertainty into the robust MIVR model under the sufficient supply scenario.

Figure 16 shows the daily performance of the robust MIVR model compared to the nominal MIVR model. Under the insufficient supply scenario, even considering a low level of uncertainty ($\rho = 0.1$) can significantly improve the performance of the robust MIVR model (better performance than the nominal MIVR model for 83% of the 65 days tested).





(a) Insufficient supply scenario with fleet size $N_v = 2000$ (the nominal MIVR model conducts 864 vehicle rebalancing trips)

(b) Sufficient supply scenario with fleet size $N_v = 3000$ (the nominal MIVR model conducts 1228 vehicle rebalancing trips)

Figure 15: Rebalancing trips for the robust MIVR model under multiple uncertain scenarios. Each cell represents the percentage of rebalancing trips under a specific level of uncertainty compared to the number rebalancing trips in the nominal MIVR model.



Figure 16: Daily robust MIVR model performance under multiple uncertain scenarios. Each cell represents the percentage of the 65 input days that the robust MIVR model performs strictly better than the nominal MIVR model under a given level of uncertainty.

When incorporating a moderate level of uncertainty ($\rho \ge 0.5$) into the model, the robust MIVR model outperforms the nominal MIVR model for every day of demand tested. When a sufficient supply of vehicles is available, the robust MIVR model performs better than the nominal MIVR model for every weekday tested over most uncertain scenarios.

Overall, the robust MIVR model generates rebalancing decisions based on out-of-sample 782 demand uncertainty defined by parameters ρ and Γ , and solutions are evaluated with real 783 demand data reflecting in-sample demand uncertainty. The parameters ρ and Γ for uncer-784 tainty sets indicate the level of demand uncertainty that ride-hailing operators are willing to 785 protect rebalancing decisions against. Based on experiment results, introducing robustness 786 into the MIVR model and protecting rebalancing decisions against demand uncertainty im-787 prove the system performance effectively under insufficient and sufficient supply cases. The 788 robust MIVR model performs even better when having sufficient number of vehicles in the 789 system. 790

⁷⁹¹ 6. Conclusions and future work

In this paper, we formulate the MIVR model, which incorporates the driver-customer 792 matching component into the consideration of vehicle rebalancing decisions made by ride-793 hailing operators, to protect rebalancing decisions against future demand uncertainty induced 794 by inaccurate demand estimates. We evaluate the performance of our model by comparing 795 against a benchmark VR model and a state-of-the-art model, named fluid-based empty-car 796 routing policy (FERP), using actual ride-hailing trip data. Comparing to the VR model, the 797 MIVR model reduces the average customer wait time and the total non-occupied VMT under 798 most scenarios. When a large fleet is available, a *Pareto* improvement can be found regarding 799 the overall non-occupied VMT, the average vehicle rebalancing trips, the average customer 800 wait time and the number of unsatisfied requests. Comparing to the FERP, the MIVR model 801 reduces the average customer wait time by generating a more proactive rebalancing strategy. 802

To further immunize solutions against demand uncertainty, we propose the robust MIVR model by introducing RO techniques. The robust MIVR is especially effective when the supply of ride-hailing vehicles is sufficient and most requests can be satisfied. Under both sufficient-supply and insufficient-supply cases, the robust MIVR model prevents rebalancing decisions from inaccurate demand estimation by rebalancing fewer vehicles. Additionally, introducing robustness into the MIVR model generates rebalancing decisions that performs better than decisions produced by the nominal MIVR model under most demand scenarios. The main limitations of this study are a result of approximations embedded in the MIVR

810 model. First, we are only able to model trips aggregated to the zonal level given the data 811 availability. While we simulate actual pickups and drop-off locations within those zones. 812 future work could incorporate disaggregate data to test rebalancing and matching at the in-813 dividual address level. The model could be improved if these data were made available. We 814 also assume static regional transition matrices estimated from the historical data. Though 815 having limited impacts on model performances, matching and rebalancing decisions-based re-816 gional transition matrices can be considered in the model to better reflect vehicle trajectories 817 across multiple time periods. 818

This paper shows how internalization of matching costs can be used to protect rebalancing 819 decisions against demand uncertainty and improve the efficiency of ride-hailing operations 820 regarding customers (satisfy more customers with shorter wait times), and under what condi-821 tions the proposed method is beneficial. Furthermore, it illustrates how robust optimization 822 complements the MIVR model by further limiting the risk of increased cost due to incor-823 rect demand estimations. Ride-hailing service operators should consider adopting the robust 824 MIVR model for improved customer outcomes, such as wait time and unsatisfied requests. 825 and reduced costs for operators. 826

There are several future research directions we identified in this paper. First, the un-827 certainty set $\overline{\mathcal{U}}^k(\Gamma)$ has a limited impact on system performance. More effective and in-828 terpretable uncertainty sets could be designed to model the uncertainty in the ride-hailing 829 system. Secondly, additional uncertainty variables could be considered besides the demand 830 uncertainty, such as travel time. Thirdly, we used the historical average as the future de-831 mand estimates in this paper. Advanced demand prediction algorithms can be incorporated 832 within the robust MIVR model to further improve operational performances. Lastly, the 833 MIVR model could be extended to solve the vehicle rebalancing problem in the shared MoD 834 system. 835

836 7. Acknowledgements

The authors would like to thank Qing Yi Wang and Nate Bailey for their comments on the paper.

839 References

- [1] S. Shaheen, A. Cohen, B. Yelchuru, S. Sarkhili, Mobility on demand operational concept report (2017).
- [2] D. Kerr, Lyft grows gangbusters in 2017, bringing competition to Uber, 2018.

- [3] H. Wang, H. Yang, Ridesourcing systems: A framework and review, Transportation
 Research Part B: Methodological 129 (2019) 122–155.
- [4] N. Agatz, A. Erera, M. Savelsbergh, X. Wang, Optimization for dynamic ride-sharing:
 A review, European Journal of Operational Research 223 (2012) 295–303.
- [5] A. Wallar, M. Van Der Zee, J. Alonso-Mora, D. Rus, Vehicle rebalancing for mobility on-demand systems with ride-sharing, in: 2018 IEEE/RSJ International Conference on
 Intelligent Robots and Systems (IROS), IEEE, pp. 4539–4546.
- [6] J. Wen, J. Zhao, P. Jaillet, Rebalancing shared mobility-on-demand systems: A reinforcement learning approach, in: 2017 IEEE 20th International Conference on Intelligent Transportation Systems (ITSC), IEEE, pp. 220–225.
- [7] F. Miao, S. Han, A. M. Hendawi, M. E. Khalefa, J. A. Stankovic, G. J. Pappas, Datadriven distributionally robust vehicle balancing using dynamic region partitions, in: Proceedings of the 8th International Conference on Cyber-Physical Systems - ICCPS
 ⁸⁵⁵ '17, ACM Press, 2017, p. 261–271.
- [8] K. Spieser, S. Samaranayake, W. Gruel, E. Frazzoli, Shared-vehicle mobility-on-demand
 systems: a fleet operator's guide to rebalancing empty vehicles, in: Transportation
 Research Board 95th Annual Meeting, 16-5987, Transportation Research Board.
- [9] R. Baldacci, V. Maniezzo, A. Mingozzi, An exact method for the car pooling problem based on lagrangean column generation, Operations Research 52 (2004) 422–439.
- [10] J. Alonso-Mora, S. Samaranayake, A. Wallar, E. Frazzoli, D. Rus, On-demand high capacity ride-sharing via dynamic trip-vehicle assignment, Proceedings of the National
 Academy of Sciences 114 (2017) 462.
- [11] Z. Xu, Z. Li, Q. Guan, D. Zhang, Q. Li, J. Nan, C. Liu, W. Bian, J. Ye, Large-scale order dispatch in on-demand ride-hailing platforms: A learning and planning approach, in: Proceedings of the 24th ACM SIGKDD International Conference on Knowledge Discovery & Data Mining, pp. 905–913.
- [12] A. Mourad, J. Puchinger, C. Chu, A survey of models and algorithms for optimizing
 shared mobility, Transportation Research Part B: Methodological 123 (2019) 323–346.
- [13] S. C. Ho, W. Y. Szeto, Y.-H. Kuo, J. M. Leung, M. Petering, T. W. Tou, A survey of dial-a-ride problems: Literature review and recent developments, Transportation Research Part B: Methodological 111 (2018) 395–421.
- ⁸⁷⁴ [14] T. A. Taylor, On-demand service platforms, Manufacturing & Service Operations
 ⁸⁷⁵ Management 20 (2018) 704–720.
- [15] J. Bai, K. C. So, C. S. Tang, X. M. Chen, H. Wang, Coordinating supply and demand on an on-demand service platform with impatient customers, Manufacturing & Service Operations Management 21 (2018) 556–570.

- [16] G. P. Cachon, K. M. Daniels, R. Lobel, The role of surge pricing on a service platform
 with self-scheduling capacity, Manufacturing & Service Operations Management 19
 (2017) 368–384.
- ⁸⁸² [17] J. Ke, H. Yang, X. Li, H. Wang, J. Ye, Pricing and equilibrium in on-demand ride-⁸⁸³ pooling markets, Transportation Research Part B: Methodological 139 (2020) 411–431.
- [18] A. A. Syed, I. Gaponova, K. Bogenberger, Neural network-based metaheuristic parameterization with application to the vehicle matching problem in ride-hailing services, Transportation Research Record 2673 (2019) 311–320.
- [19] M. Erdmann, F. Dandl, K. Bogenberger, Dynamic car-passenger matching based on tabu search using global optimization with time windows, in: 2019 8th International Conference on Modeling Simulation and Applied Optimization (ICMSAO), IEEE, pp. 1–5.
- [20] H. Yang, X. Qin, J. Ke, J. Ye, Optimizing matching time interval and matching radius
 in on-demand ride-sourcing markets, Transportation Research Part B: Methodological
 131 (2020) 84–105.
- ⁸⁹⁴ [21] G. Lyu, W. C. Cheung, C.-P. Teo, H. Wang, Multi-objective online ride-matching,
 ⁸⁹⁵ Available at SSRN 3356823 (2019).
- R. Chen, M. W. Levin, Dynamic user equilibrium of mobility-on-demand system with
 linear programming rebalancing strategy, Transportation Research Record 2673 (2019)
 447–459.
- [23] R. Zhang, F. Rossi, M. Pavone, Model predictive control of autonomous mobility-ondemand systems, in: 2016 IEEE International Conference on Robotics and Automation (ICRA), IEEE, pp. 1382–1389.
- [24] R. Iglesias, F. Rossi, K. Wang, D. Hallac, J. Leskovec, M. Pavone, Data-driven model
 predictive control of autonomous mobility-on-demand systems, arXiv:1709.07032 [cs,
 stat] (2017). ArXiv: 1709.07032.
- ⁹⁰⁵ [25] A. Braverman, J. G. Dai, X. Liu, L. Ying, Empty-car routing in ridesharing systems,
 ⁹⁰⁶ Operations Research 67 (2019) 1437–1452.
- ⁹⁰⁷ [26] L. Al-Kanj, J. Nascimento, W. B. Powell, Approximate dynamic programming for
 ⁹⁰⁸ planning a ride-hailing system using autonomous fleets of electric vehicles, European
 ⁹⁰⁹ Journal of Operational Research 284 (2020) 1088–1106.
- ⁹¹⁰ [27] F. Dandl, M. Hyland, K. Bogenberger, H. S. Mahmassani, Evaluating the impact of
 ⁹¹¹ spatio-temporal demand forecast aggregation on the operational performance of shared
 ⁹¹² autonomous mobility fleets, Transportation 46 (2019) 1975–1996.
- [28] M. Guériau, I. Dusparic, SAMoD: Shared autonomous mobility-on-demand using decentralized reinforcement learning, in: 2018 21st International Conference on Intelligent Transportation Systems (ITSC), IEEE, pp. 1558–1563.

- [29] C. Yan, H. Zhu, N. Korolko, D. Woodard, Dynamic pricing and matching in ride-hailing
 platforms, Naval Research Logistics (NRL) 67 (2020) 705–724.
- [30] D. Z. Leon Yang Chu, Zhixi Wan, Harnessing the double-edged sword via routing:
 Information provision on ride-hailing platforms, Available at SSRN 3266250 (2018).
- [31] H. Yang, C. Shao, H. Wang, J. Ye, Integrated reward scheme and surge pricing in
 a ridesourcing market, Transportation Research Part B: Methodological 134 (2020)
 126-142.
- [32] L. Zha, Y. Yin, H. Yang, Economic analysis of ride-sourcing markets, Transportation
 Research Part C: Emerging Technologies 71 (2016) 249–266.
- [33] I. Jindal, Z. T. Qin, X. Chen, M. Nokleby, J. Ye, Optimizing taxi carpool policies
 via reinforcement learning and spatio-temporal mining, in: 2018 IEEE International
 Conference on Big Data (Big Data), IEEE, pp. 1417–1426.
- [34] A. Ben-Tal, L. El Ghaoui, A. Nemirovski, Robust Optimization, Princeton Series in
 Applied Mathematics, Princeton University Press, 2009.
- [35] D. Bertsimas, D. B. Brown, C. Caramanis, Theory and applications of robust optimiza tion, SIAM Review 53 (2011) 464–501.
- [36] A. Ben-Tal, A. Nemirovski, Robust convex optimization, Mathematics of Operations
 Research 23 (1998) 769–805.
- [37] A. Ben-Tal, A. Nemirovski, Robust solutions of uncertain linear programs, Operations
 Research Letters 25 (1999) 1–13.
- ⁹³⁶ [38] D. Bertsimas, M. Sim, The price of robustness, Operations Research 52 (2004) 35–53.
- [39] P. Xiong, P. Jirutitijaroen, C. Singh, A distributionally robust optimization model for
 unit commitment considering uncertain wind power generation, IEEE Transactions on
 Power Systems 32 (2016) 39–49.
- [40] C. Ma, W. Hao, R. He, X. Jia, F. Pan, J. Fan, R. Xiong, Distribution path robust optimization of electric vehicle with multiple distribution centers, PloS One 13 (2018).
- [41] Y. Wang, Y. Zhang, J. Tang, A distributionally robust optimization approach for
 surgery block allocation, European Journal of Operational Research 273 (2019) 740–
 753.
- [42] Y. Liu, W. Skinner, C. Xiang, Globally-optimized realtime supply-demand matching in
 on-demand ridesharing, in: The World Wide Web Conference, pp. 3034–3040.
- F. Miao, S. Han, S. Lin, Q. Wang, J. A. Stankovic, A. Hendawi, D. Zhang, T. He,
 G. J. Pappas, Data-driven robust taxi dispatch under demand uncertainties, IEEE
 Transactions on Control Systems Technology 27 (2017) 175–191.

- ⁹⁵⁰ [44] D. Bertsimas, V. Gupta, N. Kallus, Data-driven robust optimization, Mathematical
 ⁹⁵¹ Programming 167 (2018) 235–292.
- ⁹⁵² [45] L. He, Z. Hu, M. Zhang, Robust repositioning for vehicle sharing, Manufacturing &
 ⁹⁵³ Service Operations Management (2019).
- ⁹⁵⁴ [46] D. Bertsimas, M. Sim, M. Zhang, Adaptive distributionally robust optimization, Man ⁹⁵⁵ agement Science 65 (2019) 604–618.
- ⁹⁵⁶ [47] B. L. Gorissen, I. Yanıkoğlu, D. d. Hertog, A practical guide to robust optimization,
 ⁹⁵⁷ Omega 53 (2015) 124–137. ArXiv: 1501.02634.
- [48] D. Bertsimas, D. den Hertog, Robust and adaptive optimization, Dynamic Ideas LLC,
 Belmont, Massachusetts, 2020.
- [49] NYC Taxi and Limousine Commission, TLC Trip Record Data, 2019. "[Online; accessed
 15-June-2020]".
- ⁹⁶² [50] Uber Technologies, Inc., Uber Movement, New York City travel speeds, 2019. "[Online;
 ⁹⁶³ accessed 15-June-2020]".
- ⁹⁶⁴ [51] E. W. Dijkstra, et al., A note on two problems in connexion with graphs, Numerische
 ⁹⁶⁵ Mathematik 1 (1959) 269–271.
- ⁹⁶⁶ [52] Gurobi Optimization, LLC, Gurobi optimizer reference manual, 2020.

⁹⁶⁷ Appendix A. Derivation of The Robust Counterpart

Given the following generic constraint

$$L(\cdot) + v^T \zeta \le c \quad \forall \zeta \in \mathcal{U}, \tag{A.1}$$

where $L(\cdot)$ indicates a function of decision variables in problem (P'), v is a vector in dimension $n\kappa$ and c is a scalar, it is equivalent to

$$L(\cdot) + \max_{\zeta \in \mathcal{U}} v^T \zeta \le c. \tag{A.2}$$

⁹⁷¹ By taking the convex conjugate of constraint (A.2) we derive the following equivalent ⁹⁷² constraint

$$L(\cdot) + \delta^*(v \mid \mathcal{U}) \le c, \tag{A.3}$$

where $\delta(v \mid \mathcal{U})$ is an indicator function such that $\delta(v \mid \mathcal{U}) = 0$ if $v \in \mathcal{U}$, otherwise $\delta(v \mid \mathcal{U}) = \infty$. $\delta^*(v \mid \mathcal{U})$ is the convex conjugate of $\delta(v \mid \mathcal{U})$. Then we introduce Lemma 1 to help with deriving the robust counterpart [48].

Lemma 1. For a constraint $\bar{a}^T x + \delta^*(P^T x \mid Z) \leq b$, let $Z_1, ..., Z_k$ be closed convex sets, such that $\bigcap_i ri(Z_i) \neq \emptyset^7$, and let $Z = \bigcap_{i=1}^k Z_i$. Then,

$$\delta^*(y \mid Z) = \min_{y^1, \dots, y^k} \{ \sum_{i=1}^k \delta^*(y^i \mid Z_i) \mid \sum_{i=1}^k y^i = y \},\$$

and the constraint becomes

$$\begin{cases} \bar{a}^T x + \sum_{i=1}^k \delta^*(y^i \mid Z_i) \le b\\ \sum_{i=1}^k y^i = P^T x \end{cases}$$

Let $\mathcal{U}_0 = \{\zeta : \|\zeta\|_{\infty} \leq \rho\}$ and $\mathcal{U}_k = \{\zeta : |e^T(\zeta \circ \Sigma^k)| \leq \Gamma\}, \forall k \in K$, where $\Sigma^k \in \mathbb{R}^{n\kappa}$ denotes a vector with (ik)-th entry equals to $\sigma_i^k, \forall i \in N$, and other entries equal to zero. The uncertainty set \mathcal{U} can be written as: $\mathcal{U} = \bigcap_{k=0}^{\kappa} \mathcal{U}_k$. By applying Lemma 1 to constraint (A.3), we develop the following robust counterpart for constraint (A.1):

$$\begin{cases} L(\cdot) + \sum_{k=0}^{\kappa} \delta^*(\theta_k \mid \mathcal{U}_k) \le c \\ \sum_{k=0}^{\kappa} \theta_k = v \end{cases}$$
(A.4)

 $^{^{7}}ri(Z_{i})$ indicates the relative interior of the set Z_{i} .

980 Which is equivalent to

$$\begin{cases} L(\cdot) + \rho \|\theta_0\|_1 + \Gamma \sum_{k=1}^{\kappa} (\eta_1^k + \eta_2^k) \le c \\ (\eta_1^{k'} - \eta_2^{k'}) \sigma_i^{k'} = \theta_{k'}^{i,k} \quad \forall i \in N, \ \forall k = k' \in K \\ \theta_{k'}^{i,k} = 0 \quad \forall i \in N, \ \forall k \neq k' \in K \\ \eta_1^k, \eta_2^k \ge 0 \quad \forall k \in K \\ \sum_{k=0}^{\kappa} \theta_k = v \end{cases}$$
(A.5)

Where $\theta_k \in \mathbb{R}^{n\kappa}$ and $\theta_{k'}^{i,k}$ represents (*ik*)-th entry of vector $\theta_{k'}, \forall k' \in K$.

⁹⁸² Appendix B. Benchmark Vehicle Rebalancing (VR) Model

In this section, we formulate a benchmark vehicle rebalancing (VR) model to test the 983 performance of our MIVR model. With similar notations to the MIVR model, we introduce 984 several additional parameters. Let P_v^k, Q_v^k be regional transition matrices regarding vacant 985 vehicles in time period k, which are learned from the historical data. $P_{v,ij}^k$ stands for the 986 probability for a vacant vehicle in sub-region i at time k to be in sub-region j at time k+1987 and becomes occupied. Similarly, $Q_{v,ij}^k$ denotes the probability for a vacant vehicle in sub-988 region i at time k to be in sub-region j at time k + 1 and remains vacant. Two regional 989 transition matrices satisfy the following condition: 990

$$\sum_{j=1}^{n} (P_{v,ij}^{k} + Q_{v,ij}^{k}) = 1, \quad \forall i \in N, \ \forall k \in K.$$

991

⁹⁹² Then the benchmark VR model is:

$$(VR) \quad \min_{x_{ij}^k} \quad \sum_{k=1}^{\kappa} \sum_{i=1}^n \sum_{j=1}^n x_{ij}^k d_{ij}^k + \alpha \cdot \sum_{k=1}^{\kappa} \sum_{i=1}^n |S_i^k - r_i^k|$$
(B.1a)

s.t.
$$S_i^k = \sum_{j=1}^n x_{ji}^k - \sum_{j=1}^n x_{ij}^k + V_i^k \quad \forall i \in N, \ \forall k \in K$$
 (B.1b)

$$V_i^{k+1} = \sum_{j=1}^n Q_{v,ji}^k S_j^k + \sum_{j=1}^n Q_{ji}^k O_j^k \quad \forall i \in N, \ \forall k \in K \setminus \{\kappa\}$$
(B.1c)

$$O_i^{k+1} = \sum_{j=1}^n P_{v,ji}^k S_j^k + \sum_{j=1}^n P_{ji}^k O_j^k \quad \forall i \in N, \ \forall k \in K \setminus \{\kappa\}$$
(B.1d)

$$\sum_{j=1}^{n} x_{ij}^{k} \le V_{i}^{k} \quad \forall i \in N, \ \forall k \in K$$
(B.1e)

$$a_{ij}^k \cdot x_{ij}^k = 0 \quad \forall i \in N, \ \forall k \in K$$
 (B.1f)

$$x_{ij}^k \in \mathbb{R}^+ \quad \forall i, j \in N, \ \forall k \in K$$
 (B.1g)

$$S_i^k, V_i^k, O_i^k \in \mathbb{R}^+ \quad \forall i \in N, \ \forall k \in K$$
 (B.1h)

Where the objective function (B.1a) consists of vehicle rebalancing cost and a service 993 availability function with a weight parameter α to minimize the difference between available 994 vehicles and estimated demand in each sub-region. Constraints (B.1b) to (B.1d) define the 995 relationship between available vehicles S_i^k , vacant vehicles V_i^k and occupied vehicles O_i^k . The 996 maximum number of available vehicles that can be rebalanced is restricted by constraints 997 (B.1e). Constraints (B.1f) impose the feasibility restrictions for rebalancing decisions, and 998 the non-negativity of integer decision variables are guaranteed by constraints (B.1g) and 999 (B.1h). To increase the computational efficiency while maintaining a satisfying solution, we 1000 further relax integer decision variables x_{ij}^k, S_i^k, V_i^k and O_i^k to positive real numbers \mathbb{R}^+ . 1001

The VR model proposed in this section is sufficient to show the benefit of integrating matching into the VR problem. When having different VR models with the area partitioning assumption, a matching-integrated version can always be constructed.

¹⁰⁰⁵ Appendix C. Optimal Assignment of Drivers to Customers

In this section, the driver-customer assignment problem implemented in the matching 1006 engine of the simulator is described. Within each matching decision time interval δ , let 1007 $\mathcal{R} = \{r_1, ..., r_n\}$ denote a set of waiting customers and $\mathcal{V} = \{v_1, ..., v_m\}$ represent a set of 1008 vacant vehicles in the system. Between a customer r_i and a vehicle v_i , let $\tau(r_i, v_i)$ indicate 1009 the minimum travel time for the vehicle to pick up the customer. The maximum pickup 1010 time for customers is denoted by \bar{w} . First, we construct a bipartite graph G = (V, E), 1011 where $V = \mathcal{R} \cup \mathcal{V}$ and $E = \{e(r_i, v_j) : \forall r_i \in \mathcal{R}, \forall v_j \in \mathcal{V}, \tau(r_i, v_j) \leq \bar{w}\}$, meaning that 1012 an edge exists between a vehicle and a customer if the customer can be picked up by the 1013 vehicle within the maximum pickup time. The cost of each edge $e(r_i, v_j)$ equals to the pickup 1014 time, i.e., $c_{e(r_i,v_j)} = \tau(r_i, v_j)$. The decision variables for the optimal assignment problem are 1015 $x_{e(r_i,v_j)} \in \{0,1\}$ for each edge $e(r_i,v_j) \in E$ in the bipartite graph G, and $y_{r_i} \in \{0,1\}$ for each 1016 customer $r_i \in \mathcal{R}$. $x_{e(r_i,v_j)} = 1$ indicates that the customer r_i will be picked up by the vehicle 1017 v_j in the optimal assignment. $y_{r_i} = 1$ implies that the customer r_i will not be assigned 1018 to any vehicles during the current decision time interval δ . Let $\mathcal{I}(r_i)$ represent the set of 1019 edges connected to a customer vertex r_i in G. Similarly, let $\mathcal{I}(v_i)$ indicate the set of edges 1020 connected to a driver vertex v_i in G. The optimal driver-customer assignment problem is: 1021

$$\min \quad \sum_{e(r_i, v_j) \in E} c_{e(r_i, v_j)} x_{e(r_i, v_j)} + \gamma \cdot \sum_{r_i \in \mathcal{R}} y_{r_i}$$
(C.1a)

s.t.
$$\sum_{e(r_i, v_j) \in \mathcal{I}(v_j)} x_{e(r_i, v_j)} \le 1 \quad \forall v_j \in \mathcal{V}$$
(C.1b)

$$\sum_{e(r_i, v_j) \in \mathcal{I}(r_i)} x_{e(r_i, v_j)} + y_{r_i} = 1 \quad \forall r_i \in \mathcal{R}$$
(C.1c)

$$x_{e(r_i,v_j)} \in \{0,1\} \quad \forall e(r_i,v_j) \in E$$
(C.1d)

$$y_{r_i} \in \{0, 1\} \quad \forall r_i \in \mathcal{R}$$
 (C.1e)

The objective function (C.1a) minimizes the summation of the total pickup time and penalties for unsatisfied requests, where γ stands for the penalty VMT for each unsatisfied customer. Constraints (C.1b) ensure that each vehicle can only be assigned to at most one ¹⁰²⁵ customer. Constraints (C.1c) guarantee that each customer is either served by a vehicle or ¹⁰²⁶ remained waiting during the current matching period. Constraints (C.1d) and (C.1e) make ¹⁰²⁷ sure that the decision variables are binary. The optimal driver-customer assignment problem ¹⁰²⁸ can be solved efficiently by the off-the-shelf ILP solvers (e.g., Gurobi) in the simulation.

¹⁰²⁹ Appendix D. Estimation of Regional Transition Matrix

In this section, the process for estimating the regional transition probability matrices for occupied and vacant vehicles, P, Q, P_v and Q_v , with the real travel time and demand data are described. There are several assumptions we made to generate these matrices:

- Given a travel time and distance between the origin and the destination of a request, the vehicle travels with a constant speed.
- Given the origin and the destination of a request, the vehicle travels along the shortest path with regards to travel time.
- For vacant vehicles within sub-regions, 100% of vehicles remain in the same sub-region.

The detailed procedure is described as follows. First, the list of sub-regions crossed by 1038 the shortest path between each origin and destination pair was determined. The time spent 1039 within each sub-region for each origin-destination pair was weighted by the total demand 1040 to get the average time spent in each sub-region across all trips. For a given starting sub-1041 region, the interzonal shortest paths, sub-region durations and origin-destination demand 1042 patterns were used to determine the likelihood of a given vehicle remaining in the starting 1043 sub-region, transitioning to a nearby sub-region or making a dropoff within a time interval. 1044 These probabilities were then used to populate P and Q. Because the taxi dataset only 1045 contains information about occupied vehicles, assumptions were made for the vacant vehicle 1046 zone transition probability matrices P_v and Q_v . 1047

Appendix E. Benchmark VR Comparison Results for Different Demand Scenarios

In this section, we provide the base case simulation results for four different demand scenarios: low demand with accurate estimation in Figure E.17, high demand with accurate estimation in Figure E.18, demand underestimation in Figure E.19 and demand overestimation in Figure E.20.



Figure E.17: Vehicle- and customer-related metrics in the simulation for the base case under the low demand with accurate estimation scenario (0 - 6).



Figure E.18: Vehicle- and customer-related metrics in the simulation for the base case under the high demand with accurate estimation scenario (6 - 10).



Figure E.19: Vehicle- and customer-related metrics in the simulation for the base case under demand underestimation scenario (11 - 17).



Figure E.20: Vehicle- and customer-related metrics in the simulation for the base case under demand overestimation scenario (20 - 24).