



1
2 **A MEAN-FIELD MARKOV DECISION PROCESS MODEL FOR SPATIAL-TEMPORAL**
3 **SUBSIDIES IN RIDE-SOURCING MARKETS**

4
5 *Zheng Zhu, zhuzheng@ust.hk*

6 *Department of Civil and Environmental Engineering, Hong Kong University of Science and*
7 *Technology, Hong Kong, China*

8 *Jintao Ke*, jke@connect.ust.hk*

9 *Department of Logistics and Maritime Studies, Hong Kong Polytechnic University, Hong*
10 *Kong, China*

11 *Hai Wang, haiwang@cmu.edu*

12 *School of Computing and Information Systems, Singapore Management University, Singapore*

13 *Heinz College of Information Systems and Public Policy, Carnegie Mellon University,*
14 *Pennsylvania, USA*

15 **Corresponding Author*

16 **ABSTRACT**

17 Ride-sourcing services are increasingly popular because of their ability to accommodate on-
18 demand travel needs. A critical issue faced by ride-sourcing platforms is the supply-demand
19 imbalance, as a result of which drivers may spend substantial time on idle cruising and picking
20 up remote passengers. Some platforms attempt to mitigate the imbalance by providing relocation
21 guidance for idle drivers—who may have their own self-relocation strategies and decline to
22 follow the suggestions. Platforms then seek to induce drivers to system-desirable locations by
23 offering them subsidies. This paper proposes a mean-field Markov decision process (MF-MDP)
24 model to depict the dynamics in ride-sourcing markets with mixed agents, whereby the platform
25 aims to optimize some objectives from a system perspective using spatial-temporal subsidies
26 with predefined subsidy rates, and a number of drivers aim to maximize their income by
27 following certain self-relocation strategies. To solve the model more efficiently, we further
28 develop a representative-agent reinforcement learning algorithm that uses a representative driver
29 to model the decision-making process of multiple drivers. This approach is shown to achieve
30 significant computational advantages, faster convergence, and better performance. Using case
31 studies, we demonstrate that by providing some spatial-temporal subsidies, the platform is well
32 able to balance a short-term objective of maximizing immediate revenue and a long-term
33 objective of maximizing service rate, while drivers can earn higher income.

34 **Keywords:** ride-sourcing, subsidy, mean-field, Markov decision process, mixed agents

1. BACKGROUND

The emergence of advanced information technologies and the surge in smartphone users enable the fast development of ride-sourcing services. Provided by transportation network companies (TNCs), such as Uber, Lyft, DiDi, and Grab, ride-sourcing services address individuals' on-demand travel needs. A ride-sourcing market is analogous to a more efficient dial-hailing taxi market, in which passengers request services with a few clicks in smartphone apps. Unlike traditional taxi services with large meeting frictions due to street-hailing behaviors between drivers and passengers, ride-sourcing services enable passengers to be matched with drivers at a certain distance. Upon receiving a travel request from a passenger, the platform assigns the passenger to a near driver who then picks up and delivers the passenger. On one hand, the efficiency of supply-demand matching makes ride-sourcing systems indispensable in modern transportation systems. On the other hand, drivers may spend significant amounts of time on idle cruising¹ (IC; i.e., waiting for dispatches) and on the way to pick up passengers. A market failure, called a "wild-goose chase" (WGC), even occurs when drivers are always dispatched to far-away passengers, and waste substantial time on picking them up. These lead to low effective earning rates of drivers and cause negative social externalities, such as exacerbating traffic congestion and increasing carbon dioxide emissions.

The main cause of inefficient IC and the WGC phenomenon is the spatial-temporal supply-demand imbalance. If there is a lack of idle drivers in one region, the platform must call remote drivers to enter the region to mitigate the loss of passengers and revenue. However, these drivers may suffer from long pick-up time (i.e., as in WGC). By contrast, if there are insufficient passengers in one region, drivers may suffer from long idle time (i.e., as in IC). To tackle the issue of supply-demand imbalance, a number of approaches have been proposed, including but not limited to order dispatching (Xu et al., 2018; Yang et al., 2020a) and surge pricing (Zha et al., 2018). In particular, with the fast development of computational power and artificial intelligence technologies, researchers are paying increasing attention to the design and optimization of idle-vehicle relocation strategies for improving supply-demand balance (Rong et al., 2016; Yu et al., 2019; Lin et al., 2018).

In practice, based on actual or predicted spatial-temporal information on supply/demand (Ke et al., 2019) and traffic conditions (Zhu et al., 2019a), idle drivers are advised/incentivized to cruise to regions with higher potential rewards. These rewards could be reflected by the saving on waiting/matching time (Hwang et al., 2015); increase in trip fares and income (Rong et al., 2016; Shou et al., 2020a); increase in vehicle occupancy/utilization rate (Gao et al., 2018); and saving on idle-cruise distance and operational costs (Lin et al., 2018; Yu et al., 2019). These studies aim to generate optimal sequential movements for idle drivers to achieve some maximal system-wide total rewards over a time horizon. Dynamic gaming approaches, such as the Markov decision process (MDP), offer a convenient framework for formulating and solving these problems. In an MDP model, one or multiple players (also referred to as agents) interact with an environment. Each agent has a set of states and a set of actions. In each time slot, each agent chooses one action after it perceives the current state. Meanwhile, by taking an action, the agent receives a reward and their state will be updated by the current state, action, and the state transition law, moving to the next state. During a time horizon, agents attempt to seek out the optimal sequence of actions (determined by a policy that maps the current state to the action) that leads to maximal total rewards. In particular, an MDP model with multiple agents is referred to as a multi-agent MDP model.

¹ We use "idle cruising" because in ride-sourcing markets some vacant vehicles are en route to pick up passengers. To distinguish this from traditional taxi markets, we note that these vehicles are vacant but not idle.

1 Although MDP-based approaches for idle-vehicle relocation have been established in
2 recent studies, research gaps remain. For instance, none of the previous studies have examined
3 the designs and analysis of spatial-temporal subsidies for ride-sourcing drivers with their own
4 relocation strategies using an MDP framework. To be more specific, each driver aims to
5 maximize their own earning by relocating to profitable regions, while the platform tries to
6 incentivize drivers' relocating behaviors to maximize overall system efficiency by subsidies.
7 Clearly, the sequential decision-making of drivers and the platform interact with each other,
8 resulting in a very complex multi-agent MDP with different (i.e., mixed) types of agents. To well
9 formulate such a complex system, we propose a *mean-field (MF)-MDP model*, which can jointly
10 analyze the platform's spatial-temporal subsidies and idle drivers' self-relocation strategies. We
11 regard the platform as a major agent that pursues the subsidy to optimize some objectives from
12 a system perspective—e.g., to maximize immediate revenue and/or the number of passengers
13 served (service rate). A number of drivers are considered as minor agents who choose their self-
14 relocation strategies to maximize their income. The decisions of the platform directly affect the
15 income and decisions of the drivers, while the decisions of the drivers, in turn, collectively affect
16 the platform's decisions via their average status (e.g., the spatial-temporal distribution of idle
17 drivers), which is captured by the MF state. By using a simple stochastic process to approximate
18 the MF state (instead of computing it based on each driver's state), we are able to reduce the
19 standard multi-agent MF-MDP model to a simplified MF-MDP model with only the platform
20 and one representative driver as agents. We then develop a representative-agent reinforcement
21 learning algorithm to solve the simplified model. We conduct a set of numerical studies to
22 examine the performance of the proposed representative-agent algorithm. By performing
23 sensitivity analysis, we further investigate the impacts of spatial-temporal subsidies on drivers'
24 self-relocation, drivers' income, the number of passengers served, and the platform's net
25 revenue. The results suggest that by providing some spatial-temporal subsidies, the platform is
26 able to achieve a higher total reward, while drivers can earn higher income.
27

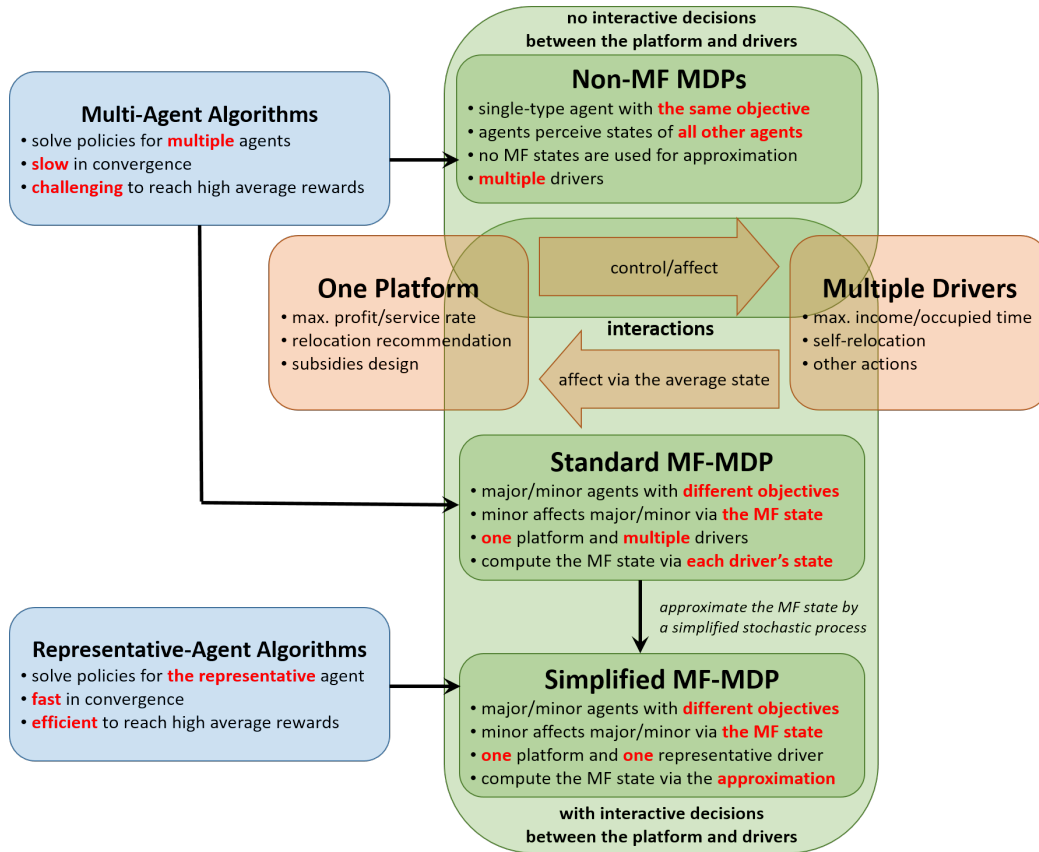


Figure 1. Summary of this paper's contributions

We use the term *non-MF-MDPs* to denote MDPs in which the MF state of minor agents is not required to compute the dynamics (e.g., transition laws and states of agents) in the environment. The main distinctions between the proposed MF-MDP model and other non-MF-MDP models, and the features of their targeting research problems are summarized in Figure 1.

In a non-MF-MDP model, each agent makes decisions by perceiving the states of all other agents, which may render the algorithm hard to converge due to the high stochasticity and instability of the environment. In an MF-MDP model, the states of agents are averaged in each zone and each time interval, and each agent makes its decisions according to the averaged (mean-field) state. This will help reduce the variance of the states and actions, and thus make the model easy to be trained. In describing a ride-sourcing system with one platform and multiple drivers, both non-MF-MDP models and standard MF-MDP models contain multiple agents (the platform and drivers). Naturally, these two models can be solved by multi-agent algorithms that treat each driver and the platform as an independent agent. The only difference is that agents in MF-MDP models could take the MF states as inputs for making actions, while non-MF-MDP models should be aware of the states of all other agents at each decision point.

Additionally, in a complex system with a large number of drivers as agents, multi-agent algorithms need to identify the optimal policy for each specific agent, and the underlying solution space (i.e., the Cartesian product of each agent's state-action set) could be so large that optimal strategies are hard to be identify. To address this critical issue, we then propose a simplified MF-MDP model that uses a representative driver to make decisions for all independent drivers. In other words, the simplified MF-MDP model only identify optimal policies for the representative driver (minor player) and the platform (major player), resulting in a much smaller solution space.

1 By developing representative-agent solution algorithms, it could be much easier for simplified
2 MF-MDP models to fast convergence to high-rewarding policies.

3 The main contributions of this paper are:

- 4 • We propose a generalized MF-MDP model to capture the interactive decisions between
5 the platform and a group of drivers with different objectives in ride-sourcing markets; in
6 contrast, previous studies in this domain generally assume that the platform has full
7 control of ride-sourcing drivers or that the platform and drivers have the same objective.
- 8 • We show theoretically that this generalized multi-agent MF-MDP model (also referred
9 to as the standard MF-MDP) can be approximated by a simplified MF-MDP model that
10 attempts to jointly identify the optimal policies of a platform and a representative driver.
11 The simplified MF-MDP model offers computational advantages for solving multi-agent
12 MDP models.
- 13 • We formulate a specific MF-MDP model to design spatial-temporal subsidies with
14 predefined subsidy rates for drivers with self-relocation strategies. A representative-agent
15 reinforcement learning algorithm is developed to solve the simplified MF-MDP model.
16 Numerical studies demonstrate the effectiveness of the proposed algorithms in a small
17 market, and examine the influences of spatial-temporal subsidies on a few key measures.

18 The rest of the paper is organized as follows. Section 2 reviews the literature on ride-
19 sourcing markets and, in particular, idle-vehicle relocation. Section 3 presents the generalized
20 ride-sourcing MF-MDP model with the platform and drivers as mixed agents. We discuss the
21 approximation of the MF state and the simplified MF-MDP model with theoretical properties
22 and a dynamic programming approach. In Section 4, we adopt the proposed MF-MDP model
23 and formulate the spatial-temporal subsidy problem. A representative-agent reinforcement
24 learning algorithm is developed. We conduct a set of numerical studies in Section 5 and
25 demonstrate the advantage of the representative-agent algorithm over conventional multi-agent
26 algorithms and the potential impacts of the subsidies on the platform and drivers. In Section 6,
27 we discuss different potential subsidy schemes. Section 7 concludes.

28 29 **2. LITERATURE REVIEW**

30 With the development and deployment of smartphone and information technologies, ride-
31 sourcing services have had substantial impacts on traditional taxis in terms of passengers' mode
32 choices and mobility efficiency, and therefore have received intensive attention from researchers
33 across fields. General research problems include optimal operating strategy designs in terms of
34 the trip fares charged to passengers and wages paid to drivers (Cachon et al., 2015; Castillo et
35 al., 2017; Zha et al., 2016; Bai et al., 2019; Taylor, 2018; Yang et al. 2020b); implications of
36 governmental policies and regulations (Yu et al., 2019); examination of the elasticities of labor
37 supply with respect to driver income level (Sun et al., 2019a; Sun et al., 2019b); on-demand
38 matching and dispatching strategies (Xu et al., 2017; Zha et al., 2018; Zhang et al., 2017; Lyu et
39 al. 2019; Yang et al. 2020a); forecasting real-time demand and supply (Ke et al., 2017; Ke et al.,
40 2021; Yao et al., 2018; Zhu et al., 2021b); equilibrium in ride-pooling services (Ke et al., 2020);
41 and the impact of ride-sourcing on public transit (Zhu et al., 2020). Readers may refer to Wang
42 and Yang (2019) for a comprehensive review.

43 One critical problem faced by ride-sourcing platforms is how to mitigate supply-demand
44 imbalance over space and time, which is commonly observed due to the stochastic arrivals and
45 heterogeneous distributions of both drivers and passengers. The supply-demand imbalance can
46 be alleviated with the help of approaches such as spatial-temporal demand prediction (Ke et al.,
47 2021); fleet-size regulation (Yang et al., 2002; Lin et al., 2018); surge/spatial pricing and rewards

1 (Yang et al., 2010; Zha et al., 2018; Zuniga Garcia, 2019; Yang et al. 2020b); driver
2 incentive/subsidy (Qian et al., 2017); efficient large-scale order dispatch (Xu et al., 2018; Li et
3 al., 2019); and idle-vehicle relocation guidance (Rong et al., 2016; Yu et al., 2019).

4 Of these methods, idle-vehicle relocation, which guides or incentivizes idle vehicles from
5 regions with extra supply to locate to regions with inadequate supply, is attracting substantial
6 attention. Braverman et al. (2019) propose a fluid-based optimization approach that controls the
7 flow of empty vehicles to optimize system-wide network utility, measured by the availability of
8 idle vehicles upon passenger arrivals. They show that the optimal utility obtained from a fluid-
9 based approach is an upper bound on the utility of a system with finite vehicles for any routing
10 policy. Lin et al. (2018) propose a multi-agent deep reinforcement learning approach that controls
11 the movements of idle vehicles. Using mobility data from DiDi, they show that their proposed
12 multi-agent model significantly outperforms benchmark algorithms. In this study, the multi-
13 agent advanced actor-critic (A2C) algorithm shows its ability to solve large-scale multi-agent
14 reinforcement learning problems based on a simulator calibrated by actual mobility data. Most
15 studies on idle-vehicle relocation assume that the platform has full control of drivers/vehicles
16 (e.g., Rong et al., 2016; Lin et al., 2018; Shou et al., 2020b). In reality, however, the ride-sourcing
17 platform and drivers have different objectives: Drivers aim to maximize their individual rewards
18 (measured by income, vehicle occupancy rate, etc.) by following certain self-relocation
19 strategies, while the platform aims to maximize overall system performances (measured by net
20 revenue, saving on matching times, number of passengers served, etc.) by using spatial-temporal
21 subsidy/guidance strategies. In this manner, incentives (e.g., subsidies or other rewards to drivers)
22 are critical to motivate drivers to move from demand-cool locations (with more supply than
23 demand) to demand-hot locations (with more demand than supply). Although
24 subsidies/incentives strategies have been implemented in some ride-sourcing companies, such
25 as DiDi, they have not been fully examined in the literature, particularly in a MDP framework.
26 Shou and Di (2020a) propose a multi-agent reinforcement learning paradigm to approximate the
27 system’s equilibrating process in a routing game among atomic selfish agents on a network. Sous
28 and Di (2020b) examine reward design scenarios with multiple drivers and a constant design
29 across zones and time periods, in which the Bayesian optimization is adopted to find the optimal
30 design strategy. Their model can help policymakers develop optimal operational and planning
31 countermeasures under different environments. The two studies also consider mean-field
32 approximation within the reinforcement learning algorithm only; in contrast, our model is a
33 mean-field “oriented” that builds the ride-sourcing simulation based on mean-field information.

34 From a modeling perspective, the difficulty in mitigating the supply-demand imbalance
35 in ride-sourcing markets lies in the complicated dynamic decision processes of the platform,
36 drivers, and passengers, as well as endogenous relationships between decisions and scenarios.
37 Specifically, the platform’s strategies, such as order dispatching, idle-vehicle relocation, and
38 dynamic pricing/subsidies, affect both supply and demand, which in turn affect the platform’s
39 decisions. A promising option for capturing the dynamics of ride-sourcing markets is the family
40 of MDP models, which can well describe the sequential interactions between agents and
41 environment. For example, Xu et al. (2018) formulate an order-dispatching process for a ride-
42 sourcing system using an MDP model, with the order dispatch as action, the numbers of idle
43 drivers and waiting passengers in each time/location as states, and the total gross merchandise
44 volume (GMV) as reward. They propose a policy that simultaneously considers the immediate
45 reward and long-term rewards, and demonstrate that the proposed policy based on the MDP
46 model can substantially improve the per-day earnings of drivers. More recently, various MDP
47 and reinforcement learning models (e.g., Wang et al., 2018; Li et al., 2019; Shou et al., 2020a;
48 Jin et al., 2019) have been developed to enhance the supply-demand balance via better
49 dispatching and idle-vehicle relocation strategies. However, as stated in Section 1, in a

1 complicated system with a huge number of drivers, it can be difficult to identify optimal policies
2 for each specific driver. In addition, in most previous studies, drivers and the platform’s
3 objectives are not necessarily coincident with each other. While these studies assume the
4 platform’s reward is equal to the summation of the rewards of all drivers (which implies that the
5 platform and drivers have the same objective), it is more interesting and realistic to ascertain the
6 platform’s and drivers’ own policies in an environment where they mutually affect each other.
7 To be more specific, drivers try to maximize their individual daily earning through self-relocation,
8 while the platform attempts to maximize system-wide efficiency by paying subsidies to drivers
9 based on location and time.

10 Inspired by the aforementioned studies and research gaps, we propose a generalized MF-
11 MDP model to analyze the dynamics in ride-sourcing markets in which the platform and multiple
12 drivers have different objectives and state-action sets. The MF-MDP model is new to
13 transportation research questions which are solved by an MDP environment with interactive
14 decisions between the mixed agents. According to the proposed model, we theoretically show
15 that efficient algorithms can be developed by only considering the platform and a representative
16 driver as agents in a simplified MF-MDP model. A specific MF-MDP model and a
17 representative-agent reinforcement learning algorithm are developed to analyze the implications
18 of spatial-temporal subsidies for drivers with self-relocation strategies. Our numerical results
19 offer insights on the interactions between the platform’s subsidy and idle drivers’ self-relocation,
20 as well as the influences of the intensity of subsidy on the platform’s spatial-temporal subsidy
21 strategy and idle drivers’ self-relocation strategies.

22 23 **3. MEAN-FIELD MARKOV DECISION PROCESS MODEL FOR RIDE-SOURCING** 24 **MARKETS**

25 In this section, we present a generalized MF-MDP model for depicting the interactive
26 decision processes of the platform and drivers in ride-sourcing markets. With generalized
27 definitions and formulas for states, actions, the MF state, state transition laws, and rewards for
28 mixed agents, we present some properties of the MF-MDP model. We also discuss simplification
29 of the model to reduce the number of agents for computational advantages.

30 **3.1 GENERAL CONCEPT OF THE MF-MDP MODEL**

31 The development of an MDP model should capture the particular feature of a research
32 problem, which is depicted by the definition of states, actions, rewards of agents and the
33 transition law (i.e., how the environment replies to agents’ actions). In a practical problem, the
34 number of states and actions for an agent can be large. For instance, in idle-vehicle relocation
35 problems, a driver’s state should include time and location and his/her actions may cover a list
36 of locations/directions. Moreover, the transition law could involve complex computations that is
37 executed based on spatial-temporal information of each agent and extra information of the
38 environment. Given the large sets of states and actions and the complex transition for each agent,
39 solving a multi-agent MDP model with a large number of agents results in a massive solution
40 space and thus can be computationally prohibitive. The scenario can become more complicated
41 when different types of agents (who may have distinct objectives) coexist in the environment,
42 resulting in an MDP with mixed agents. To capture the interactions between a major agent and
43 a number of minor agents who pursue their individual objectives, [Huang et al. \(2006\)](#) propose
44 the concept of an MF-MDP model. In an MF-MDP model, the states and actions of the major
45 agent can significantly affect the rewards and actions of minor agents. Meanwhile, each minor
46 agent has a negligible impact on the rewards and actions of another minor agent or the major
47 agent. Instead, the transitions, rewards, and actions of the major agent and a minor agent are

1 influenced by the mean-field (i.e., MF, average) state of all minor agents collectively (Gomes,
2 2014). In this manner, the major agent or a specific minor agent does not distinguish any
3 individual minor agent in the MF-MDP model, but considers the MF state when taking actions.

4 In a standard MF-MDP (with one major agent and multiple minor agents), we need to
5 compute the MF state via summarizing each minor agents' state so as to obtain the transitions
6 and rewards. This can be computationally intractable when the number of minor agents is huge.
7 To improve the efficiency, the standard MF-MDP can be simplified by approximating the MF
8 state in a stochastic process and using a representative agent to determine actions for multiple
9 minor agents with the same objective, states, and action sets (Huang et al., 2006; Huang et al.,
10 2007). Once there are a large number of minor agents in the environment, the simplified MF-
11 MDP can well approximate the dynamic nature of the standard MF-MDP. Also, it can
12 significantly reduce computational complexity and achieve more efficient solution by optimizing
13 only one policy for the representative minor agent instead of determining a group of independent
14 policies for each of the minor agents. Correspondingly, we propose representative-agent
15 dynamic programming/reinforcement learning algorithms to solve simplified MF-MDPs (see the
16 next section), while conventional DP/RL algorithms are adopted to solve standard MF-MDPs
17 and non-MF-MDPs.

18 Literature on the MF-MDP model (e.g., Huang et al., 2006; Huang et al., 2007) mainly
19 focuses on the general conception, definitions, and mathematical propositions in a simple and
20 stylized case; there is no discussion of how to configure and solve such a model when the
21 environment is complicated. Inspired by the concept of the MF-MDP model, this paper aims to
22 develop a MDP model that can well delineate the state-action transition laws in a system with
23 one platform and a group of drivers whose actions mutually affect each other. At the beginning
24 stage of MF-MDP studies, we develop a specific MF-MDP model for analyzing spatial-temporal
25 subsidies for drivers with self-relocation strategies (see Section 4.1) and an efficient solution
26 algorithm for the particular MF-MDP model (see Section 4.2). We demonstrate that the
27 algorithm achieves significant computational advantages, faster convergence, and better
28 performance on a small-scale market (see Section 5), and aim to examine the general
29 performance on actual-size problems in near future.

30 3.2 FORMULATION OF THE RIDE-SOURCING MF-MDP MODEL

31 In a ride-sourcing market, the platform's operational strategies play important roles in
32 affecting the performance (e.g., daily income, waiting time for order matches, and distances en
33 route to pick up passengers) and decisions (e.g., self-relocation and working hours) of drivers.
34 However, if the number of drivers is large, the impact of each individual driver's decisions and
35 actions on the platform or other drivers is trivial and can be ignored without causing significant
36 deviations in general. By contrast, the average (i.e., MF) state of all drivers collectively, which
37 captures the spatial-temporal supply information, will significantly influence order
38 matching/dispatching, performance (e.g., net revenue, vehicle occupied rate, and the number of
39 passengers served), and other decisions (e.g., spatial-temporal pricing and subsidy) of the
40 platform as well as those of each individual driver. Moreover, the platform sometimes chooses
41 to display heat maps of its spatial-temporal surge pricing and/or subsidy and overall demand and
42 supply to drivers on the app. In this manner, the state of the platform and the MF state of drivers
43 are public information to drivers, who then process the information and take corresponding
44 actions. Therefore, it is reasonable to describe the ride-sourcing market using an MF-MDP

1 model, in which the platform is regarded as the major agent and a number of drivers are treated
2 as minor agents².

3 Suppose there is 1 platform and M homogeneous drivers (i.e., state sets, action sets, and
4 objectives are the same for drivers), and the planning horizon consists of T time periods (i.e.,
5 time $t \in \{1, 2, \dots, T\}$). Let \mathcal{S} and \mathcal{S}_d denote finite sets of the states of the platform and drivers,
6 respectively. Let $y^t \in \mathcal{S}$ represent the state of the platform at time t ; specifically, y^t can be a
7 vector that contains time index t , the spatial-temporal pricing, subsidies, and number of waiting
8 passengers across different regions in the market at time t . We use $y_{d,i}^t \in \mathcal{S}_d$ to represent the
9 state of driver $i \in \{1, 2, \dots, M\}$ at time t , which could include time index, their location, the
10 number of loaded passengers, and the destination. Then the MF state of all drivers at any time
11 period t can be represented as a vector \mathbf{z}_d^t as follows:

$$\mathbf{z}_d^t = [z_{d,s_d}^t]_{1 \times |\mathcal{S}_d|} \quad (1)$$

$$z_{d,s_d}^t = \frac{\sum_{i=1}^M \mathbb{I}(y_{d,i}^t = s_d)}{M} \quad (2)$$

12 where $\mathbb{I}(\cdot)$ denotes the identity function and we use \mathbf{H}_d to denote the feasible domain of MF state
13 \mathbf{z}_d^t , i.e., $\mathbf{z}_d^t \in \mathbf{H}_d$. Intuitively, the MF vector \mathbf{z}_d^t represents the distribution of drivers' states. For
14 instance, if a driver's state contains their current location and the occupancy of their vehicle,
15 then the MF state captures the spatial distribution of all vacant vehicles and occupied vehicles.

16 Let \mathbf{A} and \mathbf{A}_d denote finite sets of the actions of the platform and drivers, respectively.
17 We use $x^t \in \mathbf{A}$ and $x_{d,i}^t \in \mathbf{A}_d$, respectively, to denote the actions of the platform and driver i at
18 time t . The actions of the platform can include pricing or subsidy strategies (e.g., 1 for
19 subsidizing and 0 for not offering subsidy), and the actions of a driver are their self-relocation
20 directions.

21 Following the conventions in discrete-time MDPs, the transition probability of a major
22 or minor agent in the MF-MDP is determined by their current state and action and the MF state
23 of the minor agents. Specifically, the state transition laws for the platform and a specific driver
24 are denoted as $Q(\cdot | \cdot)$ and $Q_d(\cdot | \cdot)$ in Eqs. (3)–(4), where $P(\cdot)$ denotes the probability operator³.

$$Q(s' | s, \mathbf{h}_d, a) = P(y^{t+1} = s' | y^t = s, \mathbf{z}_d^t = \mathbf{h}_d, x^t = a) \quad (3)$$

$$Q_d(s'_d | s_d, \mathbf{h}_d, a_d) = P(y_{d,i}^{t+1} = s'_d | y_{d,i}^t = s_d, \mathbf{z}_d^t = \mathbf{h}_d, x_{d,i}^t = a_d) \quad (4)$$

25 The platform or a driver takes sequential actions to maximize their total rewards in T
26 time periods, which can be measured by the net revenue, the number of passengers served, and
27 so on. Let r denote the reward of the platform, which is a function of the platform's current state

² Note that in real ride-sourcing markets, market conditions have strong time-varying patterns with peak and off-peak hours, which indicate the nonstationary states and transitions in a day. However, if we consider a certain period of 2 to 3 hours, market conditions are more stable and thus can be approximately described using stationary states and transitions. Readers can refer to Figures 4 and 5 in [Lyu et al. \(2019\)](#) for demonstrations of daily temporal distributions of demand and supply in a real ride-sourcing market. If we consider a certain period—e.g., 8 am to 10 am during peak hours or 2 pm to 4 pm during off-peak hours—market conditions are quite stable and thus can be modeled as stationary MDP, with different transition matrices, respectively.

³ In this paper, we use y^t and $y_{d,i}^t$ (also y_d^t) to represent random variables of states, x^t and $x_{d,i}^t$ (also x_d^t) random variables of actions, and \mathbf{z}_d^t (also $\hat{\mathbf{z}}_d^t$) random variables of MF states in the MF-MDP model. We use s and s_d to represent values of random states, a and a_d values of random actions, and \mathbf{h}_d values of random MF states.

1 and action and the MF state of drivers. For a particular driver, the reward r_d could be measured
 2 by their income, saving on idle-cruise distance, saving on operational costs, etc., and it is a
 3 function of their current state and action, the current state of the platform, and the current MF
 4 state⁴. The total rewards of the platform and a specific driver, which are also referred to as value
 5 functions, are given by

$$V^\pi(s, \mathbf{h}_d) = \mathbb{E}_\pi(\sum_{t=1}^T (\rho)^t r(y^t, \mathbf{z}_d^t, x^t) | y^1 = s, \mathbf{z}_d^1 = \mathbf{h}_d) \quad (5)$$

$$V_d^{\pi_{d,i}}(s, s_d, \mathbf{h}_d) = \mathbb{E}_{\pi_{d,i}}(\sum_{t=1}^T (\rho)^t r_d(y^t, y_{d,i}^t, \mathbf{z}_d^t, x_{d,i}^t) | y^1 = s, y_{d,i}^1 = s_d, \mathbf{z}_d^1 = \mathbf{h}_d) \quad (6)$$

6 where V^π and $V_d^{\pi_{d,i}}$ denote the total rewards for the platform and driver i given some specific
 7 initial states (i.e., $y^1 = s$, $y_{d,i}^1 = s_d$, and $\mathbf{z}_d^1 = \mathbf{h}_d$), respectively; π and $\pi_{d,i}$ denote the policies
 8 (a mapping from states to actions) of the platform and the i -th driver respectively; $x^t =$
 9 $\pi(y^t, \mathbf{z}_d^t)$ and $x_{d,i}^t = \pi_{d,i}(y^t, y_{d,i}^t, \mathbf{z}_d^t)$ represent the actions following the corresponding policies;
 10 $\rho \in (0,1)$ is the discount factor that measures how the policy balances the trade-off between
 11 immediate reward and long-term rewards; and $\mathbb{E}_\pi(\cdot)$ and $\mathbb{E}_{\pi_{d,i}}(\cdot)$ are the expectation operators
 12 under policies π and $\pi_{d,i}$, respectively.

13 Given specific formulas for rewards and state transition laws, a straightforward approach
 14 to solving the ride-sourcing MF-MDP model is to regard the platform and each driver as an
 15 agent, then try to solve the problem with a decentralized multi-agent MDP approach. However,
 16 the decentralized multi-agent MDP is generally hard to solve, especially when there are many
 17 agents. In reality, we will have a large number of minor agents (drivers). The distinct objectives
 18 of the major agent (platform) and minor agents (drivers) also render the solution-seeking process
 19 more unstable and intractable. Alternatively, we approximate the random MF vector \mathbf{z}_d^t as a
 20 stationary process and optimize an aggregate policy for all drivers. Namely, as $M \rightarrow \infty$, we have
 21 $\mathbf{z}_d^t \xrightarrow{a.s.} \hat{\mathbf{z}}_d^t$. Similar approximations of asymptotic processes of homogeneous decision-makers
 22 have been adopted in studies of day-to-day traffic dynamics (Hazelton and Watling, 2004; Zhu
 23 et al., 2019b; Zhu et al., 2020a). In the simplified MF-MDP model, the platform takes actions
 24 according to policy π (i.e., $x^t = \pi(y^t, \hat{\mathbf{z}}_d^t)$), and the decision processes of all drivers are
 25 determined by policy π_d (i.e., $x_d^t = \pi_d(y^t, y_{d,i}^t, \hat{\mathbf{z}}_d^t)$) of a representative driver⁵. The MF state at
 26 the next time period depends on the current MF state and the platform's state, which is simplified
 27 as an updating rule $\hat{\mathbf{z}}_d^{t+1} = l_d(y^t, \hat{\mathbf{z}}_d^t)$. Note that the updating rule also incorporates the policies
 28 (for action taking) of the platform and the representative driver. Therefore, the standard multi-
 29 agent MF-MDP model with $1 + M$ agents can be reduced to a simplified MF-MDP model with
 30 only 2 agents:

- 31 • The ride-sourcing platform that acts as a major agent to design the optimal policy to
 32 maximize its total rewards. The optimal value function is defined as $V^*(s, \mathbf{h}_d) =$
 33 $\max_\pi \mathbb{E}_\pi(\sum_{t=1}^T \rho^t r(y^t, \mathbf{z}_d^t, x^t) | y^1 = s, \mathbf{z}_d^1 = \mathbf{h}_d)$.
- 34 • A representative driver who acts as a representative minor agent to pursue the optimal
 35 policy and maximize their total rewards. The total reward is regarded as the average total

⁴ With specific research problems in ride-sourcing markets, we sometimes need to incorporate the previous state (i.e., y^{t-1} , $y_{d,i}^{t-1}$, and \mathbf{z}_d^{t-1}) into the formulas for rewards (i.e., r and r_d). This is because the before-and-after changes in states may affect the reward. For instance, if a subsidy is offered to a driver upon a new match with a passenger, we must check the driver's previous state and include the subsidy in the reward only if the current state is "matched/dispatched" and the previous state is "idle".

⁵ For convenience and clarity, we use notation without a driver index to denote the state (y_d^t), action (x_d^t), and policy (π_d) of the representative driver in the simplified MF-MDP model.

1 rewards of all drivers. The optimal value function is defined as $V_d^*(s, s_d, \mathbf{h}_d) =$
 2 $\max_{\pi_d} \mathbb{E}_{\pi_d} (\sum_{t=1}^T \rho^t r_d(y^t, y_d^t, \mathbf{z}_d^t, x_d^t) | y^1 = s, y_d^1 = s_d, \mathbf{z}_d^1 = \mathbf{h}_d)$.

3 The form of function $l_d(y^t, \hat{\mathbf{z}}_d^t)$ determines the consistency between the approximated
 4 MF state $\hat{\mathbf{z}}_d^t$ and the exact MF state \mathbf{z}_d^t (i.e., Eqs. (1)–(2)). A consistent approximation of the MF
 5 states is a critical requirement, such that the simplified MF-MDP model is able to represent the
 6 complex state transition and decision dynamics characterized by the standard MF-MDP. We
 7 discuss the consistency requirement in Section 3.3.

8 3.3 OPTIMAL POLICIES AND THE CONSISTENCY REQUIREMENT

9 An MDP model is generally solved by Bellman equations. We first illustrate the Bellman
 10 equations of the simplified MF-MDP model. An arbitrary MF state updating rule $l_d(y^t, \hat{\mathbf{z}}_d^t)$ is
 11 adopted without checking the consistency between $\hat{\mathbf{z}}_d^t$ and \mathbf{z}_d^t . The following propositions are
 12 necessary to obtain the optimal policies with Bellman equations:

13 **Proposition 1.** \mathbf{H}_d is a continuous and compact set.

14 **Proposition 2.** Given a continuous reward function $r(y^t, \hat{\mathbf{z}}_d^t, x^t)$ on \mathbf{H}_d , the value function
 15 $V^\pi(s, \mathbf{h}_d)$ is continuous on \mathbf{H}_d .

16 **Proposition 3.** Given a continuous reward function $r_d(y^t, y_d^t, \hat{\mathbf{z}}_d^t, x_d^t)$ on \mathbf{H}_d , the value function
 17 $V_d^{\pi_d}(s, s_d, \mathbf{h}_d)$ is continuous on \mathbf{H}_d .

18 where **Proposition 1** is straightforward because \mathbf{z}_d^t is continuous as M goes to infinity, and given
 19 specific policies π and π_d , the reward functions (i.e., r and r_d) and the corresponding value
 20 functions (i.e., V^π and $V_d^{\pi_d}$) are continuous, leading to **Propositions 2** and **3**.

21 The optimal policy for the platform can be solved based on the following Bellman
 22 equation:

$$V^*(s, \mathbf{h}_d) = \max_{a \in A} \left\{ r(s, \mathbf{h}_d, a) + \rho \sum_{s' \in \mathcal{S}} Q(s'|s, \mathbf{h}_d, a) V(s', \mathbf{h}'_d) \right\} \quad (7)$$

23 where $\mathbf{h}' = l_d(s, \mathbf{h}_d)$. In light of **Proposition 2**, the existence of an optimal policy π^* for Eq.
 24 (7) is guaranteed. Suppose the optimal policy π^* has been implemented in the simplified MF-
 25 MDP model. The Bellman equation for the representative driver is given by

$$V_d^*(s, s_d, \mathbf{h}_d) = \max_{a_d \in A_d} \left\{ \rho \sum_{\substack{s' \in \mathcal{S}, \\ s'_d \in \mathcal{S}_d}} Q(s'|s, \mathbf{h}_d, a) Q_d(s'_d|s_d, \mathbf{h}_d, a_d) V_d(s', s'_d, \mathbf{h}'_d) + r_d(s, s_d, \mathbf{h}_d, a_d) \right\} \quad (8)$$

26 where $a = \pi^*(s, \mathbf{h}_d)$.

27 Similarly, based on **Proposition 3**, given π^* , the optimal policy π_d^* exists for Eq. (8). In
 28 other words, there is an optimal policy group (π^*, π_d^*) that simultaneously satisfies Eqs. (7)–(8).

29 Next, we seek the specific formula of $l_d(y^t, \hat{\mathbf{z}}_d^t)$ for a consistent approximation of the
 30 MF state. The basic idea is to identify some updating rule of the exact MF state \mathbf{z}_d^t in the
 31 simplified MF-MDP model, then adapt this rule to the approximated MF state $\hat{\mathbf{z}}_d^t$. Based on Eq.
 32 (2), we obtain the asymptotic \mathbf{z}_d^t as M goes to infinity:

$$\lim_{M \rightarrow \infty} z_{d,s_d}^t = \lim_{M \rightarrow \infty} \frac{\sum_{i=1}^M I(y_{d,i}^t = s_d)}{M} \xrightarrow{a.s.} P(y_d^t = s_d) \quad (9)$$

1 To examine the asymptotic property of \mathbf{z}_d^t under the optimal policy group (π^*, π_d^*) , we
 2 introduce the following theorem, which is valid for any function $\hat{\mathbf{z}}_d^{t+1} = l_d(y^t, \hat{\mathbf{z}}_d^t)$.

3 **Theorem 1.** Let policy group (π^*, π_d^*) denote the optimal policies of Eqs. (7)–(8); the underlying
 4 vector $(y^t, y_d^t, \hat{\mathbf{z}}_d^t)$ forms a stationary Markov process.

5 **Proof.** The policy group provides stationary mapping from states to actions, such that $x^t =$
 6 $\pi^*(y^t, \hat{\mathbf{z}}_d^t)$ and $x_d^t = \pi_d^*(y^t, y_d^t, \hat{\mathbf{z}}_d^t)$. The state transition probability from state (s, s_d, \mathbf{h}_d) to
 7 state $(s', s'_d, \mathbf{h}'_d)$ is given by

$$\begin{aligned} & P(y^{t+1} = s', y_d^{t+1} = s'_d, \hat{\mathbf{z}}_d^{t+1} = \mathbf{h}'_d | y^t = s, y_d^t = s_d, \hat{\mathbf{z}}_d^t = \mathbf{h}_d) \\ & = Q(s' | s, \mathbf{h}_d, \pi^*(s, \mathbf{h}_d)) Q_d(s'_d | s_d, \mathbf{h}_d, \pi_d^*(s, s_d, \mathbf{h}_d)) I(\mathbf{h}'_d = l_d(s, \mathbf{h}_d)) \end{aligned} \quad (10)$$

8 where the LHS only depends on the current state (s, s_d, \mathbf{h}_d) .

9 ■

10 In light of **Theorem 1**, the asymptotic z_{d,s_d}^t and $P(y_d^t = s_d)$ are also Markov processes.
 11 Based on the transition law, the formula of $P(y_d^{t+1} = s'_d)$ is given by

$$P(y_d^{t+1} = s'_d) = \sum_{s_d \in \mathcal{S}_d} P(y_d^t = s_d) Q_d(s'_d | s_d, \hat{\mathbf{z}}_d^t, \pi_d^*(y^t, s_d, \hat{\mathbf{z}}_d^t)) \quad (11)$$

12 Eq. (11) is summarized as a matrix product form:

$$\mathbf{z}_d^{t+1} = \mathbf{z}_d^t \hat{\mathbf{Q}}_d(y^t, \hat{\mathbf{z}}_d^t) \quad (12)$$

13 where $\hat{\mathbf{Q}}_d(y^t, \hat{\mathbf{z}}_d^t) = [Q_d(s'_d | s_d, \hat{\mathbf{z}}_d^t, \pi_d^*(y^t, s_d, \hat{\mathbf{z}}_d^t))]_{|s_d| \times |s_d|}$ is a probability transition matrix of
 14 the MF state. Let $\hat{\mathbf{z}}_d^1 = \mathbf{z}_d^1$ and $l_d(y^t, \hat{\mathbf{z}}_d^t) = \hat{\mathbf{z}}_d^t \hat{\mathbf{Q}}_d(y^t, \hat{\mathbf{z}}_d^t)$; for any $t \in \mathbf{T}$, we can obtain the
 15 following equation by iteratively substituting Eq. (12) and $l_d(y^t, \hat{\mathbf{z}}_d^t)$.

$$\hat{\mathbf{z}}_d^{t+1} = \hat{\mathbf{z}}_d^1 \prod_{t'=1}^t \hat{\mathbf{Q}}_d(y^{t'}, \hat{\mathbf{z}}_d^{t'}) = \mathbf{z}_d^1 \prod_{t'=1}^t \hat{\mathbf{Q}}_d(y^{t'}, \hat{\mathbf{z}}_d^{t'}) = \mathbf{z}_d^{t+1} \quad (13)$$

16 Therefore, we conclude that $E_{\pi^*, \pi_d^*}(\hat{\mathbf{z}}_d^t) = \mathbf{z}_d^t$ and the consistency requirement for the
 17 approximation of MF states reduces to the following updating rule:

$$\hat{\mathbf{z}}_d^{t+1} = l_d^\#(y^t, \hat{\mathbf{z}}_d^t) = \hat{\mathbf{z}}_d^t \hat{\mathbf{Q}}_d(y^t, \hat{\mathbf{z}}_d^t) \quad (14)$$

18 where superscript # means that the updating rule is consistent.

19 We refer to the combination of the optimal policies for the platform and the representative
 20 driver and the consistent updating rule for MF states, i.e., $(\pi^*, \pi_d^*, l_d^\#(y^t, \hat{\mathbf{z}}_d^t))$, as a consistent
 21 optimal solution of the simplified MF-MDP model. Note that $(\pi^*, \pi_d^*, l_d^\#(y^t, \hat{\mathbf{z}}_d^t))$ satisfies Eqs.
 22 (7), (8), and (14) simultaneously. The consistency of the stochastic process depicted in Eq. (14)

1 requires a soft policy for the representative driver, i.e., $a_d | \pi_d \sim P(x_d^{t+1} = a_d | \pi_d(y^t, y_d^t, \hat{z}_d^t))$. In
 2 contrast to a “hard” policy that selects a deterministic action given the observed state, a “soft”
 3 policy is a probabilistic distribution over the action set, and the agent stochastically selects an
 4 action according to the distribution given any observed state. The design of the soft policy
 5 enables the model to use a single policy to represent the aggregate actions of all drivers, rather
 6 than determining different policies for each driver.

7 The generalized MF-MDP model and its theoretical guarantees in terms of simplification
 8 and optimal solution seeking allow us to formulate and solve a variety of research questions in
 9 ride-sourcing markets. For instance, we can use this model to delineate the dynamics of a ride-
 10 sourcing market in which drivers (minor agents) have self-relocating behaviors for maximizing
 11 their individual earnings while a platform tries to achieve a more efficient system by imposing
 12 spatial-temporal pricing/subsidy strategies (see Section 4). As described in this subsection, the
 13 simplified MF-MDP model contributes to the solution algorithm of MDPs with multiple mixed
 14 agents. However, due to complex interactions between the platform and drivers, the formulas of
 15 $Q(\cdot | \cdot)$, $Q_d(\cdot | \cdot)$, r and r_d can be complicated. Moreover, the solution space with respect to
 16 states, actions, and time periods can be extremely large. Therefore, it is generally difficult to
 17 obtain exact optimal policies via solving the Bellman equations. A typical method is to use
 18 simulations to approximate the interactions between the environment and agents, and attempt to
 19 find close-optimal policies through reinforcement learning-based algorithms (Wang et al., 2018;
 20 Li et al., 2019; Jin et al., 2019; Shou et al., 2020b). The idea of a soft policy⁶ for the representative
 21 driver and the consistent updating rule for the approximated MF state are valuable for designing
 22 computationally efficient simulation processes.

24 4 DESIGN AND ANALYZE SUBSIDIES FOR DRIVERS WITH SELF-RELOCATION

25 In this section, we substantialize the proposed generalized MF-MDP model in a specific
 26 research problem in which a platform tries to better allocate spatial-temporal subsidies for drivers,
 27 while drivers attempt to maximize their individual earnings by self-relocation. The formulation
 28 (in terms of states, actions, and rewards) of this model is introduced in Section 4.1, while a
 29 representative-agent reinforcement learning algorithm for solving the model is developed in
 30 Section 4.2.

31 4.1 FORMULATION OF THE SPECIFIC RIDE-SOURCING MF-MDP MODEL

32 In this subsection, we provide definitions and intuitive explanations of the states, actions,
 33 and rewards of the platform and drivers, and introduce a matching rule between passengers and
 34 drivers. Readers can refer to Appendix A for detailed mathematical formulations of the state
 35 transition laws, order-matching probabilities, and rewards.

36 In a ride-sourcing market with spatial-temporal imbalance between demand and supply.
 37 We use a hexagonal zone system, which has been used in some previous studies (Ke et al., 2019;
 38 Xu et al. 2018; Lin et al., 2018) and DiDi’s ride-sourcing simulator (Xu et al. 2017). There are
 39 O hexagonal zones, passenger demand is exogenous, and driver supply can be characterized by
 40 the MF state z_d^t of M drivers (also by the approximate state \hat{z}_d^t in the simplified MF-MDP
 41 model). Other notation is the same as in the generalized model in Section 3. The platform could
 42 lose passengers in zones with high demand and insufficient idle vehicles. To increase net revenue

⁶ The soft policy maps a state to a probability distribution over all possible actions. Given a MF state, the representative driver takes a stochastic action based on a probability distribution that is determined by the soft policy. In this way, we can approximate the collective behavior/decision of a group of drivers by using one representative driver at the expense of a small measuring error, especially when the number of drivers is large.

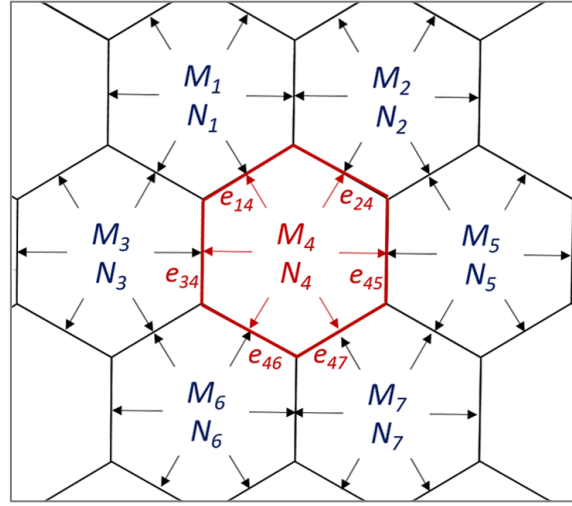
1 and the number of passengers served, the platform may offer spatial-temporal (time- and zone-
 2 based) subsidies to incentivize idle drivers to move from demand-cool zones to demand-hot
 3 zones. Meanwhile, drivers design their self-relocation strategies to increase their own income.

4 We begin with the states and actions of the platform. To model the subsidy strategy across
 5 different zones, the state of the platform is characterized by an $1 + O$ dimensional vector $\mathbf{s} =$
 6 $[t, s_1, \dots, s_O]$, where $s_o \in \{0, \beta\}$ denotes the subsidy in zone $o \in \{1, 2, \dots, O\}$ such that 0 and β
 7 refer to “no subsidy” and “offering a subsidy,” respectively, and β is a predefined amount of
 8 subsidy per ride (i.e., subsidy rate). In theory, we allow β to be non-positive values in the model,
 9 and $\beta = 0$ means a “non-subsidy” strategy and $\beta < 0$ indicates that drivers pay an extra
 10 “charge” rather than get subsidies. If zone o is subsidized at time t , drivers who are matched and
 11 dispatched to passengers originating from zone o at time t will be offered the same amount of
 12 subsidy (such a scheme is referred to as a uniform subsidy scheme). The platform’s action with
 13 respect to subsidy is represented by a vector $\mathbf{a} = [a_1, \dots, a_O]$, where $a_o \in \{0, \beta\}$. Therefore, we
 14 have $\mathbf{S} = \{[t, s_1, \dots, s_O] | s_o \in \{0, \beta\}, o \in \{1, 2, \dots, O\}, t \in \{1, 2, \dots, T\}\}$ and $\mathbf{A} =$
 15 $\{[a_1, \dots, a_O] | a_o \in \{0, \beta\}, o \in \{1, 2, \dots, O\}\}$. In this manner, both the state vector \mathbf{s} and the action
 16 vector \mathbf{a} represent the spatial distribution of subsidy, and the state of the platform for the next
 17 time period is identical to its current action. Eq (A.1) in Appendix A gives a mathematical
 18 formula of the state transition law $Q(\mathbf{s}' | \mathbf{s}, \mathbf{a})$.

19 Next, we introduce the states and actions of drivers. To comprehensively describe a
 20 driver’s different status (e.g., idle, on the way to pick up passengers, delivering passengers), we
 21 formulate a driver’s state as a six-dimensional vector $\mathbf{s}_d = [t, s_{d1}, s_{d2}, s_{d3}, s_{d4}, s_{d5}]$. Here, $s_{d1} \in$
 22 $\{0, 1, 2\}$ denotes the current task of the driver/vehicle with 0, 1, and 2 representing “idle,”
 23 “picking up a passenger,” and “delivering a passenger,” respectively; $s_{d2} \in \{1, 2, \dots, T\}$ denotes
 24 the remaining time periods before finishing the current status; $s_{d3} \in \mathbf{O}$ denotes the idling zone
 25 in which the driver is idle and waiting for a match; and $s_{d4} \in \mathbf{O}$ and $s_{d5} \in \mathbf{O}$ denote the origin
 26 and destination of the current passenger order, respectively. The value of s_{d2} depends on the
 27 travel time on the zone network. The action of the driver is the destination zone of self-relocation.
 28 Given the driver’s current idling zone o (i.e., $s_{d3} = o$), the set of their actions (i.e., $\mathbf{A}_d | o$) is
 29 represented by set \mathbf{J}_o , which is the set of adjacent zones of o plus o itself. The action $a_d = o$
 30 means the driver will stay in the current zone, and other actions $a_d \in \mathbf{J}_o / o$ indicate that the driver
 31 will relocate to an adjacent zone. A self-relocation action is only needed when the driver has no
 32 picking-up or delivering tasks and is not on the way of cruising to an adjacent zone, i.e., when
 33 $s_{d1} = 0$ and $s_{d2} = 0$ (referred to as a “purely idle” state). Therefore, one only optimizes policies
 34 in the “purely idle” state-action space (i.e., the Cartesian product of set $\{\mathbf{s}_d \in \mathbf{S}_d | s_{d1} = 0, s_{d2} =$
 35 $0\}$ and the action set). Given a large number of drivers, drivers with the “purely idle” state are in
 36 different pairs of (t, s_{d3}) that could uniformly distribute across the spatial-temporal domain of the
 37 scenario, making the state-action-reward transitions in the learning process non-sparse.
 38 Furthermore, a state with $s_{d1} = 0$ and $s_{d2} > 0$ indicates a “self-relocating” state with a
 39 relocation destination such that no action is needed until he/she arrives at the destination and
 40 becomes purely idle. Following Eq. (4), the state transition law for a driver, i.e.,
 41 $Q_d(\mathbf{s}'_d | \mathbf{s}_d, \mathbf{h}_d, a_d)$, depends on its current state and action as well as the MF state of drivers;
 42 detailed formulas are given by Eqs. (A.2)–(A.5) in Appendix A.

43 Next, we introduce the order-matching rule between drivers and passengers. If a driver is
 44 in a “purely idle” or “self-relocating” state, they have a chance to be matched with a passenger.
 45 Therefore, the state transition probability of the driver is substantially affected by the matching
 46 rule. Generally, the platform considers a maximal matching radius that only prevents passengers
 47 from being matched with far away drivers. A larger radius allows a larger flexibility in matching;
 48 namely, a larger pool of candidate idle drivers is generated for each passenger, and thus the

1 matching rate becomes larger; this also indicates a smaller passengers' expected waiting time.
 2 However, since some distant drivers may be matched to passengers, a larger matching radius
 3 will increase the average pick-up time. MDP-based models in the literature usually adopt a small
 4 matching radius so that drivers and passengers can be matched only if they are in the same zone
 5 (e.g., Shou et al., 2020b). Such matching rules ignore the cross-region dispatching and picking-
 6 up events that are commonly observed in ride-sourcing services, and thus are more suitable for
 7 taxi markets rather than for ride-sourcing markets. To allow cross-region matching between
 8 passengers and drivers, in this paper we propose an edge-based matching rule in calculating
 9 transition laws (Figure 2). Termination "edge-based" means that we allow drivers in zone o to
 10 be matched with passengers in zones $o' \in J_o$; if the driver is matched with a passenger in zone
 11 $o' = o$, they immediately pick up the passenger and start the delivering task; if the driver is
 12 matched with a passenger in an adjacent zone $o' \in J_o/o$, they must spend some time on the
 13 picking up task before delivering the passenger. For simplicity, we assume that drivers and
 14 passengers in each hexagonal zone is uniformly distributed (which does not mean that they are
 15 uniform across the network with many zones). The edge-based rule matches drivers and
 16 passengers near each common edge between zone o and its adjacent zones $o' \in J_o/o$.



17
 18 Figure 2. Edge-based matching.

19 Let M_o denote the number of idle drivers in zone o (i.e., with state $s_{d1} = 0$ and $s_{d3} = o$);
 20 N_o the number of passengers with origin in zone o ; and $e_{oo'}$ the common edge between two
 21 hexagonal zones o and o' . At each time period, M_o is obtained from the MF state and N_o is
 22 observable and thus exogenously given. We illustrate the number of matches near edge $e_{oo'}$
 23 using a simple example. Taking the zone indices in Figure 2, for instance, $J_4 = \{2,5,7,6,3,1,4\}$,
 24 $J_4/\{4\} = \{2,5,7,6,3,1\}$, and we match drivers and passengers near edge e_{14} . With uniformly
 25 distributed demand and supply, there are $\frac{M_4}{6}$ idle drivers and $\frac{N_4}{6}$ passengers near e_{14} in zone 4,
 26 and $\frac{M_1}{6}$ idle drivers and $\frac{N_1}{6}$ passengers near this edge in zone 1. Therefore, there are a total of
 27 $\frac{M_1+M_4}{6}$ idle drivers and $\frac{N_1+N_4}{6}$ passengers to be matched near edge e_{14} ; for these passengers and
 28 drivers, we allow a driver/passenger in zone 4 to be matched with passengers/drivers in either
 29 zone 1 or zone 4. Based on a matching rule in Yu et al. (2019), the number of matches near e_{14}
 30 is approximated as $\min\left\{\frac{M_1+M_4}{6}, \frac{N_1+N_4}{6}\right\}$. Similar to this example, we can compute the number of
 31 matches near an arbitrary edge.

32 To approximate the matching probabilities of a driver, we assume that the numbers of
 33 matched passengers and drivers are proportional to the corresponding demand and supply near
 34 the common edge. To continue with the example above, if a driver is in zone 1, the probability

1 that they are near edge e_{14} equals $\frac{1}{6}$, and the probability that they get a passenger order near edge
2 e_{14} is $\frac{1}{6} \frac{\min\{\frac{M_1+M_4}{6}, \frac{N_1+N_4}{6}\}}{\frac{M_1+M_4}{6}} = \frac{\min\{\frac{M_1+M_4}{6}, \frac{N_1+N_4}{6}\}}{M_1+M_4}$. Detailed formulas for calculating the number of
3 matches and driver-side matching probabilities are given in Eqs. (A.6)–(A.9) in Appendix A.
4 We need the MF state to compute M_o (also denoted as $M_o(\mathbf{h}_d)$ in Eq. (A.7)) and then the
5 matching probabilities; this explains why the state transition law of a driver depends on the MF
6 state, i.e., $Q_d(\mathbf{s}'_d | \mathbf{s}_d, \mathbf{h}_d, a_d)$. To our best knowledge, this paper is among the first idle vehicle
7 relocation studies that consider cross-zone matches with the proposed edge-based matching rule.

8 Last, we discuss the rewards for the platform and drivers. We consider that the platform's
9 objective is to maximize a weighted sum of the net revenue and the service rate, which is defined
10 by the number of passengers served divided by the total passenger demand. Intuitively, the
11 service rate reflects passengers' satisfaction, and a low service rate may cause a decrease in
12 passenger demand in the long run and affect the platform's market share. Our motivation to set
13 this objective structure is that the platform usually needs to make a trade-off between net revenue
14 (short-term benefits) and customer service rate (long-term interests). To be more specific, the
15 reward (also referred to as the objective value) of the platform is formulated by
16 $r(\mathbf{y}^t = \mathbf{s}, \mathbf{z}_d^t = \mathbf{h}_d, \mathbf{z}_d^{t-1} = \mathbf{h}'_d) = r_1(\mathbf{h}_d) - r_2(\mathbf{s}, \mathbf{h}_d, \mathbf{h}'_d) + \mu r_3(\mathbf{s}, \mathbf{h}_d, \mathbf{h}'_d)$, where r_1 refers to
17 the commission withheld from trip fares by the platform; r_2 is the amount of subsidies offered
18 to drivers; r_3 is the service rate; and μ denotes the weight of the service rate for the platform⁷. In
19 addition to the MF state of drivers and the state of the platform, the calculation of r involves the
20 following predefined variables: the ride-sourcing trip fare per time period (i.e., trip fare rate) α ,
21 commission rate for the platform η , and total passenger demand across the entire operational
22 horizon N . Readers can refer to Eqs. (A.10)–(A.13) in Appendix A for detailed formulas for r_1 ,
23 r_2 , and r_3 .

24 A driver's objective is to maximize the total income. The reward (referred to as income)
25 for a particular driver is the sum of the trip fare and subsidy offered by the platform, i.e.,
26 $r_d(\mathbf{y}^t = \mathbf{s}, \mathbf{y}_{d,i}^t = \mathbf{s}_d, \mathbf{y}_{d,i}^{t-1} = \mathbf{s}'_d) = r_{d1}(\mathbf{s}_d) + r_{d2}(\mathbf{s}, \mathbf{s}_d, \mathbf{s}'_d)$, where r_{d1} and r_{d2} denote the
27 income from the trip fare and the subsidy, respectively. The income from fare is provided
28 gradually during the delivery task, while the subsidy is a one-time reward upon a match if the
29 origin of the passenger is subsidized. Similar to r , we need α and η to compute r_d . Detailed
30 formulas for r_{d1} and r_{d2} are given in Eqs. (A.14)–(A.16) in Appendix A.

31 To sum up, we provide the formulation of a specific standard MF-MDP model, in which
32 the MF state of drivers is mainly used to compute the order-matching probabilities. As a result,
33 the MF state directly determines drivers' state transition law and the platform's reward; it also
34 affects the matching outcome of an individual driver and their reward. Approximation of the MF
35 state and simplification of the MF-MDP model play an important role in solution-finding. For
36 the simplified MF-MDP model, the states, actions, transition laws, and rewards for the
37 representative driver are the same as for an arbitrary driver in the standard model. In light of Eq.
38 (14), the consistent updating rule for the approximated MF state is $\widehat{Q}_d(\mathbf{y}^t, \widehat{\mathbf{z}}_d^t) = \widehat{Q}_d(\widehat{\mathbf{z}}_d^t)$, which
39 can be summarized based on Eqs. (A.4)–(A.5).

40 4.2 SOLUTION ALGORITHMS

41 Given the state transition laws for the platform and drivers in Eqs. (A.1)–(A.5), the order-
42 matching probabilities in Eqs. (A.6)–(A.9), and the rewards in Eqs. (A.10)–(A.16), we have a
43 specific formulation of a standard MF-MDP model with multiple mixed agents. Although the

⁷ As stated in footnote 3, we consider the previous state \mathbf{z}_d^{t-1} in the formulation of r and consider $\mathbf{y}_{d,i}^{t-1}$ in r_d .

1 rewards of the platform and drivers and the state transition law of the platform have deterministic
 2 formulas, drivers' state changes are dependent on the stochastic order-matching process with
 3 numerous possible outcomes (e.g., different origins and destinations of passenger orders),
 4 making it difficult to obtain an exact solution via Bellman equations. As discussed in Section 3,
 5 reinforcement learning algorithms can be adopted to solve such a multi-agent MF-MDP model
 6 with large state and action sets, complex state transition laws, and reward formulas. We consider
 7 two solution-seeking approaches below.

- 8 • The multi-agent approach, which solves the standard MF-MDP model with 1 platform
 9 and M drivers. In a reinforcement learning algorithm, each of the $1 + M$ agents learns
 10 their own decision policy⁸, which can be characterized by Q-tables, neural networks, etc.
 11 An agent takes actions based on their policy and updates the parameters of the policy via
 12 their experiences under the actions. The MF states \mathbf{z}_a^t are summarized based on the states
 13 of all drivers (Eq. (2)).
- 14 • The representative-agent approach, which solves the simplified MF-MDP model with
 15 one platform and one representative driver. In a reinforcement learning algorithm, we
 16 create two decision policies: one for the platform and the other for the representative
 17 driver. The representative driver takes actions based on a soft policy and updates the
 18 parameters via experiences under the actions. We adopt the approximated MF state $\hat{\mathbf{z}}_a^t$,
 19 which is updated according to the previous approximated MF state ($\hat{\mathbf{Q}}_a(\hat{\mathbf{z}}_a^t)$) summarized
 20 based on Eqs. (A.4)–(A.5)).

21 The multi-agent approach is proposed as a benchmark that solves the standard MF-MDP
 22 model. With a large M , the MF space is continuous and compact, and **Propositions 1** to **3**
 23 are valid for the representative-agent approach. Comparing with the benchmark, the representative-
 24 agent approach meets the consistency requirement with respect to the MF state and could be
 25 faster in terms of computation and identifying the optimal policies.

26 In this paper, the two approaches are implemented via the A2C algorithm, one of the most
 27 popular reinforcement learning algorithms (Mnih et al., 2016). For each agent, the A2C
 28 algorithm establishes two networks (also referred to as a group of networks): one policy network
 29 (or critic network) that observes the current states and generates policy, and one value network
 30 (or actor network) that evaluates the performance of the policy. Both networks are parameterized
 31 multi-layer neural networks, and their parameters (e.g., θ_p for the policy network and θ_v for the
 32 value network) are updated iteratively. The parameters of the value network θ_v can be updated
 33 by minimizing a loss function $L(\theta_v)$ defined as follows (Lin et al., 2018):

$$L(\theta_v) = \left[V_{\theta_v}(y^t) - \left(r(y^t, x^t) + \rho V_{\theta'_v}(y^t) \right) \right]^2 \quad (15)$$

34 where θ_v denote the parameters of the value network to be updated, θ'_v denote the parameters of
 35 the targeted value network, and $r(y^t, x^t)$ be the current reward. As for the policy network,
 36 parameters θ_p are updated using a gradient descent rule $\theta_p \leftarrow \theta_p + \delta \nabla_{\theta_p} G(\theta_p)$, where δ is the
 37 learning rate and $\nabla_{\theta_p} G(\theta_p)$ represents the gradient given below.

$$\nabla_{\theta_p} G(\theta_p) = \nabla_{\theta_p} \log \pi_{\theta_p}(x^t | y^t) \left[x^t + \rho V_{\theta'_v}(y^{t+1}) - V_{\theta'_v}(y^t) \right] \quad (17)$$

⁸ Note that although the drivers are homogeneous, their optimal policies can differ. That is, idle drivers who are in the same zone might have different self-relocation destinations; otherwise, they would relocate to the same destination and compete with each other for passengers, which would result in a small matching probability and a low average income.

1 where $\pi_{\theta_p}(x^t|y^t)$ refers to the action taken by the agent given state y^t according to the policy
 2 π parameterized by θ_p , $x^t + \rho V_{\theta'_v}(y^{t+1}) - V_{\theta'_v}(y^t)$ is an advantage function used to reduce the
 3 variance of the value function and approximate the policy gradient⁹.

4 The algorithm for the multi-agent approach is shown in Algorithm 1, which is referred to
 5 as the multi-agent actor-critic (MAC) algorithm (i.e., the benchmark algorithm). As mentioned
 6 previously, in the MAC algorithm each driver or the platform has a specific group of networks
 7 to characterize their own policy, which leads to a total of $1 + M$ groups of networks. The
 8 algorithm for the representative-agent approach is shown in Algorithm 2, which is referred to as
 9 the representative-agent actor-critic (RAC) algorithm. In the RAC algorithm, we propose two
 10 groups of networks: one for the platform and the other for the representative driver. In both MAC
 11 and RAC algorithms, N_E , N_S , and D denote the maximal number of learning epochs, the
 12 maximal number of learning samples, and the replay memory, respectively¹⁰. Note that in the
 13 MAC algorithm, index i denotes driver ID. We must simulate the transitions of each individual
 14 driver (i.e., steps 3.2 to 3.7, including states, actions, and rewards) based on the individual policy
 15 $\pi_{d,i}$; in the RAC algorithm, we simulate the transitions of the representative driver (i.e., steps
 16 3.2 to 3.7) according to the soft policy π_d and use index j to represent the indices of the
 17 simulated transitions regardless of the ID of a specific driver who experiences the transition. In
 18 steps 4.4 to 4.6 in the MAC algorithm, each driver learns and updates the parameters of their
 19 own value network and policy network, via a mini-batch of samples extracted from the replay
 20 memory. In the RAC algorithm, the representative driver learns and updates the soft policy in
 21 steps 4.4 to 4.6.

22 Theoretically, given unlimited computational power, the MAC algorithm might learn the
 23 optimal self-relocation policy for each driver by conducting extensive simulations and sampling
 24 sufficient transitions over the huge solution space, which contains all feasible policies for each
 25 driver and the platform. However, computational resources are generally limited in practice, and
 26 thus it is nearly impossible to generate a massive number of samples for all possible transitions.
 27 In addition, one driver’s self-relocation policy will affect other drivers’ rewards, which is
 28 reflected by the impact of the MF state on the matching outcomes of each time period. Limited
 29 computational power and complicated competitive relationships between drivers make it
 30 difficult for the MAC to find the right pathway and obtain close-optimal policies for all drivers.
 31 By contrast, the RAC algorithm has a smaller solution space (the Cartesian product of the
 32 platform’s and the representative driver’s state-action set), and thus could identify a close-
 33 optimal solution more efficiently, which may provide solutions better than those obtained with
 34 the MAC.

35 Algorithm 1. Multi-Agent Actor-Critic (MAC) to solve the standard MF-MDP model

1. Initialize the value network with a fixed value table.
 For $n_e = 1$ to N_E do:
 2. Reset simulator, get initial state \mathbf{y}^1 , $\{\mathbf{y}_{d,i}^1\}_M$, and \mathbf{z}_d^1 .
 3. Stage one: collecting experience.
 For $t = 1$ to T do:
 3.1. Decide action x^t based on policy π , and execute x^t .
 For $i = 1$ to M do:

⁹ Note that Eqs. (15)–(18) give the basic formulas of the A2C algorithm. When adopting the A2C algorithm in the MF-MDP model, r , $V_{\theta'_v}$, V_{target} are calculated based on the definitions and formulas in Section 3 and Section 4.1.

¹⁰ These are general terminations in reinforcement learning. A learning epoch refers to an iteration for the algorithm to simulate the transitions of the agents and update the parameters of their policies, and a replay memory is used to store and sample the simulated transitions.

3.2. Decide action $x_{d,i}^t$ based on policy $\pi_{d,i}$, and execute $x_{d,i}^t$.
 End for.
3.3. Based on $\{\mathbf{y}_{d,i}^1\}_M$ and \mathbf{z}_d^t , compute $Q_d(\cdot | \cdot)$ for the simulator.
3.4 Run the simulator and observe next state \mathbf{y}^{t+1} and $\{\mathbf{y}_{d,i}^{t+1}\}_M$.
3.5. Summarize \mathbf{z}_d^{t+1} based on $\{\mathbf{y}_{d,i}^{t+1}\}_M$.
 End for.
3.6. Observe reward $r^t(\mathbf{y}^t, \mathbf{z}_d^t, \mathbf{z}_d^{t-1}), \{r_d^t(\mathbf{y}^t, \mathbf{y}_{d,i}^t, \mathbf{y}_{d,i}^{t-1})\}_M$.
3.7. Store transitions $(\mathbf{y}^t, x^t, \mathbf{y}^{t+1}, r^t(\cdot))$ and $\left\{ \left(\mathbf{y}_{d,i}^t, x_{d,i}^t, \mathbf{y}_{d,i}^{t+1}, r_d^t(\cdot) \right) \right\}_M$ to D .
 End for.
4. Stage two: learning the experiences.
 For $n_s = 1$ to N_s do:
4.1. Sample a mini-batch of transitions $(\mathbf{y}^t, x^t, \mathbf{y}^{t+1}, r^t(\cdot))$ from D .
4.2. Update the platform's value networks by minimizing $L(\theta_v)$.
4.3. Update the platform's policy networks as $\theta_p \leftarrow \theta_p + \delta \nabla_{\theta_p} G(\theta_p)$.
 For $i = 1$ to M do:
4.4. Sample a mini-batch of transitions $(\mathbf{y}_{d,i}^t, x_{d,i}^t, \mathbf{y}_{d,i}^{t+1}, r_d^t(\cdot))$ from D .
4.5. Update the i th driver's value networks by minimizing $L(\theta_v)$.
4.6. Update the i th driver's policy networks as $\theta_p \leftarrow \theta_p + \delta \nabla_{\theta_p} G(\theta_p)$.
 End for.
 End for.
 End for.
5. Finish.

1 Algorithm 2. Representative-Agent Actor-Critic (RAC) to solve the simplified MF-MDP model

1. Initialize the value network with a fixed value table.
 For $n_e = 1$ to N_E do:
2. Reset simulator, get initial state $\mathbf{y}^1, \{\mathbf{y}_{d,j}^1\}_M$, and $\hat{\mathbf{z}}_d^1 = \mathbf{z}_d^1$.
3. Stage one: collecting experience.
 For $t = 1$ to T do:
3.1. Decide action x^t based on policy π , and execute x^t .
3.2. Decide actions $\{x_{d,j}^t\}_M$ based on soft policy π_d , and execute $\{x_{d,j}^t\}_M$.
3.3. Based on $\{\mathbf{y}_{d,j}^1\}_M$ and $\hat{\mathbf{z}}_d^t$, compute $Q_d(\cdot | \cdot)$ for the simulator.
3.4. Observe next state \mathbf{y}^{t+1} and $\{\mathbf{y}_{d,j}^{t+1}\}_M$.
3.5. Calculate $\hat{\mathbf{z}}_d^{t+1}$ based on $Q_d(\cdot | \cdot)$ and $\hat{\mathbf{z}}_d^t$.
3.6. Observe reward $r^t(\mathbf{y}^t, \hat{\mathbf{z}}_d^t, \hat{\mathbf{z}}_d^{t-1}), \{r_d^t(\mathbf{y}^t, \mathbf{y}_{d,j}^t, \mathbf{y}_{d,j}^{t-1})\}_M$.
3.7. Store transitions $(\mathbf{y}^t, x^t, \mathbf{y}^{t+1}, r^t)$ and $\left\{ \left(\mathbf{y}_{d,j}^t, x_{d,j}^t, \mathbf{y}_{d,j}^{t+1}, r_d^t(\cdot) \right) \right\}_M$ to D .
 End for.
4. Stage two: learning the experiences.
 For $n_s = 1$ to N_s do:
4.1. Sample a mini-batch of transitions $(\mathbf{y}^t, x^t, \mathbf{y}^{t+1}, r^t(\cdot))$ from D .
4.2. Update the platform's value networks by minimizing $L(\theta_v)$.
4.3. Update the platform's policy networks as $\theta_p \leftarrow \theta_p + \delta \nabla_{\theta_p} G(\theta_p)$.
4.4. Sample a mini-batch of transitions $\left\{ \left(\mathbf{y}_{d,j}^t, x_{d,j}^t, \mathbf{y}_{d,j}^{t+1}, r_d^t(\cdot) \right) \right\}_M$ from D .

4.5. Update the representative driver’s value network by minimizing $L(\theta_v)$.
 4.6. Update the representative driver’s policy network as $\theta_p \leftarrow \theta_p + \delta \nabla_{\theta_p} G(\theta_p)$.

End for.

End for.

5. Finish.

1

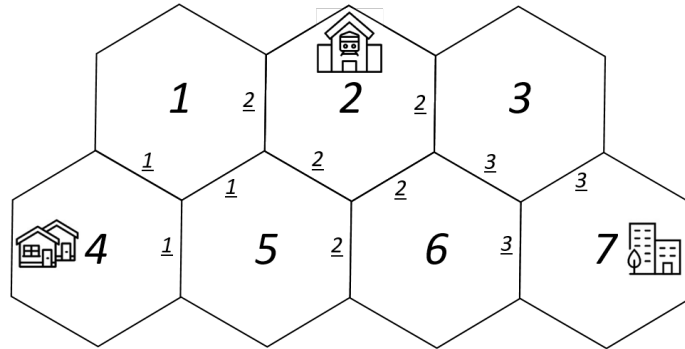
2 5. NUMERICAL STUDY

3 In this section, we conduct a set of numerical experiments with the MF-MDP model and
 4 algorithms developed in Section 4. We show (1) the computation time, converging speed for
 5 learning policies, and converged total rewards for the agents of the RAC algorithm compared
 6 with the benchmark MAC algorithm; (2) the impact of spatial-temporal subsidies on drivers’
 7 relocation strategies; and (3) the platform’s different spatial-temporal subsidy strategies that
 8 balance the trade-offs between net revenue and service rate.

9 5.1 SCENARIO SETTINGS

10 The zone network of the ride-sourcing market is illustrated in Figure 3. Zone IDs are
 11 shown in the center of the hexagons, and travel times (number of time periods) between adjacent
 12 zones are shown via underlined numbers near edges. For instance, a driver needs 2 time periods
 13 to travel between zone 1 and zone 2. In this small town with 7 zones, we assume that zone 4 is
 14 a residential area, zone 7 is a business area, and zone 2 has a railway station. Due to the huge
 15 computational costs, numerical studies with a small network were usually adopted in
 16 reinforcement learning-related studies. For example, [Mao et al. \(2020\)](#) divide Manhattan into 8
 17 zones and examine drivers’ optimal repositioning among these zones. [Braverman et al. \(2019\)](#)
 18 use a nine-region network with parameters calibrated by DiDi data to evaluate their proposed
 19 empty-car routing policy. Moreover, by using a small network, we can observe clear patterns of
 20 drivers’ sequential actions and better understand how self-relocation is affected by subsidies.

21



22

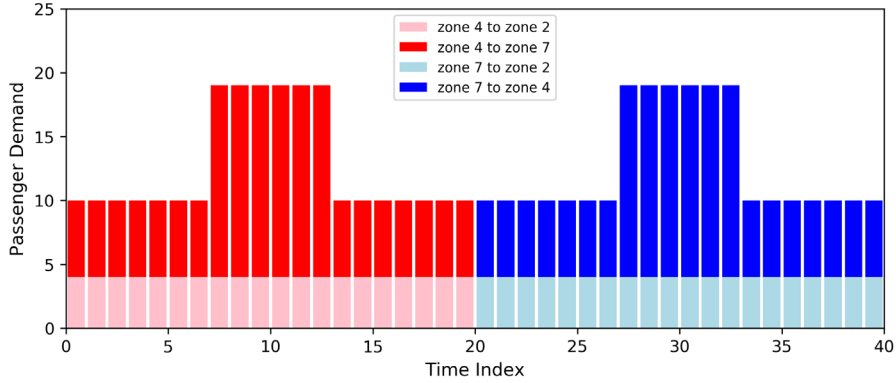
23

Figure 3. Network of the numerical study.

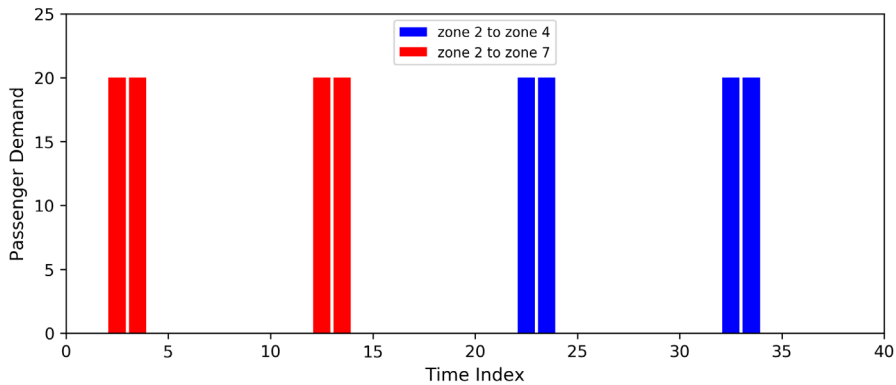
24 We consider a general ride-sourcing market that consists of both tidal and periodic
 25 passenger demand. The pattern of deterministic passenger demand is shown in Figure 4¹¹. We
 26 consider a total of 40 time periods in the operational horizon, and each period represents 5
 27 minutes. First, there is tidal demand between zone 4 and zone 7 (see Figure 4(a)), such that
 28 passengers go from the residential area to the business center at time periods 1–20 (red bars) and

¹¹ Note that the ticks on the horizontal axis, i.e., time index, refer to “time point,” while demand is generated during time periods. Therefore, for instance, time period 15 refers to the period between time index 15 and time index 16. The same representations are adopted in Figure 6.

1 return at time periods 21–40 (blue bars). Each tidal demand has a peak period—i.e., during time
 2 periods 8–13 and 28–33. Second, also in Figure 4(a), we assume that since some passengers
 3 live in zone 4 but work outside the small town, there is ride-sourcing demand from zone 4 to
 4 zone 2 at time periods 1–20 (light pink bars); and since some passengers work in zone 7 but live
 5 out of town, there is ride-sourcing demand from zone 7 to zone 2 at time periods 21–40 (light
 6 blue bars). Third, in Figure 4(b), for every 10 time periods (50 minutes), a train arrives at zone
 7 2, and passengers from the train either go to work (red bars before time index 20) or home (blue
 8 bars after time index 20). We use this demand setting because it reflects general scenarios with
 9 both demand-hot areas, demand-cold areas, tide demand, and periodic demand. Based on the
 10 edge-based matching rule, idle drivers in zones 1, 3, 5, 6 can also get passenger orders. However,
 11 the matching probability at zones 1, 3, 5, 6 would be much lower than that at zones 2, 4, and 7.
 12 In this case, there are trade-offs in the market with competing drivers: a driver can idly cruise to
 13 zone 4 to ensure a high matching probability; alternatively, the driver can stay in zones 1 or 5
 14 such that he/she has a low probability of getting an order from zone 4 or zone 2 (at time periods
 15 when a train arrives).



(a) Tidal demand in zones 2, 4 and 7



(b) Periodic demand in zone 2

16 Figure 4. Passenger demand.

17 Other exogenous parameters are set as follows: commission rate $\eta = 0.20$, trip fare rate
 18 $\alpha = 10$ CNY per time period, and discount factor $\rho = 0.8$. Based on the numerical settings and
 19 definitions of states and actions, the cardinality of the platform's state-action set is 2,560; the

1 cardinality of the state-action set for a driver to take actions (i.e., in a purely idle state) is $1,160$ ¹².
 2 With multiple drivers, since agents take actions independently, the cardinality of the solution
 3 space (containing the state-action sets for all the agents) can be as large as $2,560 \times 1,160^M$,
 4 which makes it difficult to solve via the MAC algorithm. In contrast, for the RAC algorithm, the
 5 cardinality of the solution space reduces to $2,560 \times 1,160$.

6 The hyperparameters of the algorithms for all subsequent numerical studies are as
 7 follows. In the RAC algorithm, for the platform agent, we establish a simple three-layer fully
 8 connected network with 24 neurons in the hidden layer for the value network and a three-layer
 9 fully connected network with 24 neurons in the hidden layer for the policy network. Similarly,
 10 for the representative driver agent, we use a three-layer fully connected network with 24 neurons
 11 in the hidden layer for both the value and policy networks. The activations of all hidden units are
 12 ReLu, while output layers of the value function approximation networks and policy networks
 13 use Linear and Softmax activations, respectively. The same policy and value network structures
 14 are used for each driver agent and the platform agent in the MAC algorithm. For both algorithms,
 15 the learning rate of the policy network is set at 0.001, and the learning rate of the value network
 16 is set at 0.01.

17 We consider two experiments. First, the platform implements a non-subsidy strategy (i.e.,
 18 $\beta = 0$). We use different numbers of drivers (i.e., $M = 1, 10, 50, \text{ and } 100$) in the market to
 19 evaluate the performance (in terms of achieved total rewards) and efficiency (in terms of
 20 computation time) of the MAC and RAC algorithms. Second, we assume there are 100 drivers
 21 serving in the market and the platform must design the spatial-temporal subsidy strategy to
 22 maximize its total objective value, which is a weighted sum of net revenue and service rate (see
 23 Eq. (A.10)). A range of subsidy rates (i.e., $\beta = 0, 2, 4, 6, \text{ and } 8$ CNY per ride) and weights of
 24 the service rate (i.e., $\mu = 0, 20,000, \text{ and } 40,000$ CNY) are tested¹³. This is to investigate how
 25 the spatial-temporal subsidy affects the self-relocation strategies of drivers, as well as the supply-
 26 demand situation, and examine how the platform’s subsidy strategy varies with different weights
 27 for service rate. The execution programming codes for the two experiments are the same except
 28 for the settings of number of drivers and subsidy rates. Therefore, the specific amount of
 29 subsidies have no impact on the computational time and the results in Section 5.2 well support
 30 the performance of the proposed algorithms.

31 5.2 PERFORMANCE OF THE REPRESENTATIVE-AGENT ALGORITHM

32 In the first experiment, we test the computation time for the RAC and MAC algorithms
 33 for drivers to pursue high-rewarding self-relocation strategies without subsidies. Simulation of
 34 the environment (i.e., calculating the matching probabilities and sampling the matchings,
 35 rewards, and transitions) and reinforcement learning algorithms are conducted on an HP Z4G4
 36 workstation with 12 Inter I7-7800 processors and four 16-GB rams.

¹² For the platform, since three zones (2, 4, and 7) have passenger demand, we can ignore zones without demand in set \mathcal{S} ; therefore, there are 2^3 possible states and 2^3 possible actions, and the cardinality of the state-action set is $2^3 \times 2^3 \times 40 = 2560$. For a driver, we consider purely idle states such that a driver must take a relocation action. If the driver is in zone 4 or 7, then they have 3 relocation destinations; once they are in zone 1 or 3, then they have 4 relocation destinations; and if the driver is in zone 2, 5, or 6, they have 5 relocation destinations. Therefore, the cardinality of a driver’s state-action set is $(3 \times 2 + 4 \times 2 + 5 \times 3) \times 40 = 1,160$.

¹³ In both experiments, we assume that drivers are purely idle and uniformly distributed in zones at the beginning of the simulation. The platform would not allow a high subsidy rate that causes low net revenue for a single ride. Since the platform’s net revenue for an order from zone 4 to zone 7 is $0.2 \times 10 \times 6 = 12$ CNY ($\alpha = 10$ CNY per time period and $\eta = 0.2$), the maximal subsidy rate is set as 8 CNY per ride and the net revenue after offering a subsidy is $12 - 8 = 4$ CNY per ride.

In Table 1, we present the computation time for one learning epoch, which consists of simulating the order matches, actions, and states in the environment; storing agents' transitions; and updating the knowledge and value networks of agents (i.e., steps 3 and 4 in Algorithms 1 and 2). We note that the RAC algorithm is slower than the MAC algorithm when $M = 1$; this is because the RAC algorithm must derive a comprehensive soft policy that covers all of the state-action sets of the representative driver; in contrast, the MAC algorithm updates 1 driver's policy based on their own experienced states and actions, ignoring policies that are conditional on unvisited states. As M increases, the time for computing the MF state and updating each agent's policy and value networks will get longer, so that the computation time notably increases for the MAC algorithm. By contrast, the RAC algorithm only computes an approximated MF state and updates the policy and value networks for the representative driver, and the computation time gets much longer as M increases. As a result, with $M = 50$ or 100 , we note that the RAC algorithm is significantly faster than the MAC algorithm.

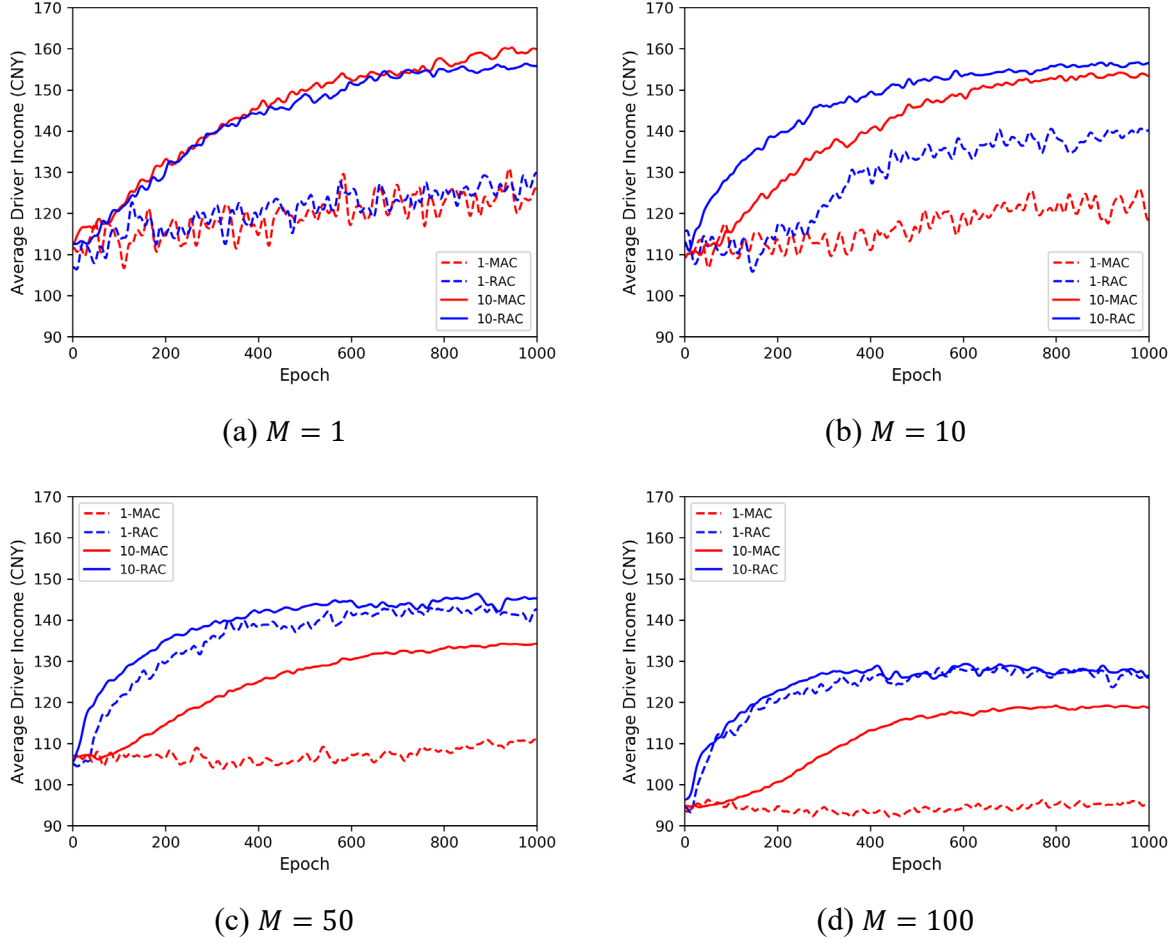
Table 1. Computation Time for One Learning Epoch (seconds)

	$M = 1$	$M = 10$	$M = 50$	$M = 100$
MAC	1.98	27.19	347.3	1654
RAC	2.59	8.17	35.5	85.9

The performance of the two algorithms can be measured by the increase in average driver income (i.e., the average value of the total income for M drivers) over learning epochs. In the experiment, the number of total epochs is 1,000; within each epoch, we conduct either 1 simulation or 10 simulations (i.e., for each n_e , to repeat step 3 in Algorithms 1 and 2 for 10 times before going to step 4, such that more samples of transitions can be generated) to update the policy and value networks. Although global optimality is not guaranteed with reinforcement learning algorithms, the 10-simulation case provides much faster convergence and higher total rewards than the 1-simulation case. The disadvantage is that the computation time for each epoch will be longer as the number of simulations increases.

We illustrate the performance of the two algorithms with different numbers of drivers in Figure 5. When $M = 1$, the MAC algorithm with 10 simulations (referred to as 10-MAC) results in higher average driver income than the RAC algorithm with 10 simulations (10-RAC; see Figure 5(a)). This reflects the ineffectiveness of a soft policy in the simplified MF-MDP model when the number of drivers is small. In addition, average driver income grows slowly with the 1-MAC and 1-RAC algorithms (i.e., by conducting 1 simulation within each epoch), because it is difficult for the transitions observed in 1 simulation to cover a large state-action set of the driver (see Figure 5(a)). When $M = 10$, the RAC algorithm begins to demonstrate its advantages. In Figure 5(b), we see that the 10-RAC algorithm converges much faster and leads to higher average driver income than the 10-MAC algorithm. Unlike the $M = 1$ scenario, the 1-RAC algorithm under $M = 10$ can also achieve high average driver income because more transitions are sampled and used to update the policy and value networks of the representative driver. As M increases to 50, on one hand, the 10-MAC algorithm encounters its bottleneck (at around 132 CNY per driver) and can barely increase average income further (see Figure 5(c)). On the other hand, the 10-RAC algorithm still outperforms the other algorithms, but its gaps in both convergence speed and final average income from the 1-RAC algorithm become smaller. This is because with a large M , 1 simulation during an epoch can generate sufficient samples of transitions. Finally, in Figure 5(d) under $M = 100$, the 1-MAC algorithm never improves average driver income due to insufficient transitions sampled over the large solution space, and we still observe the bottleneck for the 10-MAC algorithm (at around 128 CNY per driver). Also, the advantage of the 10-RAC algorithm over the 1-RAC algorithm becomes negligible. Based

1 on these findings, we conclude that with a large number of drivers, the RAC algorithm is capable
 2 of identifying policies to further improve average total rewards compared with the MAC
 3 algorithm on a small-scale network. In addition, a small number of simulations within one epoch
 4 is sufficient for the RAC algorithm to quickly converge to a policy that leads to high average
 5 total rewards. We aim to examine the MAC and RAC algorithms on real-world networks in near
 6 future studies.



7 Figure 5. Performance of the MAC and RAC algorithms

8 5.3 SPATIAL-TEMPORAL SUBSIDIES AND DRIVERS' SELF-RELOCATION

9 In the second experiment, we retain the ride-sourcing market that contains 100 drivers
 10 and let the platform to pursue rewardable spatial-temporal subsidy strategies with some
 11 predefined β . As described in Section 5.1, we assume the platform assigns weights to the service
 12 rate in the objective and adopts different subsidy rates. Namely, a zero weight (e.g., $\mu = 0$) for
 13 the service rate indicates that the platform is only concerned with net revenue; a medium weight
 14 (e.g., $\mu = 20,000$) implies a balance between net revenue and the service rate; and a large weight
 15 (e.g., $\mu = 40,000$) indicates that the platform mainly focuses on the service rate in the objective.
 16 As stated previously, the subsidy rates β range from 0 to 8 CNY per ride in steps of 2 CNY per
 17 ride. For each combination of parameters (i.e., μ and β), we pursue the platform's optimal
 18 subsidy strategy with respect to the total objective value in terms of drivers' self-relocation. We
 19 use the RAC algorithm to solve the simplified MF-MDP model for different combinations of μ
 20 and β . To balance computational cost and performance (i.e., for both the platform and the

1 representative driver), the number of epochs is set at 1,000 and 2 simulations are conducted
 2 within each epoch (i.e., a 2-RAC algorithm is adopted).

3 Table 2. Results with Spatial-temporal Subsidies for the Entire Horizon

4 (a) Total objective value (weighted sum of net revenue and service rate) for the platform

	$\beta = 0$	$\beta = 2$	$\beta = 4$	$\beta = 6$	$\beta = 8$
$\mu = 0$	3,189	3,186	3,179	3,190	3,194
$\mu = 20,000$	13,616	13,735	13,859	13,930	13,689
$\mu = 40,000$	24,112	24,490	24,677	24,830	24,964

5 (b) Average driver income (CNY per driver)

	$\beta = 0$	$\beta = 2$	$\beta = 4$	$\beta = 6$	$\beta = 8$
$\mu = 0$	127.55	127.49	127.35	127.72	127.87
$\mu = 20,000$	127.84	131.61	134.27	134.64	128.80
$\mu = 40,000$	127.68	132.43	136.55	136.16	136.10

6 (c) Total net revenue for the platform (CNY)

	$\beta = 0$	$\beta = 2$	$\beta = 4$	$\beta = 6$	$\beta = 8$
$\mu = 0$	3,189	3,186	3,179	3,190	3,194
$\mu = 20,000$	3,196	3,095	3,159	3,110	3,189
$\mu = 40,000$	3,192	3,050	3,077	2,950	2,804

7 (d) Service rate in the market

	$\beta = 0$	$\beta = 2$	$\beta = 4$	$\beta = 6$	$\beta = 8$
$\mu = 0$	51.5%	51.3%	50.8%	51.1%	51.3%
$\mu = 20,000$	52.1%	53.2%	53.5%	54.1%	52.5%
$\mu = 40,000$	52.3%	53.6%	54.0%	54.7%	55.4%

8 (e) Total subsidies offered by the platform (i.e., earned by all drivers) (CNY)

	$\beta = 0$	$\beta = 2$	$\beta = 4$	$\beta = 6$	$\beta = 8$
$\mu = 0$	0.0	0.6	4.2	2.4	3.8
$\mu = 20,000$	0.0	117.7	158.4	210.0	29.6
$\mu = 40,000$	0.0	208.3	269.2	356.2	479.2

9 (f) Average subsidies per matched order (CNY)

	$\beta = 0$	$\beta = 2$	$\beta = 4$	$\beta = 6$	$\beta = 8$
$\mu = 0$	0.00	0.00	0.01	0.01	0.01
$\mu = 20,000$	0.00	0.33	0.44	0.58	0.08
$\mu = 40,000$	0.00	0.58	0.75	0.97	1.29

10 (g) Average number of driver delivery/pickup time periods

	$\beta = 0$	$\beta = 2$	$\beta = 4$	$\beta = 6$	$\beta = 8$
$\mu = 0$	15.9/3.0	15.9/3.0	15.9/2.9	16.0/2.9	16.0/2.9
$\mu = 20,000$	16.0/3.0	16.3/3.0	16.6/3.1	16.6/3.1	16.1/3.0
$\mu = 40,000$	16.0/2.9	16.3/3.1	16.7/3.2	16.6/3.2	16.4/3.1

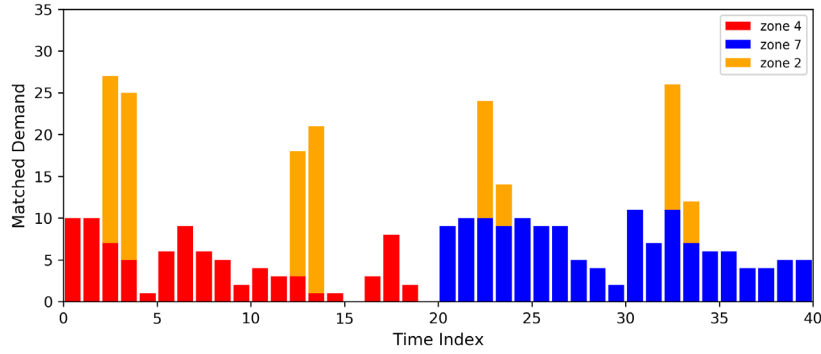
1
2 Results of this experience are illustrated in Table 2. From the 1,000 epochs with fixed μ
3 and β in the 2-RAC algorithm, we select 5 epochs with the highest platform objective values
4 and summarize the average metrics¹⁴. The relation between β and the metrics with a fixed μ can
5 be obtained; the total subsidies offered (i.e., Table 2(e)) reflect the platform’s subsidy strategy.

6 When $\mu = 0$, the total amount of subsidies offered to drivers under different values of β
7 is small (i.e., Table 2(e)); this indicates that the platform prefers not to provide subsidies, and
8 drivers pursue high-rewarding self-relocation strategies without subsidies. As a result, the total
9 objective value for the platform, total net revenue for the platform, average driver income, and
10 service rate (i.e., Tables 2(a)–2(d)) across different values of β are more or less the same. Note
11 that the reinforcement learning algorithm cannot guarantee the global optimal, and the small
12 differences between the values are mainly due to simulation noise. The “non-subsidy” result
13 under $\mu = 0$ is foreseeable for two reasons: (1) the commission rate is low, so that the platform,
14 which only retains a positive but small net revenue for an order, might not afford a large subsidy
15 rate; and (2) a small subsidy rate might not motivate enough drivers to change their self-
16 relocation strategies in order to gain sufficient benefits from alleviating supply-demand
17 imbalance. Therefore, the increased commission withheld by the platform cannot cover the
18 subsidies offered to drivers, which leads to a loss in net revenue while implementing subsidy
19 strategies.

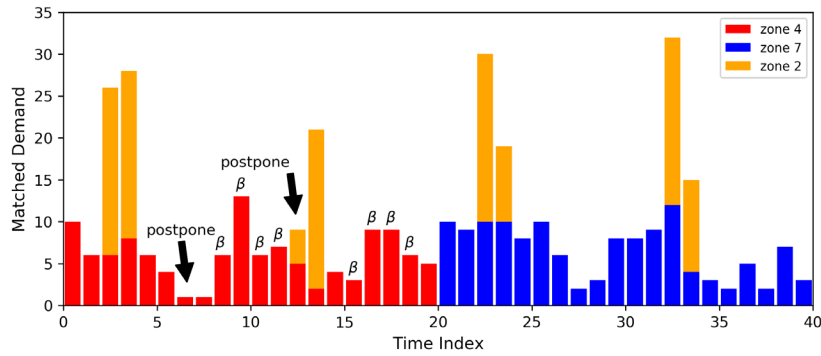
20 If the platform has a balanced weight with $\mu = 20,000$, its total objective value first
21 increases and then decreases with β (see Table 2(a)). Based on Table 2(d) and Table 2(e), we
22 note that as β increases from 0 CNY to 6 CNY per ride, the increased subsidy improves the
23 service rate; however, once $\beta = 8$ CNY per ride, the subsidy becomes cost-ineffective because
24 the benefits gained from the enhanced weighted service rate cannot offset the revenue loss caused
25 by subsidy provisions, and thus the platform is inclined to adopt a non-subsidy strategy.

26 If the platform mainly prefers a high service rate (i.e., $\mu = 40,000$), a subsidy less than
27 or equal to 8 CNY per ride is always cost-effective for improving the total objective value by
28 reshaping drivers’ self-relocation strategies and increasing the service rate (see Table 2(a), Table
29 2(d), and Table 2(e)).

¹⁴ Table 2(a), 2(c) and 2(d) shows the total reward, net revenue, and service rate of the platform, respectively, 2(b) the average income (net revenue) of drivers (i.e., total income divided by 100 drivers), 2(g) the average time periods in delivery/pickup task for drivers (i.e., total number of delivery/pickup times divided by 100 drivers), 2(e) and 2(f) are subsidy related metrics. A value of 0.6 in Table 2(e) under $\mu = 0$ and $\beta = 2$ is obtained. This is because for the 5 epochs (i.e., a total of 10 simulations) with the highest platform’s objective values, the total amount of subsidies offered by the platform to all the drivers in the entire horizon is 6 CNY; dividing 6 by 10 simulations, we get 0.6.



(a) Baseline scenario



(b) Subsidy scenario

1 Figure 6. Served Rides across the Time Horizon and Zones

2 We refer to $\mu = 0$ and $\beta = 0$ as the baseline scenario, in which the platform's objective

3 consists of only net revenue, and refer to $\mu = 40,000$ CNY and $\beta = 8$ CNY per ride as the

4 subsidy scenario in which the platform focuses more on the service rate. Comparing the subsidy

5 scenario with the baseline scenario, spatial-temporal subsidies can lead to a 6.7% and 7.6%

6 increase in the average driver income and the service rate, respectively (see Table 2(b) and Table

7 2(d)). We show the spatial-temporal number of passengers served (i.e., matched demand) in

8 Figure 6. Time periods with subsidies are denoted using β in Figure 6(b): The platform provides

9 subsidies at zone 4 at time periods 9–12 and 16–19. Motivated by the subsidy, some drivers

10 “postpone” service by idle cruising before the target (i.e., demand-hot) zone is subsidized but

11 relocating to (and thus arriving at) the target zone in time periods with subsidies. As denoted by

12 “postpone” in Figure 6(b), fewer passengers are served at time periods 7–8 and 13 due to the

13 idle cruising and postponing phenomenon; instead, more passengers are served during time

14 periods with subsidies. There are 148 passengers served during time periods 1–15 under both

15 scenarios. In contrast, due to the postponement of services, the number of passengers served after

16 time index 15 notably increases from 195 to 223. These results imply that a platform with an

17 emphasis on service rate has the foresight to mitigate the imbalance between driver supply and

18 passenger demand.

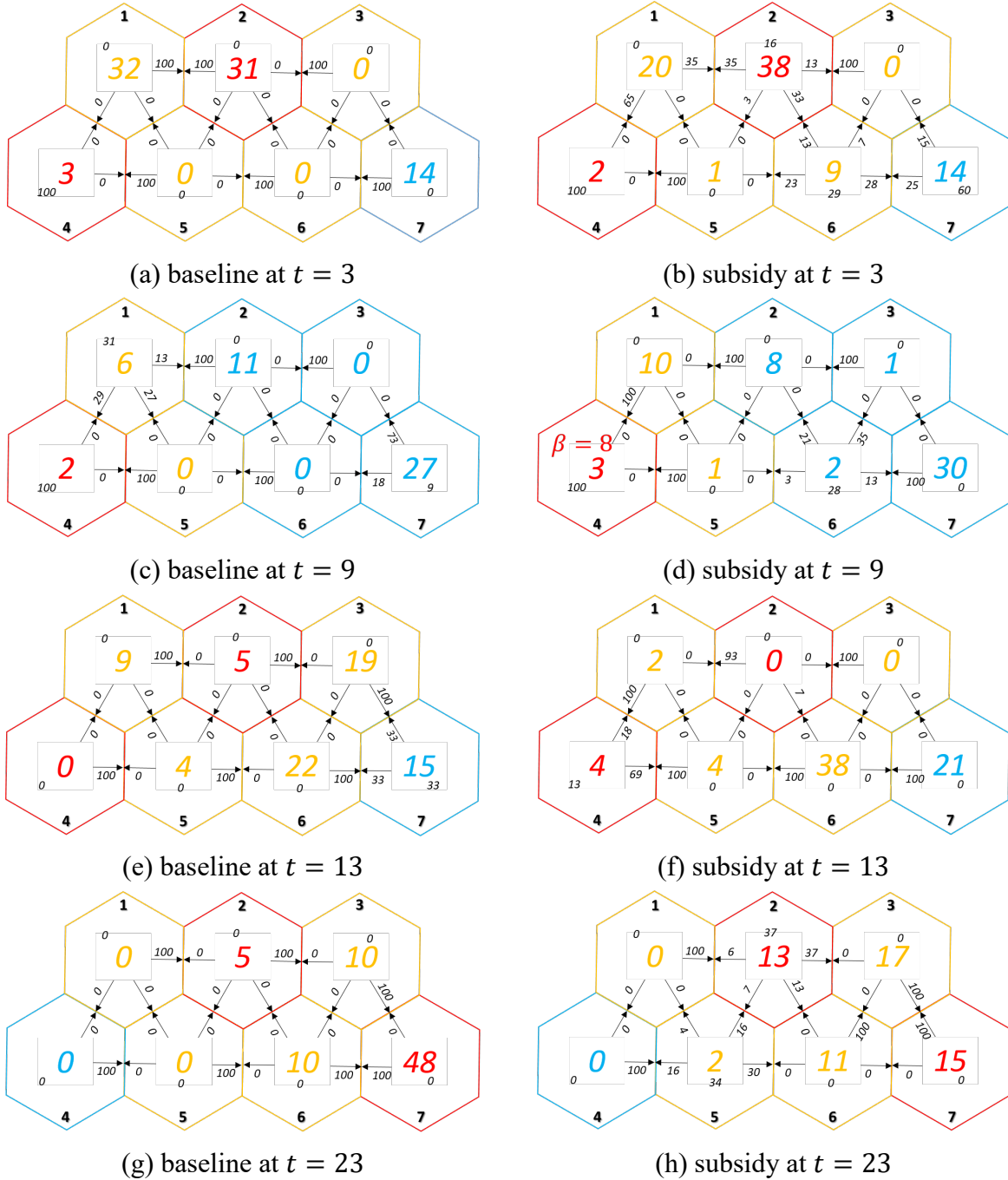


Figure 7. Spatial-temporal idle supply and self-relocation policies.

1
2 To better understand how drivers' self-relocation strategies are affected in the subsidy
3 scenario, we provide the spatial-temporal idle driver supply and soft self-relocation policies in
4 the simplified MF-MDP model in Figure 7. Black numbers at the top/bottom of the zones
5 represent zone IDs; colored numbers at the center of each zone denote the numbers of idle drivers;
6 and arrows with small underlined numbers denote relocation destinations and corresponding
7 proportions in percentage (i.e., the soft policy). Blue, yellow, and red represent zones with a low
8 matching probability (in demand-cold zones), a medium matching probability (in zones adjacent
9 to demand-hot zones), and a high matching probability (in demand-hot zones), respectively. Note
10 that the instances reported in Figure 7 are from a single simulation with the highest objective
11 value for the platform. At $t = 3$, idle drivers in the baseline scenario move to either zone 2 or
12 zone 4 to pick up passengers (see Figure 7(a)); in the subsidy scenario, some drivers in zones 6

1 and 7 have diverse relocation directions (e.g., drivers in zone 7 have a 60% chance of staying).
2 One reason might be that they first cruise around and wait, then try to arrive at zone 4 at $t \in$
3 $\{9,10,11,12\}$ to earn subsidies. Because of this phenomenon, at $t = 9$, the number of idle drivers
4 adjacent to zone 4 in the subsidy scenario is notably higher than in the baseline scenario. The
5 postponing phenomenon also happens at $t = 13$, when some drivers perceive the upcoming
6 subsidies at $t \in \{16,17,18,19\}$ and decide not to immediately serve passengers in zone 2 (see
7 Figure 7(f) and Figure 7(b)). The benefits of the postponing phenomenon can be partially
8 observed in Figures 7(g)–7(h). In the baseline scenario, drivers in the first half of the operational
9 horizon keep relocating to zones 2 and 4 to serve passengers. This results in a shortage of supply
10 during time periods 9–11 and 14–17 (see Figure 6(a)), as well as peak arrivals of idle drivers at
11 zone 7 at an early time in the second half of the operational horizon (e.g., time period 23 in
12 Figure 7(g)). In contrast, in the subsidy scenario, some drivers are inclined to idly cruise and
13 postpone their services so that the supply becomes smooth across zones and time periods,
14 especially during the second half of the operational horizon. Consequently, the supply-demand
15 imbalance in this case study is alleviated due to the implementation of spatial-temporal subsidies.

17 **6 DISCUSSION OF SUBSIDY SCHEMES**

18 The numerical studies in Section 5 offer in-depth insights for ride-sourcing platforms
19 about the effectiveness of a uniform subsidy scheme in addressing supply-demand imbalance. In
20 such a scheme, the platform predetermines the amount of subsidy per order (i.e., subsidy rate)
21 and offers this amount of subsidy to drivers once they are matched with passengers whose origins
22 are the subsidized zones. The strategy of spatial-temporal subsidies under the uniform scheme is
23 largely affected by the objective of the platform. If the platform only cares about the immediate
24 net revenue, the effectiveness of this subsidy scheme could be limited. This is because for a
25 single order, the subsidy rate generally does not exceed the commission withheld by the platform
26 (otherwise the platform earns negative net revenue for an order). Thus, the amount of subsidy
27 offered to a driver is much smaller than they earn from the trip fare and is unattractive to drivers,
28 who may not be motivated to move to the designated zones. In this case, offering subsidies could
29 cause a loss in net revenue because the subsidy provision is higher than the commission gain;
30 therefore, the platform would prefer a non-subsidy strategy. By contrast, if the platform pursues
31 a high service rate (i.e., number of passengers served), it would like to offer a subsidy sufficient
32 to stimulate drivers to demand-hot zones despite the reduction in immediate revenue. The latter
33 case might occur when a platform expands its business and competes with others. An example
34 is the price war between DiDi and Uber in mainland China in 2016 before they consolidated.

35 However, the uniform subsidy scheme is not superior in improving the service rate and
36 net revenue of the platform simultaneously. This is because the platform provides the same
37 amount of subsidies to different drivers: (1) those who already have desirable self-relocation
38 strategies such that they could relocate to demand-hot zones even without subsidies, and (2)
39 those who are incentivized by subsidies but would not relocate to demand-hot zones if no
40 subsidies were provided. Therefore, subsidies offered to drivers belonging to the first type would
41 not generate more net revenue for the platform, and only subsidies offered to drivers in the
42 second type would help mitigate the supply-demand imbalance and improve the service rate.

43 To thoroughly examine the designs of spatial-temporal subsidies for drivers with certain
44 self-relocation strategies, a number of tasks must be relayed left to the future studies. Different
45 subsidy schemes must be examined and evaluated based on the ride-sourcing MF-MDP model.
46 Below, we provide a few sample schemes with the potential advantages and feasibility:

- 1 • Surge subsidy (or zone-based) scheme, in which the platform provides higher subsidies
2 at zones with greater supply-demand imbalance. Once there is super large passenger
3 demand in a hot area, the platform could offer irresistible subsidies to drivers who
4 relocate to and then serve passengers in the area. The revenue loss due to high subsidy
5 provision could be offset by the improvement in the service rate, since sufficient drivers
6 will be attracted to the hot area to accommodate the high service needs.
- 7 • Distance-based subsidy scheme, in which the subsidy rate is proportional to the travel
8 distance of the ride order or a subsidy is applied only if the travel distance of the order
9 exceeds some threshold. Such a scheme could be beneficial once there is an insufficient
10 supply of long-distance passenger demand.
- 11 • Origin-destination-based subsidy scheme, in which the platform offers heterogeneous
12 subsidies based on the origin and destination of the ride order. For instance, the platform
13 could provide a high subsidy for trips that originate at a demand-cold area and terminate
14 at a demand-hot area. Consequently, the supply-demand imbalance could be improved
15 as the overall driver supply at demand-cold areas is incentivized to relocate to demand-
16 hot areas. We can employ the MF-MDP model to determine the critical rules with respect
17 to subsidy rates and the characteristics of origins/destinations.
- 18 • Performance-based (or driver-based) subsidy scheme, in which the platform offers
19 subsidies according to drivers' performance and behaviors. For instance, the platform
20 could only offer subsidies to drivers who would not relocate to demand-hot zones without
21 incentives. Although this subsidy scheme can reduce the subsidy provided to drivers with
22 high-rewarding self-relocation strategies, it could be controversial due to potential
23 discrimination concerns.

24 All of these subsidy schemes merit analysis using the MF-MDP model to gain
25 comprehensive insights into the pros and cons of diverse spatial-temporal subsidies in ride-
26 sourcing markets. Furthermore, we would examine spatial-temporal subsidies in more realistic
27 scenarios in terms of a large-scale zone network, passenger demand derived from actual data,
28 and a flexible setting of subsidy levels. Real-world public datasets can be used to generate large-
29 scale ride-sourcing scenarios. Although the current edge-based matching rules (Eqs. A.(6)–A.(9))
30 analytically capture the cross-zone matching feature of ride-sourcing markets, it can be
31 computationally inefficient for large-scale analyses when calculating the joint probability
32 distribution of matching results. We aim to improve the efficiency of the edge-based matching
33 rule for real-world scenarios in future studies. In addition, the predefined subsidy rates (i.e.,
34 either 0 or β) in this paper could underestimate the effectiveness of a subsidy due to the
35 inflexibility of implementing heterogeneous subsidies at different locations and with different
36 traveling distances. Instead, we can predefine a few subsidy levels and apply reinforcement
37 learning (e.g., the RAC) algorithms to pursue the optimal subsidy level in each time period to
38 maximize the total rewards.

40 7 CONCLUSIONS

41 In this paper, we propose a generalized MF-MDP model to capture sequential and
42 interactive decision processes in a ride-sourcing environment with the platform as the major
43 agent and multiple drivers as minor agents. The MF-MDP model is particularly suitable for
44 research questions in which the major agent (platform) and minor agents (drivers) have distinct
45 objectives. The decisions/actions of the platform can directly affect the drivers' states, while the
46 drivers' actions can influence the platform's state and drivers' average state, which is referred to
47 as the MF state. An approximation of the MF state is employed to simplify the model, such that
48 we only need to optimize the policies for the platform and one representative driver instead of

1 the policies for the platform and all individual drivers (as in the standard MF-MDP model).
2 Consequently, computational complexity can be notably reduced when there are a large number
3 of drivers in the environment.

4 In particular, we adopt the MF-MDP model to design the platform’s spatial-temporal
5 subsidy strategies with a predefined subsidy rate for drivers who have self-relocation strategies.
6 A representative-agent reinforcement learning algorithm is proposed to solve the MF-MDP
7 model. Using numerical studies, we demonstrate that due to the significant reduction of the
8 number of agents and solution space, the representative-agent algorithm demonstrates significant
9 computational advantages and fast convergence and achieves higher rewards, compared with the
10 conventional multi-agent algorithm. In addition, we investigate the potential impact of spatial-
11 temporal subsidies on drivers’ self-relocation strategies and the resulting platform’s objective
12 values and drivers’ income. Based on a uniform subsidy scheme, our results suggest that
13 subsidies can improve the service level (number of passengers served) by incentivizing idle
14 drivers to locations with overfull passenger demand and insufficient driver supply. On one hand,
15 if the platform only pursues net revenue (measured commission withheld from trip fares by the
16 platform minus the amount of subsidies offered to drivers), a subsidy strategy with predefined
17 subsidy levels is cost-ineffective due to a large reduction in net revenue from a single order
18 versus a small increase in the number of passengers served. On the other hand, when the platform
19 pays more attention to the number of passengers served (in order to improve the customer
20 satisfaction rate), it is more willing to offer sufficient subsidies to stimulate drivers to demand-
21 hot zones and achieve a better supply-demand balance. In this case, the spatial-temporal subsidy
22 strategy leads to a win-win situation in which both average driver income and the platform’s
23 total objective value are notably improved.

24 This paper makes three major contributions to the literature. First, unlike previous ride-
25 sourcing MDP models that assume the platform has full control of drivers, the proposed MF-
26 MDP model considers interactive decision processes between the platform and drivers. Second,
27 the simplified MF-MDP model with the representative-agent algorithm is shown to be a good
28 alternative to multi-agent MDP models, which are generally computationally expensive. The
29 mean-field approximation not only saves computational resources but also achieves higher
30 rewards with a faster convergence in our research problem. This is mainly due to the special
31 characteristic of a ride-sourcing market, in the sense that the number of drivers is so large that
32 the platform can consider the mean-field state of all drivers without tracking the individual state
33 of each driver. Third, our model describes the interactions between the platform’s spatial-
34 temporal subsidy strategies and drivers’ self-relocation strategies; in contrast, most previous
35 studies investigate either the impact of platform-side incentives or idle-vehicle relocation
36 problems separately without considering their interplay.

37 There are several important directions for future research. First, some deep learning-
38 based algorithms and MF simulation approaches can be developed to further enhance
39 performance and reduce computational complexity. We are particularly interested in developing
40 edge-based matching rules that are both capable of depicting cross-zone matching processes and
41 powerful for large-scale multi-agent problems in practically relevant scenarios. Second, based
42 on the generalized ride-sourcing MF-MDP model, we will examine the impacts of other subsidy
43 schemes, such as surge subsidy schemes over time and zones, distance-based subsidy schemes,
44 and origin-destination-based subsidy schemes. These subsidy schemes are expected to mitigate
45 supply-demand imbalance more efficiently than the uniform subsidy scheme that offers the same
46 amount of subsidy to drivers upon matches with passengers from subsidized regions. Third, the
47 framework can be extended to investigate ride-sourcing markets coupled with public transit
48 services, and identify optimal coordination between ride-sourcing drivers who aim to improve
49 their earnings by self-relocation and public transit operators who attempt to design transit

1 schedules to improve transit usage. For example, the platform’s knowledge of bus services’
2 timeline could incentivize drivers to relocate to transit stations at the appropriate time, as a result
3 of which the cooperation and substituting effect between ride-sourcing and public transit services
4 could be enhanced.

6 ACKNOWLEDGEMENT

7 This work is partially supported by the Hong Kong Research Grants Council under
8 projects HKUST16208920 and NHKUST627/18, and is partially supported by the Hong Kong
9 University of Science and Technology–Didi Chuxing (HKUST-DiDi) Joint Laboratory. The
10 third author gratefully acknowledges support by the Lee Kong Chian (LKC) Fellowship awarded
11 by Singapore Management University. The opinions in this paper do not necessarily reflect the
12 official views of the HKUST-DiDi Joint Laboratory. The authors are responsible for all
13 statements. The authors would also like to acknowledge all of the reviewers for their constructive
14 comments.

17 REFERENCES

- 18 Bai, J., So, K.C., Tang, C.S., Chen, X. and Wang, H., 2019. Coordinating supply and demand on
19 an on-demand service platform with impatient customers. *Manufacturing & Service*
20 *Operations Management*, 21(3), pp. 556-570.
- 21 Braverman, A., Dai, J.G., Liu, X. and Ying, L., 2019. Empty-car routing in ridesharing systems.
22 *Operations Research*, 67(5), pp. 1437-1452.
- 23 Cachon, G.P., Daniels, K.M. and Lobel, R., 2017. The role of surge pricing on a service platform
24 with self-scheduling capacity. *Manufacturing & Service Operations Management*, 19(3),
25 pp. 368-384.
- 26 Castillo, J.C., Knoepfle, D. and Weyl, G., 2017. Surge pricing solves the wild goose chase. In
27 *Proceedings of the 2017 ACM Conference on Economics and Computation* (pp. 241-
28 242). ACM.
- 29 Gao, Y., Jiang, D. and Xu, Y., 2018. Optimize taxi driving strategies based on reinforcement
30 learning. *International Journal of Geographical Information Science*, 32(8), pp. 1677-
31 1696.
- 32 Gomes, D.A., 2014. Mean field games models—a brief survey. *Dynamic Games and*
33 *Applications*, 4(2), pp. 110-154.
- 34 Hazelton, M.L. and Watling, D.P., 2004. Computation of equilibrium distributions of Markov
35 traffic-assignment models. *Transportation Science*, 38(3), pp. 331-342.
- 36 Huang, M., 2010. Large-population LQG games involving a major player: the Nash certainty
37 equivalence principle. *SIAM Journal on Control and Optimization*, 48(5), pp. 3318-3353.
- 38 Huang, M., 2012, May. Mean field stochastic games with discrete states and mixed players.
39 In *International Conference on Game Theory for Networks* (pp. 138-151). Springer,
40 Berlin, Heidelberg.
- 41 Huang, M., Caines, P.E. and Malhamé, R.P., 2007. Large-population cost-coupled LQG
42 problems with nonuniform agents: individual-mass behavior and decentralized ϵ -Nash
43 equilibria. *IEEE Transactions on Automatic Control*, 52(9), pp. 1560-1571.

- 1 Huang, M., Malhamé, R.P. and Caines, P.E., 2006. Large population stochastic dynamic games:
2 closed-loop McKean-Vlasov systems and the Nash certainty equivalence
3 principle. *Communications in Information & Systems*, 6(3), pp. 221-252.
- 4 Hwang, R.H., Hsueh, Y.L. and Chen, Y.T., 2015. An effective taxi recommender system based
5 on a spatio-temporal factor analysis model. *Information Sciences*, 314, pp. 28-40.
- 6 Jin, J., Zhou, M., Zhang, W., Li, M., Guo, Z., Qin, Z., Jiao, Y., Tang, X., Wang, C., Wang, J.
7 and Wu, G., 2019. CoRide: Joint order dispatching and fleet management for multi-scale
8 ride-hailing platforms. In *Proceedings of the 28th ACM International Conference on*
9 *Information and Knowledge Management* (pp. 1983-1992).
- 10 Ke, J., Qin, X., Yang, H., Zheng, Z., Zhu, Z. and Ye, J., 2021. Predicting origin-destination ride-
11 sourcing demand with a spatio-temporal encoder-decoder residual multi-graph
12 convolutional network. *Transportation Research Part C: Emerging Technologies*, 122,
13 102858.
- 14 Ke, J., Yang, H., Li, X., Wang, H. and Ye, J., 2020. Pricing and equilibrium in on-demand ride-
15 pooling markets. *Transportation Research Part B: Methodological*, 139, pp. 411-431.
- 16 Ke, J., Yang, H., Zheng, H., Chen, X., Jia, Y., Gong, P. and Ye, J., 2019. Hexagon-based
17 convolutional neural network for supply-demand forecasting of ride-sourcing services.
18 *IEEE Transactions on Intelligent Transportation Systems*, 20(11), pp. 4160-4173.
- 19 Ke, J., Zheng, H., Yang, H. and Chen, X., 2017. Short-term forecasting of passenger demand
20 under on-demand ride services: a spatio-temporal deep learning approach.
21 *Transportation Research Part C: Emerging Technologies*, 85, pp. 591-608.
- 22 Li, M., Qin, Z., Jiao, Y., Yang, Y., Wang, J., Wang, C., Wu, G. and Ye, J., 2019, May. Efficient
23 ridesharing order dispatching with mean field multi-agent reinforcement learning. In *The*
24 *World Wide Web Conference* (pp. 983-994). ACM.
- 25 Lin, K., Zhao, R., Xu, Z. and Zhou, J., 2018, July. Efficient large-scale fleet management via
26 multi-agent deep reinforcement learning. In *Proceedings of the 24th ACM SIGKDD*
27 *International Conference on Knowledge Discovery & Data Mining* (pp. 1774-1783).
28 ACM.
- 29 Lyu, G., Cheung, W.C., Teo, C.P. and Wang, H., 2019. Multi-objective online ride-
30 matching. *Available at SSRN 3356823*.
- 31 Mnih, V., Badia, A. P., Mirza, M., Graves, A., Lillicrap, T., Harley, T., Silver, D and
32 Kavukcuoglu, K. 2016, June. Asynchronous methods for deep reinforcement learning. In
33 *International Conference on Machine Learning* (pp. 1928-1937).
- 34 Mao, C., Liu, Y., & Shen, Z. J. M., 2020. Dispatch of autonomous vehicles for taxi services: A
35 deep reinforcement learning approach. *Transportation Research Part C: Emerging*
36 *Technologies*, 115, 102626.
- 37 Qian, X., Zhang, W., Ukkusuri, S.V. and Yang, C., 2017. Optimal assignment and incentive
38 design in the taxi group ride problem. *Transportation Research Part B:*
39 *Methodological*, 103, pp. 208-226.
- 40 Rong, H., Zhou, X., Yang, C., Shafiq, Z. and Liu, A., 2016, October. The rich and the poor: A
41 Markov decision process approach to optimizing taxi driver revenue efficiency.
42 In *Proceedings of the 25th ACM International on Conference on Information and*
43 *Knowledge Management* (pp. 2329-2334). ACM.

- 1 Shou, Z. and Di, X., 2020a. Multi-Agent Reinforcement Learning for Dynamic Routing Games:
2 A Unified Paradigm. *arXiv preprint arXiv:2011.10915*.
- 3 Shou, Z. and Di, X., 2020b. Reward design for driver repositioning using multi-agent
4 reinforcement learning. *Transportation research part C: Emerging Technologies*, 119,
5 102738.
- 6 Shou, Z., Di, X., Ye, J., Zhu, H., Zhang, H., and Hampshire, R., 2020. Optimal passenger-seeking
7 policies on E-hailing platforms using Markov decision process and imitation learning.
8 *Transportation Research Part C: Emerging Technologies*, 111, 91-113.
- 9 Sun, H., Wang, H. and Wan, Z., 2019a. Model and analysis of labor supply for ride-sharing
10 platforms in the presence of sample self-selection and endogeneity. *Transportation
11 Research Part B: Methodological*, 125, pp. 76-93.
- 12 Sun, H., Wang, H. and Wan, Z., 2019b. Flexible labor supply behavior on ride-sourcing
13 platforms. Available at SSRN:
14 https://papers.ssrn.com/sol3/papers.cfm?abstract_id=3357365.
- 15 Taylor, T.A., 2018. On-demand service platforms. *Manufacturing & Service Operations
16 Management*, 20(4), pp. 704-720.
- 17 Wang, H. and Yang H., 2019. Ride-sourcing systems: a framework and review. *Transportation
18 Research Part B: Methodological*, 129, pp. 122-155.
- 19 Wang, Z., Qin, Z., Tang, X., Ye, J. and Zhu, H., 2018, November. Deep reinforcement learning
20 with knowledge transfer for online rides order dispatching. In *2018 IEEE International
21 Conference on Data Mining (ICDM)* (pp. 617-626). IEEE.
- 22 Xu, Z., Li, Z., Guan, Q., Zhang, D., Li, Q., Nan, J., Liu, C., Bian, W. and Ye, J., 2018, July.
23 Large-scale order dispatch in on-demand ride-hailing platforms: a learning and planning
24 approach. In *Proceedings of the 24th ACM SIGKDD International Conference on
25 Knowledge Discovery & Data Mining* (pp. 905-913). ACM.
- 26 Xu, Z., Yin, Y. and Zha, L., 2017. Optimal parking provision for ride-sourcing services.
27 *Transportation Research Part B: Methodological*, 105, pp. 559-578.
- 28 Yang, H., Fung, C.S., Wong, K.I. and Wong, S.C., 2010. Nonlinear pricing of taxi
29 services. *Transportation Research Part A: Policy and Practice*, 44(5), pp. 337-348.
- 30 Yang, H., Qin, X., Ke, J. and Ye, J., 2020a. Optimizing matching time interval and matching
31 radius in on-demand ride-sourcing markets. *Transportation Research Part B:
32 Methodological*, 131, pp. 84-105.
- 33 Yang, H., Shao, C., Wang, H. and Ye, J., 2020b. Integrated reward scheme and surge pricing in
34 a ridesourcing market. *Transportation Research Part B: Methodological*, 134, pp. 126-
35 142.
- 36 Yang, H., Wong, S.C. and Wong, K.I., 2002. Demand–supply equilibrium of taxi services in a
37 network under competition and regulation. *Transportation Research Part B:
38 Methodological*, 36(9), pp. 799-819.
- 39 Yao, H., Wu, F., Ke, J., Tang, X., Jia, Y., Lu, S., Gong, P., Ye J. and Li, Z., 2018. Deep multi-
40 view spatial-temporal network for taxi demand prediction. In *Thirty-Second AAAI
41 Conference on Artificial Intelligence*.
- 42 Yu, X., Gao, S., Hu, X. and Park, H., 2019. A Markov decision process approach to vacant taxi
43 routing with e-hailing. *Transportation Research Part B: Methodological*, 121, pp. 114-
44 134.

- 1 Zha, L., Yin, Y., Xu, Z., 2018. Geometric matching and spatial pricing in ride-sourcing
2 markets. *Transportation Research Part C: Emerging Technologies*, 92, pp. 58-75.
- 3 Zha, L., Yin, Y. and Yang, H., 2016. Economic analysis of ride-sourcing
4 markets. *Transportation Research Part C: Emerging Technologies*, 71, pp. 249-266.
- 5 Zhang, L., Hu, T., Min, Y., Wu, G., Zhang, J., Feng, P., Gong, P. and Ye, J., 2017. A taxi order
6 dispatch model based on combinatorial optimization. In *Proceedings of the 23rd ACM*
7 *SIGKDD International Conference on Knowledge Discovery and Data Mining* (pp. 2151-
8 2159). ACM.
- 9 Zhu, Z., Mardan, A., Zhu, S. and Yang, H., 2021a. Capturing the interaction between travel time
10 reliability and route choice behavior based on the generalized Bayesian traffic model.
11 *Transportation Research Part B: Methodological*, 143, 48-64.
- 12 Zhu, Z., Tang, L., Xiong, C., Chen, X. and Zhang, L., 2019a. The conditional probability of
13 travel speed and its application to short-term prediction. *Transportmetrica B: Transport*
14 *Dynamics*, 7(1), pp. 684-706.
- 15 Zhu, Z., Sun, L., Chen, X. and Yang, H., 2021b. Integrating probabilistic tensor factorization
16 with Bayesian supervised learning for dynamic ridesharing pattern analysis. Accepted by
17 *Transportation Research Part C: Emerging Technologies*.
- 18 Zhu, Z., Qin X., Ke, J., Zheng, Z. and Yang, H., 2020. Analysis of multi-modal commute
19 behavior with feeding and competing ridesplitting services. *Transportation Research*
20 *Part A: Policy and Practice*, 132, 713-727.
- 21 Zhu, Z., Zhu, S., Zheng, Z. and Yang, H., 2019b. A generalized Bayesian traffic
22 model. *Transportation Research Part C: Emerging Technologies*, 108, pp. 182-206.
- 23 Zuniga Garcia, N., 2019. *Spatial pricing empirical evaluation of ride-sourcing trips using the*
24 *graph-fused lasso for total variation denoising* (doctoral dissertation).

25
26

27 **APPENDIX A: FORMULAS FOR THE SPECIFIC MF-MDP MODEL**

28 In this appendix we provide detailed formulas for state transition laws, matching
29 probability, and rewards in the MF-MDP model developed in Section 4.

30 First, we illustrate the state transition laws of the platform and a driver in the standard
31 MF-MDP model (or the representative driver in the simplified MF-MDF model). With the state
32 vector $\mathbf{s} = [t, s_1, \dots, s_o]$ and the action vector $\mathbf{a} = [a_1, \dots, a_o]$, the state transition law for the
33 platform is given by

$$Q(\mathbf{s}'|\mathbf{s}, \mathbf{a}) = \begin{cases} 1 & , \mathbf{s}' = [t + 1, a_1, \dots, a_o], \forall \mathbf{a} \in \mathbf{A} \\ 0 & , otherwise \end{cases} \quad (\text{A.1})$$

34 An intuitive explanation of Eq. (A.1) is that the state of the platform for the next time index
35 equals its action vector.

36 With a state vector $\mathbf{s}_d = [t, s_{d1}, s_{d2}, s_{d3}, s_{d4}, s_{d5}]$ and an action $a_d \in \mathbf{J}_{s_{d3}}$, the state
37 transition law for a driver has different formulas according to the current task s_{d1} and remaining
38 time s_{d2} . Note that for the following formulas, we need $\mathbf{s}_d, \mathbf{s}'_d \in \mathbf{S}_d$, and $\mathbf{h}_d \in \mathbf{H}_d$.

- 1 • The driver is picking up a passenger, i.e., $\mathbf{s}_d = [t, 1, \tau, o, o', o'']$, $\tau > 0$, $o, o', o'' \in$
2 $\{1, 2, \dots, O\}$, and $o \neq o' \neq o''$

$$Q_d(\mathbf{s}'_d | \mathbf{s}_d, \mathbf{h}_d, a_d) = \begin{cases} 1 & , \mathbf{s}'_d = [t + 1, 1, \tau - 1, o, o', o''], \tau > 1 \\ 1 & , \mathbf{s}'_d = [t + 1, 2, \tau_{o'o''}, o', o', o''], \tau = 1 \\ 0 & , \text{otherwise} \end{cases} \quad (\text{A.2})$$

3 where $\tau_{o'o''}$ denotes the average travel time from zone o' to zone o'' , which is
4 exogenous. The first line means that the driver still needs more than one time period to
5 finish the current picking-up task; therefore the new remaining time s'_{d2} decreases by one,
6 while the other dimensions of the state vector remain unchanged. The second line means
7 that the driver is about to finish a picking-up task and will immediately change the task
8 to delivering the passenger; then the new remaining time s'_{d2} becomes the average travel
9 time between the origin zone and the destination zone, the new task s'_{d1} becomes 2, and
10 the new idling zone becomes the current picking-up destination, which is identical to the
11 origin of the passenger (i.e., $s'_{d3} = s_{d4} = o'$).

- 12 • The driver is delivering a passenger, i.e., $\mathbf{s}_d = [t, 2, \tau, o, o, o']$, $\tau > 0$, $o, o' \in$
13 $\{1, 2, \dots, O\}$, and $o \neq o'$.

$$Q_d(\mathbf{s}'_d | \mathbf{s}_d, \mathbf{h}_d, a_d) = \begin{cases} 1 & , \mathbf{s}'_d = [t + 1, 2, \tau - 1, o, o, o'], \tau > 1 \\ 1 & , \mathbf{s}'_d = [t + 1, 0, 0, o', o', o'], \tau = 1 \\ 0 & , \text{otherwise} \end{cases} \quad (\text{A.3})$$

14 The first line means that the driver needs more than one time period to finish the current
15 delivering task and the new remaining time s'_{d2} decreases by one, while the other
16 dimensions of the state vector remain the same. The second line indicates that if the driver
17 is about to finish a delivering task, the new state becomes a purely idle state (i.e., $s'_{d1} =$
18 0 and $s'_{d2} = 0$). Since there is no passenger order, we let $s'_{d3} = s'_{d4} = s'_{d5} = o'$ for
19 convenience.

- 20 • The driver is purely idle, i.e., $\mathbf{s}_d = [t, 0, 0, o, o, o]$, and $o \in \{1, 2, \dots, O\}$.

$$Q_d(\mathbf{s}'_d | \mathbf{s}_d, \mathbf{h}_d, a_d) = \begin{cases} m_{o,o',o''}(\mathbf{h}_d) & , \mathbf{s}'_d = [t + 1, 1, \tau_{oo'}, o, o', o''] \\ m_{o,o,o''}(\mathbf{h}_d) & , \mathbf{s}'_d = [t + 1, 2, \tau_{oo''}, o, o, o''] \\ u_o(\mathbf{h}_d) & , \mathbf{s}'_d = [t + 1, 0, 0, o, o, o], a_d = o \\ u_o(\mathbf{h}_d) & , \mathbf{s}'_d = [t + 1, 0, \tau_{oa_d} - 1, o, a_d, a_d], a_d \neq o \\ 0 & , \text{otherwise} \end{cases} \quad (\text{A.4})$$

21 where $m_{o,o',o''}(\mathbf{h}_d)$ denotes the probability of getting a matched passenger order at zone
22 o with the origin and destination of the passenger being o' and o'' , respectively, and
23 $u_o(\mathbf{h}_d)$ denotes the probability of not being matched at zone o . Generally, both
24 $m_{o,o',o''}(\mathbf{h}_d)$ and $u_o(\mathbf{h}_d)$ depend on the MF state of drivers and the exogenous passenger
25 demand. More specific formulas for the probabilities are given later in this appendix. The
26 first line in the equation implies that the origin of the newly matched passenger is
27 different from the driver's current idling zone (i.e., $s'_{d4} = o' \neq s_{d3} = o$); therefore, a
28 picking-up task is needed and we have $s'_{d1} = 1$ and $s'_{d2} = \tau_{oo'}$. In the second line, the
29 matched passenger and the driver are in the same zone (i.e., $s'_{d4} = s_{d3} = o$) and we
30 assume the picking-up process can be ignored; therefore, the driver will directly start to

1 deliver the passenger (i.e., $s'_{d1} = 2$). The third line indicates that the driver is still not
2 matched and their action is to stay in the current idling zone (i.e., $a_d = o$); therefore the
3 state of the driver will remain unchanged. In the fourth line, the driver is not matched and
4 will relocate to zone a_d ; we let $s'_{d4} = s'_{d5} = a_d$ for convenience, and let $s'_{d2} = \tau_{oa_d} - 1$
5 because we assume the driver is already in the middle of the self-relocating state (i.e., no
6 time is wasted by stopping the vehicle to load or drop off passengers).
7 • The driver is in a self-relocating state, i.e., $s_d = [t, 0, \tau, o, o', o']$, $\tau > 0$, and $o \neq o'$.

$$Q_d(s'_d | s_d, \mathbf{h}_d, a_d) = \begin{cases} m_{o,o'',o'''}(\mathbf{h}_d) & , s'_d = [t + 1, 1, \tau_{oo''}, o, o'', o'''] \\ m_{o,o',o'''}(\mathbf{h}_d) & , s'_d = [t + 1, 1, \tau - 1, o, o', o'''], \tau > 1 \\ m_{o,o',o'''}(\mathbf{h}_d) & , s'_d = [t + 1, 2, \tau_{o'o'''}, o', o', o'''], \tau = 1 \\ m_{o,o,o'''}(\mathbf{h}_d) & , s'_d = [t + 1, 2, \tau_{oo'''}, o, o, o'''] \\ u_o(\mathbf{h}_d) & , s'_d = [t + 1, 0, \tau - 1, o, o', o'], \tau > 1 \\ u_o(\mathbf{h}_d) & , s'_d = [t + 1, 0, 0, o', o', o'], \tau = 1 \\ 0 & , otherwise \end{cases} \quad (\text{A.5})$$

8 In the first line, the driver is matched with a passenger whose origin $s'_{d4} = o''$ is different
9 from either the driver's current idling zone $s_{d3} = o$ or the self-relocation destination
10 $s_{d4} = o'$; therefore, the driver begins a picking-up task and moves to zone o'' (i.e., $s'_{d1} =$
11 1 and $s'_{d4} = o''$). The second line indicates that if the self-relocation destination
12 coincides with the origin zone of the matched passenger (i.e., $s'_{d4} = s_{d4} = o'$), the new
13 remaining time s'_{d2} decreases by one because we regard the driver as already in the
14 middle of the picking-up task (i.e., $s'_{d1} = 1$). To continue with the case $s'_{d4} = s_{d4} = o'$,
15 the third line means that if drivers are leaving the self-relocating state (i.e., $s_{d2} = \tau = 1$),
16 they immediately load the passenger and begin the delivering task (i.e., $s'_{d1} = 2$). In the
17 fourth line, both the matched passenger and driver are in zone o (i.e., $s'_{d4} = s_{d3} = o$) and
18 a delivering task starts. The fifth and sixth lines indicate that the driver is not matched
19 with passengers, such that he/she either remains in the self-relocating state (the fifth line)
20 or becomes purely idle in zone o' (i.e., $s'_{d1} = 0$, $s'_{d2} = 0$, and $s'_{d3} = o'$ for the sixth line).

21
22 Note that the MF vector \mathbf{h}_d is used in $Q_d(s'_d | s_d, \mathbf{h}_d, a_d)$ to calculate the matching
23 probabilities $m_{o,o',o''}(\mathbf{h}_d)$ and $u_o(\mathbf{h}_d)$. Next, we provide detailed formula related to the number
24 of matches and matching probabilities. According to the edge-based matching rule in Section
25 4.1, the number of matches near an arbitrary edge $e_{oo'}$, which is denoted by $k_{e_{oo'}}$, is given by

$$k_{e_{oo'}}(\mathbf{h}_d) = \min \left\{ \frac{M_o(\mathbf{h}_d)}{E_o} + \frac{M_{o'}(\mathbf{h}_d)}{E_{o'}}, \frac{N_o}{E_o} + \frac{N_{o'}}{E_{o'}} \right\} \quad (\text{A.6})$$

$$M_o(\mathbf{h}_d) = M \left(h_{d,[t,0,0,0,o,o]} + \sum_{\tau > 0, o' \in J_{o/o}} h_{d,[t,0,\tau,o,o',o']} \right) \quad (\text{A.7})$$

26 where E_o denote the number of edges of zone o , and assuming hexagonal zones, the value of E_o
27 is 6 unless the zone is located at the boundary of the network; and h_{d,s_d} is the scalar value in \mathbf{h}_d
28 and represents the proportion of drivers in state s_d (i.e., $z_{d,s_d}^t = h_{d,s_d}$ in Eq. (2)), such that
29 $h_{d,[t,0,0,0,o,o]}$ denotes the proportion of purely idle drivers and $h_{d,[t,0,\tau,o,o',o']}$ the proportion of

1 drivers in self-relocating states at time t (see Eq. (A.5)). Since $k_{e_{oo'}}$ is only used to calculate
 2 matching probabilities in this paper, we allow the value of $k_{e_{oo'}}$ to be a non-integer.

3 Since the numbers of matched passengers and drivers are proportional to the demand and
 4 supply near the common edge, and drivers and passengers are uniformly distributed in the zones
 5 (see Section 4.1), we have the formulas for $m_{o,o',o''}(\mathbf{h}_d)$ and $u_o(\mathbf{h})$ as follows:

$$m_{o,o',o''}(\mathbf{h}_d) = \begin{cases} \frac{1}{E_o} \frac{k_{e_{oo'}}(\mathbf{h}_d)}{M_o(\mathbf{h}_d) + \frac{M_{o'}(\mathbf{h}_d)}{E_{o'}}} \frac{N_{o'}}{N_o + N_{o'}} \frac{N_{o'o''}}{N_{o'}}, o' \in J_o/o \\ \left(\sum_{o''' \in J_o/o} \frac{1}{E_o} \frac{k_{e_{oo'''}}(\mathbf{h}_d)}{M_o(\mathbf{h}_d) + \frac{M_{o'''}(\mathbf{h}_d)}{E_{o'''}}} \frac{N_o}{N_o + N_{o'''}} \right) \frac{N_{oo''}}{N_o}, o' = o \end{cases} \quad (\text{A.8})$$

$$u_o(\mathbf{h}_d) = 1 - \sum_{o' \in J_o, o'' \in O} m_{o,o',o''}(\mathbf{h}_d) \quad (\text{A.9})$$

6 where $N_{o'o''}$ denotes exogenous passenger demand from zone o' to zone o'' . In the first line we
 7 calculate the probability of matching an adjacent passenger in zone o (i.e., picking-up is needed);
 8 the term $\frac{1}{E_o}$ denotes the probability that the driver is near edge $e_{oo'}$; the term $\frac{k_{e_{oo'}}}{M_o/E_o + M_{o'}/E_{o'}}$ the
 9 chance of getting a match for drivers who are near edge $e_{oo'}$; the term $\frac{N_{o'}}{N_o + N_{o'}}$ the probability that
 10 the origin of the matched passenger is zone o' ; and the term $\frac{N_{o'o''}}{N_{o'}}$ the chance that the matched
 11 passenger's destination is o'' . The second line denotes the probability of matching a local
 12 passenger in zone o (i.e., direct delivery without picking-up), in which we sum the matching
 13 probabilities from all common edges between zone o and its adjacent zones (i.e., the summation
 14 term for $o''' \in J_o/o$); the explanation for each term is similar to that for the first line. The
 15 unmatched probability equals one minus all matched probabilities in zone o (i.e., Eq. (A.9)).

16 Last, we show the detailed calculation of $r(\mathbf{s}, \mathbf{h}_d, \mathbf{h}'_d)$ and $r_d(\mathbf{s}, \mathbf{s}_d, \mathbf{s}'_d)$. The
 17 decomposition of the platform reward $r(\mathbf{s}, \mathbf{h}_d, \mathbf{h}'_d)$ is given by

$$r(\mathbf{s}, \mathbf{h}_d, \mathbf{h}'_d) = r_1(\mathbf{h}_d) - r_2(\mathbf{s}, \mathbf{h}_d, \mathbf{h}'_d) + \mu r_3(\mathbf{s}, \mathbf{h}_d, \mathbf{h}'_d) \quad (\text{A.10})$$

$$r_1(\mathbf{h}_d) = \eta \alpha M \sum_{\tau, o, o'} h_{d, [t, 2, \tau, o, o'] } \quad (\text{A.11})$$

$$\begin{aligned} & r_2(\mathbf{s}, \mathbf{h}_d, \mathbf{h}'_d) \\ = & \beta M \left(\sum_{o, o', o'' | s_o = \beta} h_{d, [t, 1, \tau, o', o', o, o'']} + \sum_{o, o', o'' | s_o = \beta} \left(h_{d, [t, 2, \tau, oo'', o, o, o'']} - h'_{d, [t-1, 1, 1, o', o, o'']} \right) \right) \\ & + \sum_{o, o', o'', \tau > 1 | s_o = \beta} \left(h_{d, [t, 1, \tau-1, o', o, o'']} - h'_{d, [t-1, 1, \tau, o', o, o'']} \right) \end{aligned} \quad (\text{A.12})$$

$$\begin{aligned}
& r_3(\mathbf{s}, \mathbf{h}_d, \mathbf{h}'_d) \\
= & \frac{M}{N} \left(\sum_{o, o', o''} h_{d, [t, 1, \tau, o', o, o'']} + \sum_{o, o', o''} \left(h_{d, [t, 2, \tau, o, o, o'']} - h'_{d, [t-1, 1, 1, o', o, o'']} \right) \right) \\
& + \sum_{o, o', o'', \tau > 1} \left(h_{d, [t, 1, \tau-1, o', o, o'']} - h'_{d, [t-1, 1, \tau, o', o, o'']} \right)
\end{aligned} \tag{A.13}$$

1 In Eq. (A.11), the commission withheld during one time period is calculated based on the
2 proportion of drivers who are performing delivery task $h_{d, [t, 2, \tau, o, o, o']}$ (i.e., the proportion of
3 drivers in state $[t, 2, \tau, o, o, o']$); the total number of drivers M ; commission rate η ; and trip fare
4 rate α . In Eq. (A.12), the proportion of drivers who are offered subsidies consists of three terms:
5 $h_{d, [t, 1, \tau, o', o', o, o']}$ denotes newly matched/dispatched drivers who are currently neither in nor self-
6 relocating to the subsidized zones but will pick up passengers there; $h_{d, [t, 2, \tau, o, o, o]} -$
7 $h'_{d, [t-1, 1, 1, o', o, o]}$ denotes newly matched/dispatched drivers who are currently in the subsidized
8 zones; and $h_{d, [t, 1, \tau-1, o', o, o]} - h'_{d, [t-1, 1, \tau, o', o, o]}$ ($\tau > 1$) denotes newly matched/dispatched
9 drivers who are coincidentally in the process of self-relocating to the subsidized zones. In Eq.
10 (A.13), the service rate is calculated via the number of drivers M , the proportion of drivers who
11 are newly matched, and the total number of passenger demand N . Note that in Eqs. (A.12)–
12 (A.13), the term $h_{d, [t, 2, \tau, o, o, o]}$ also includes previously dispatched drivers who just finished
13 picking-up tasks at subsidized zone o (i.e., the term $h'_{d, [t-1, 1, 1, o', o, o]}$); therefore, we need a
14 subtraction, $h_{d, [t, 2, \tau, o, o, o]} - h'_{d, [t-1, 1, 1, o', o, o]}$, to only count newly matched drivers.
15 Similarly, the term $h_{d, [t, 1, \tau-1, o', o, o]}$ also includes previously dispatched drivers who are on the
16 way to pick up passengers in zone o (i.e., the term $h'_{d, [t-1, 1, \tau, o', o, o]}$, $\tau > 1$), and we need a
17 subtraction to exclude these drivers¹⁵.

18 The decomposition of a driver's one-step reward $r_d(\mathbf{s}, \mathbf{s}_d, \mathbf{s}'_d)$ is as follows:

$$r_d(\mathbf{s}, \mathbf{s}_d, \mathbf{s}'_d) = r_{d1}(\mathbf{s}_d) + r_{d2}(\mathbf{s}, \mathbf{s}_d, \mathbf{s}'_d). \tag{A.14}$$

$$r_{d1}(\mathbf{s}_d) = \begin{cases} (1 - \eta)\alpha & , s_{d1} = 2 \\ 0 & , otherwise \end{cases} \tag{A.15}$$

$$r_{d2}(\mathbf{s}, \mathbf{s}_d, \mathbf{s}'_d) = \begin{cases} \beta & , s_{d1} = 2, s'_{d1} = 0, s_{d3} = o, s_o = \beta \\ \beta & , s_{d1} = 1, s'_{d1} = 0, s_{d3} = o, s_o = \beta \\ 0 & , otherwise \end{cases} \tag{A.16}$$

19 In the first line of Eq. (A.15), we assume the fare is uniformly collected when the driver is
20 delivering the passenger, i.e., $s_{d1} = 2$. For instance, if a driver delivers a passenger from time t_1
21 to t_2 (i.e., $s_{d1} = 2$ for $t \in \{t_1, \dots, t_2\}$), the driver will receive an income of $(1 - \eta)\alpha$ for each
22 time period during the delivering task, and the total income from trip fare equals
23 $(t_2 - t_1 + 1)(1 - \eta)\alpha$. In Eq. (A.16), a subsidy β for drivers is executed immediately after the

¹⁵ Based on Eq. (A.2), for drivers who are in state $[t - 1, 1, 1, o', o, o']$ at time $t - 1$, their states become $[t, 2, \tau, o, o, o']$ at time t . Therefore, these drivers are counted in the term $h_{d, [t, 2, \tau, o, o, o]}$. Still based on Eq. (A.2), for drivers who are in state $[t - 1, 1, \tau, o', o, o']$ ($\tau > 1$) at time $t - 1$, their states become $[t, 1, \tau - 1, o', o, o']$, and these drivers are counted in the term $h_{d, [t, 1, \tau-1, o', o, o]}$.

1 task switches from “idle” to “picking-up” or “delivering” (i.e., $s'_{d1} = 0$ and $s_{d1} \neq 0$) and the
2 origin of the matched passenger is subsidized (i.e., $s_{d3} = 0$ and $s_o = \beta$); the first line indicates
3 that the matched passenger is local (i.e., $s_{d3} = s_{d4}$); and the second line that the matched
4 passenger is in an adjacent zone (i.e., $s_{d3} \neq s_{d4}$).

5