

Dynamic Mean-Field Control for Network MDPs with Exogenous Demand

Anonymous authors

Paper under double-blind review

Abstract

1 This paper studies the network control problems with exogenous demand, where net-
 2 work controller must dynamically allocate resources to satisfy exogenous demands with
 3 unknown distributions. We formalize the problem using Networked Markov Decision
 4 Processes with Exogenous Demands (Exo-NMDPs), where the system states are de-
 5 coupled into endogenous states and stochastic exogenous demands. However, Exo-
 6 NMDPs pose three main challenges: scalability in large-scale networks; stochasticity
 7 from fluctuating exogenous demands; and delayed feedback of scheduling actions. To
 8 address these issues, we propose the Dynamic Mean-Field Control (DMFC) algorithm,
 9 a scalable and computationally efficient approach for matching exogenous demands.
 10 Specifically, DMFC transforms the high-dimensional actual states of the Exo-NMDP
 11 into low-dimensional mean-field states, and dynamically optimizes the policy by solv-
 12 ing a mean-field control problem at each time step. This enables DMFC to capture
 13 spatiotemporal correlations between demand and system state, while remaining robust
 14 against demand fluctuations and action execution delay. We validate DMFC on two
 15 representative scenarios: supply-chain inventory management and vehicle routing. Our
 16 experimental results show that DMFC adapts well to various demand patterns and out-
 17 performs state-of-the-art baselines in both scenarios.

18 1 Introduction

19 Network control problems with exogenous demand has a broad application in real-world scenarios,
 20 including: supply chain management [Bellamy & Basole \(2013\)](#); [Zhang et al. \(2014\)](#); [Aminzade-
 21 gan et al. \(2019\)](#), scheduling in robotic systems [Rus; \(2012\)](#); [Pavone; \(2016\)](#), and vehicle routing
 22 in mobility-on-demand systems [Bullo et al. \(2011\)](#); [Holler et al. \(2019\)](#); [Gammelli et al. \(2021\)](#).
 23 The control policies are the operational backbone of these systems, enhancing service reliability
 24 and driving cost reduction through dynamic and adaptive agent scheduling. However, designing
 25 such policies is challenging due to: 1) scalability in large-scale networks, 2) stochasticity from fluc-
 26 tuating exogenous demand, and 3) delayed feedback of control actions. These challenges lead to
 27 spatiotemporal mismatches between agents and demand, requiring a control policy π that adapts to
 28 evolving demand while incorporating past decisions for effective coordination.

29 We formulate network control problems with exogenous demand as a Networked Markov Decision
 30 Processes with Exogenous Demands (Exo-NMDPs), extending Exo-MDPs [Sinclair et al. \(2023\)](#) into
 31 networking settings. To address Exo-NMDPs, we introduce a mean-field Exo-NMDP formulation
 32 to transform the high-dimensional actual state into low-dimensional mean-field state, capturing both
 33 the endogenous system dynamics and the exogenous demand signals. Based on this formulation,
 34 we develop the Dynamic Mean-Field Control (DMFC) framework, which operates in two stages at
 35 each time step t : First, DMFC constructs an ideal mean-field state by incorporating predicted future
 36 demand into the current endogenous system state. Second, it solves a linear program that yields the
 37 control policy π_t by optimizing the system objective subject to mean-field dynamics and constraints.

38 The major contributions of our paper are listed below.

- 39 • **Problem Formulation.** We model the network control problems with exogenous demand as a
40 networked Markov decision process with exogenous demands (Exo-NMDP), where the system
41 states are decoupled into endogenous states and stochastic exogenous demands. The decomposi-
42 tion enables a more efficient and concise representation for policy design.
- 43 • **Algorithm Design.** We propose Dynamic Mean-Field Control (DMFC) algorithm for Exo-
44 NMDP. DMFC leverages historical information to infer ideal mean-field states and then synthe-
45 sizes a control policy to align system states towards the target states. This addresses “the curse of
46 networked agents” and intricate spatio-temporal correlations.
- 47 • **Experimental Evaluation.** We evaluate DMFC in two real-world scenarios: supply chain in-
48 ventory management and vehicle routing in mobility-on-demand systems. Experimental results
49 show that our algorithm outperforms the state-of-the-art baseline in both applications. Our code
50 is available at <https://anonymous.4open.science/r/DMFC-1247-F3E4T38D/>.

51 2 Related Work

52 Network control problems with exogenous demand have been studied from three main perspectives:
53 RL-based (Reinforcement Learning), MPC-based (Model Predictive Control), and queueing-based
54 methods. Here, we provide an overview of each approach. As mobility-on-demand systems are the
55 representative example, we include them in the discussion.

56 **RL-based methods:** (Deep) reinforcement learning methods [Sutton & Barto \(2018\)](#); [Mnih et al. \(2013\)](#); [Ladosz et al. \(2022\)](#) offers promising solutions for network control problems with exoge-
57 nous demand, particularly in mobility-on-demand systems [Qin et al. \(2022\)](#); [Wen et al. \(2024\)](#).
58 Recent hybrid approaches integrate RL with optimization: [Wei et al. \(2024\)](#) combine the mean-field
59 optimization formulation in [Braverman et al. \(2019\)](#) with temporal difference (TD) learning, but
60 relies heavily on historical demand trajectories; [Gammelli et al. \(2023\)](#) develop a Graph RL frame-
61 work under bi-level optimization that first predicts reward-driven desired next states then solves for
62 policy from linear programs, yet struggles with scalability and generality in complex networks. Be-
63 yond single-agent paradigms, multi-agent RL (MARL) approaches encounter inherent dimensionality
64 challenges. For instance, [Lin et al. \(2019\)](#) models vehicles as discrete-action agents, where com-
65 plexity grow exponentially with agent counts. To address scalability limitations, [Liu et al. \(2022\)](#)
66 employs regional agents while [Wang et al. \(2024\)](#) implements dynamic parameter sharing. How-
67 ever, communication and coordination between agents remain unsolved. These challenges motivate
68 mean-field control (MFC) [Gast et al. \(2012\)](#); [Bäuerle \(2023\)](#) and mean-field RL (MFRL) [Carmona
69 et al. \(2023\)](#); [Pásztor et al. \(2023\)](#); [Jusup et al. \(2024\)](#), which simplify multi-agent interactions
70 through representative agents operating within aggregated distributions. A practical demonstration
71 comes from [Jusup et al. \(2025\)](#), who apply MFC and MFRL to vehicle rebalancing via fleet dynam-
72 ics modeling to avoid direct vehicle-to-vehicle coordination. However, their framework assumes
73 instantaneous task completion, which is a critical limitation given operational delays observed in
74 physical systems.

76 **MPC-based methods:** Model Predictive Control (MPC) [Kouvaritakis & Cannon \(2016\)](#); [Borrelli
77 et al. \(2017\)](#) has become predominant in network control problems with exogenous demand, partic-
78 ularly in vehicle repositioning scenarios [Iglesias et al. \(2018\)](#); [Tsao et al. \(2019\)](#); [Aalipour & Khani
79 \(2024\)](#). However, MPC performance critically depends on model accuracy, and developing a high-
80 fidelity model is both time and data-intensive. Furthermore, as the look-ahead steps increase, the
81 computational complexity of MPC also increases, especially for a large-scale network with substan-
82 tial agents.

83 **Queueing-based methods:** Queueing theory [Shortle et al. \(2018\)](#) provides a framework for net-
84 work resource allocation under exogenous demand constraints. Prior works [Iglesias et al. \(2016\)](#);
85 [Banerjee et al. \(2022\)](#) investigate the problem within a queueing network, deriving routing poli-
86 cies by analyzing the stationary distribution of a Markov chain. One of the most closely related

studies in this area is [Braverman et al. \(2019\)](#), which formulates the problem as a BCMP queueing network [Baskett et al. \(1975\)](#) and employs a mean-field optimization approach to determine the optimal mean-field control policy. However, these studies primarily focus on the steady-state solutions under static environments, lacking adaptability to non-stationary or time-varying demands, where system states evolve over time.

We propose the mean-field Exo-NMDP framework to model network control problems with exogenous demand, and introduce DMFC, a scalable and robust algorithm that coordinates agents via mean-field control while explicitly accounting for operational delays. Experiments show that DMFC offers superior adaptability to dynamic demand variations and outperforms the state-of-the-art methods in supply chain inventory management and vehicle routing in mobility-on-demand systems.

3 Problem Formulation

In this paper, we study a large-scale network control problem with exogenous demand over a network $\mathcal{G} = (\mathcal{V}, \mathcal{E})$, where \mathcal{V} and \mathcal{E} are the sets of nodes and edges. At each time step, the system makes decisions to dispatch agents from node u to node v for optimizing exogenous demand satisfaction. We formulate the problem as a Networked Markov Decision Process with Exogenous Demand.

3.1 Networked MDPs with Exogenous Demands (Exo-NMDPs)

We propose Networked Markov Decision Processes with Exogenous Demand (Exo-NMDPs), which extend Exo-MDPs [Sinclair et al. \(2023\)](#) to multi-agent systems where agents are dynamically scheduled in response to exogenous demand. An Exo-NMDP is defined by the tuple $(\mathcal{N}, \mathcal{G}, \mathcal{S}, \Xi, \mathcal{A}, \mathcal{P}_{\mathcal{S}}, \mathcal{P}_{\Xi}, \mathcal{R}, \gamma)$. Here, $\mathcal{N} = \{1, 2, \dots, n\}$ is the set of agents operating on network topology \mathcal{G} . The endogenous state space is $\mathcal{S} := \prod_{k \in \mathcal{N}} \mathcal{S}_k$, where each agent k has a local state $s_{k,t} \in \mathcal{S}_k$. The exogenous demand process is modeled as a stochastic sequence $\Xi := \{\xi_t\}_{t \geq 0}$, where each $\xi_t = \{\xi_{v,t}\}_{v \in \mathcal{V}} \subset \mathbb{N}_+^{|\mathcal{V}|}$ represents node-specific, time-varying demands. In our setting, ξ_t is revealed at the start of time t , as opposed to [Sinclair et al. \(2023\)](#), where it becomes known only after actions are taken. Given ξ_t , the system selects joint actions $a_t = (a_{1,t}, \dots, a_{n,t}) \in \mathcal{A} := \prod_{k \in \mathcal{N}} \mathcal{A}_k$ via a policy $\pi_t : \mathcal{S} \times \Xi^t \rightarrow \Delta(\mathcal{A})$, which maps current state s_t and demand history $\xi_{1:t}$ to a distribution over actions. The system endogenous state evolves according to $s_{t+1} \sim \mathcal{P}_{\mathcal{S}}(\cdot | s_t, a_t, \xi_t)$, while the exogenous demand follows a stochastic process $\xi_{t+1} \sim \mathcal{P}_{\Xi}(\cdot | \xi_{1:t})$. The reward function $\mathcal{R} : \mathcal{S} \times \Xi \times \mathcal{A} \rightarrow \mathbb{R}$ evaluates system performance at each time step, for example, based on revenue from satisfied demand. Future rewards are discounted by a factor $\gamma \in [0, 1)$. We assume that both $\mathcal{P}_{\mathcal{S}}$ and \mathcal{R} are known, and that uncertainty arises only from \mathcal{P}_{Ξ} . The objective is to learn a policy that maximizes the expected long-term cumulative reward.

Exo-NMDPs as Multi-agent Semi-Markov Decision Processes. In Exo-NMDPs, agent actions may involve delays, such as travel time between network nodes. To model these temporal extensions, we adopt the option framework [Sutton et al. \(1999\)](#), where each option represents a temporally extended action that may span multiple time steps. For each individual agent, the option framework transforms its local MDP into a semi-Markov decision process (SMDP) [Ross \(1992\)](#), as the original action set \mathcal{A} is replaced by a fixed set of options \mathcal{O} (without loss of generality, in the subsequent discussion, we will abuse the term “actions” to also refer to “options”). While each agent operates under its own SMDP, the overall Exo-NMDP cannot be reduced to a single SMDP due to the asynchronous nature of agent decisions and variable option durations. Instead, we model the system as a Multiagent SMDP (MSMDP) [Ghavamzadeh & Mahadevan \(2004\)](#), where decision epochs are aligned with fixed-length time intervals (e.g., every 5 minutes). At each decision epoch t , only the subset of agents whose previous actions have just completed make new decisions, while the remaining agents continue executing their current actions. The MSMDP formulation of Exo-NMDPs enables tractable analysis across agents; however, the challenges of multi-agent scalability and communication remain unresolved.

Mean-Field Representation. To address the scalability and communication limitations inherent in Exo-NMDPs, we employ a mean-field formulation that aggregates individual agent behaviors into node-level dynamics, treating network nodes as the primary decision units. Let $\ell_{k,t} \in \mathcal{L}$ denote the location of agent k at time t , which is embedded in its full state $s_{k,t} := (\ell_{k,t}, \dots)$. The system mean-field state is defined as: $\mu_t = \{\mu_t^v(i)\}_{i \in \mathcal{V}} \cup \{\mu_t^e(u, v)\}_{(u,v) \in \mathcal{E}}$, where $\mu_t^v(i) := \frac{1}{n} \sum_{k=1}^n \mathbb{I}(\ell_{k,t} = v_i)$ denotes the fraction of agents currently idle at node i , and $\mu_t^e(u, v) := \frac{1}{n} \sum_{k=1}^n \mathbb{I}(\ell_{k,t} = e_{uv})$ captures the fraction of agents in transit along edge $e_{uv} \in \mathcal{E}$. The former corresponds to agents available for new decisions at time t , while the latter represents those with ongoing actions. The exogenous demand is similarly normalized as $\xi_t := \{\lambda_{v,t}, \Phi_{vv',t}\}$, where $\lambda_{v,t} := \frac{1}{n} \xi_{v,t}$ reflects the demand intensity at node v , and $\Phi_{vv',t} \in \Delta(\mathcal{V})$ is the empirical conditional distribution over destinations v' given an origin v , estimated from observed demands. When destination information is unavailable or unnecessary—such as in inventory systems where demand does not induce agent relocation and only local stock levels of nodes matter— Φ_t can be omitted. Based on this representation, joint actions a_t are sampled from a mean-field policy $\pi_t(a_t \mid \mu_t, \xi_{1:t})$, which maps the current mean-field state and demand history to control decisions for each node. This abstraction transforms the system high-dimensional actual state into low-dimensional mean-field state, which greatly simplifies the system as the number of agent grows.

3.2 Networked Control Problems with Exogenous Demand

We formulate the networked control problems with exogenous demand as a mean-field Exo-NMDP and further adopt a fluid modeling perspective (similar to the approach in Braverman et al. (2019)) to characterize the system macroscopic evolution through flow dynamics. Here, $\mu_t^v(i)$ denotes the idle agent density (or stock) at node i , $\mu_t^e(u, v)$ captures the in-transit flow of agents along edge (u, v) , while actions a_t represent newly initiated outflows. This fluid abstraction yields a tractable, low-dimensional system that supports scalable and robust control policy design.

Flow-Based Mean-Field Evolution. Two types of flows drive system evolution: demand flow and reposition flow. The demand flow from node i is defined as $f_{i,t}^D = \sum_j f_{ij,t}^D = \sum_j \alpha_{i,t} \lambda_{i,t} \Phi_{ij,t}$, where $\alpha_{i,t} \in [0, 1]$ is the fulfillment rate, $\lambda_{i,t}$ is the normalized demand intensity, and $\Phi_{ij,t}$ is the conditional probability distribution of demand from i being routed to j . Both $\lambda_{i,t}$ and $\Phi_{ij,t}$ are components of ξ_t . Depending on the system, fulfilled demand may either trigger agent relocation (e.g., in ride-sharing or delivery tasks), or remove agents from the system entirely (e.g., in inventory consumption). The reposition flow describes proactive agent movement based on a routing policy to balance supply and demand. It is defined as $f_{i,t}^R = \sum_j f_{ij,t}^R = \sum_j q_{ij,t} \cdot \mu_t^v(i)$, where $\mu_t^v(i)$ is the current density of idle agents at node i , and $q_{ij,t} \in \Delta(\mathcal{V})$ is the routing policy that specifies the probability of routing an idle agent from node i to j .

Our framework naturally extends to settings where agents have multiple types or internal states, indexed by a finite set \mathcal{K} . In this case, the mean-field state at each node i is now a vector $\mu_t^v(i) \in \mathbb{R}^{\mathcal{K}}$ where each element tracks the density of agents of that type at node i . To capture internal transitions between types (e.g., status changes), we define a conversion flow $f_{i,t}^C = C_{i,t} \cdot \mu_t^v(i)$, where $C_{i,t}$ is the type transition matrix at node i .

In practice, flows f^D , f^R , and f^C often incur delays. In other words, both agent type transitions and repositioning require some time, which we denote by a delay parameter τ . We incorporate this into a flow conservation model as follows:

$$\max \quad \mathbb{E} \left[\sum_{t=1}^T \mathcal{R}(\mu_t, \xi_t, a_t) \right] \quad (1)$$

$$\text{s.t.} \quad f_{i,t}^{\text{out}} = f_{i,t}^D + f_{i,t}^R + f_{i,t}^C \quad (2)$$

$$f_{i,t}^{\text{in}} = \sum_k (f_{ki,t-\tau}^D + f_{ki,t-\tau}^R) + C_{i,t-\tau} f_{i,t-\tau}^C \quad (3)$$

$$\mu_{i,t+1} = \mu_{i,t} + f_{i,t}^{\text{in}} - f_{i,t}^{\text{out}} \quad (4)$$

At each step t , we assume that demand-driven flows f^D take place first, then agent repositioning f^R , and finally local conversion f^C . Equation (2) represents the total outflow and Equation (3) captures the delayed inflow. The system's mean-field state evolves according to the flow conservation dynamics in Equation (4). The objective (1) tries to maximize system long-term revenue.

4 Method

The stochastic nature of exogenous demand renders direct solutions to (1)–(4) intractable. To address this, we propose the Dynamic Mean-Field Control (DMFC) framework for stepwise policy optimization. At each time t , DMFC first predicts a ideal mean-field state μ_t^* with historical information, and then solves for the policy π_t from a mean-field control problem (6) - (7) via linear programming. The overall closed-loop interaction of DMFC with the environment is presented in Algorithm 1.

Algorithm 1 Dynamic Mean-Field Control Loop with Environment Interaction

Require: network \mathcal{G} , time horizon T

```

1: for  $t = 1, \dots, T$  do
2:   Receive inflows  $f_t^{in}$  from external sources or previous steps (environment update)
3:   Observe endogenous system state  $s_t$  and exogenous demand  $\xi_t$ 
4:   Scale to the mean-field level: compute  $\mu_t$  and demand  $\xi_t = (\lambda_t, \Phi_t)$ 
5:   Match supply ( $\mu_t$ ) with demand ( $\xi_t$ ) for demand flow  $f_t^D$ 
6:   Forecast future demand  $\hat{\xi}_{t+1}$  based on historical observations  $\xi_{1:t}$ 
7:   Predict future supply  $\hat{\mu}_{t+1}$  by aggregating current state and delayed inflows
8:   Construct the ideal mean-field state  $\hat{\mu}_{t+1}^*$  by demand-proportional allocation
9:   Solve the mean-field control problem (6)–(7) with ideal mean-field state  $\hat{\mu}_{t+1}^*$  to obtain  $\pi_t$ 
10:  Execute policy  $\pi_t$  in the environment to generate flows  $f_t^R$  and  $f_t^C$ 
11: end for
    
```

4.1 The Ideal Mean-Field State

The ideal mean-field state $\hat{\mu}_{t+1}^*$ serves as a data-driven intermediate state that guides the stepwise policy optimization, which is constructed through the following steps. First, DMFC forecasts next-step demand $\hat{\xi}_{t+1}$ using a weighted history: $\hat{\xi}_{t+1} = \mathbf{w}_{1:t} \cdot \xi_{1:t}$, where $\mathbf{w}_{1:t}$ serves as learnable weight parameters controlling the influence of past observations $\xi_{1:t}$ on the prediction $\hat{\xi}_{t+1}$. To account for agent latency, projected future agent availability is estimated as $\hat{\mu}_{i,t+1} = \mu_{i,t} + \sum_{\tau=t+1}^{t+\Delta t} f_{i,\tau}^{in}$. The ideal mean-field state $\hat{\mu}_{t+1}^*$ is then constructed by demand-proportional allocation:

$$\hat{\mu}_{i,t+1}^* = \frac{\hat{\xi}_{i,t+1}}{\sum_k \hat{\xi}_{k,t+1}} \sum_k \hat{\mu}_{k,t+1}, \quad \forall i \in \mathcal{V} \quad (5)$$

The ideal mean-field state $\mu_{i,t}^*$ improves both short and long-term performance by balancing supply with demand. In the short term, the proportional allocation scheme ensures high demand fulfillment and immediate reward maximization. In the long term, despite flow delays, proportional allocation allows the system to gradually concentrate supply where demand is high, enhancing robustness and sustaining long-term efficiency.

4.2 Dynamic Mean-Field Control

With the ideal mean-field state $\hat{\mu}_{t+1}^*$ defined, we seek a control policy π_t that moves the system from its current state μ_t toward this target by solving:

$$\max_{a_t} \mathbb{E} [\mathcal{R}(\mu_t, \xi_t, a_t) - \mathbf{d}(\mu_{t+1}, \hat{\mu}_{t+1}^*)] \quad (6)$$

$$\text{s.t.} \quad \mu_{t+1} = \mu_t + f_t^{in} - f_t^{out} \quad (7)$$

Here, the reward function $\mathcal{R}(\mu_t, \xi_t, a_t)$ evaluates immediate performance through demand fulfillment, while the discrepancy term $\mathbf{d}(\mathbf{m}'_t, \mu_t^*)$ penalizes deviations from the forecasted ideal state. The constraint in (7) models the system dynamics via flow conservation. In practice, however, due to stochastic disturbances, delayed flows, and model approximation errors, the policy π_t may not perfectly align the evolved system state μ_{t+1} with the ideal state $\hat{\mu}_{t+1}^*$. Instead, the system tends to stabilize around the ideal trajectory, maintaining a dynamic equilibrium that supports robust performance across varying conditions.

5 Experimental evaluation

In this section, we implemented DMFC and conducted experiments in two representative scenarios: (1) **supply chain inventory management**, following the experimental setup of the GRL framework Gammelli et al. (2023); and (2) **vehicle routing in mobility-on-demand systems**, where we adopt the original GRL configuration and further introduce augmented datasets with increased demand variability.

5.1 Supply Chain Inventory Management (SCIM)

We consider the Supply Chain Inventory Management (SCIM) problem as a mean-field Exo-NMDP, enabling scalable optimization of commodity production, transportation, and inventory control under uncertainty. In this framework, the supply chain network is modeled as a directed graph $\mathcal{G} = (\mathcal{V}, \mathcal{E})$, comprising factories (\mathcal{V}_F) and stores (\mathcal{V}_S), where each node acts as a control agent managing aggregate flows.

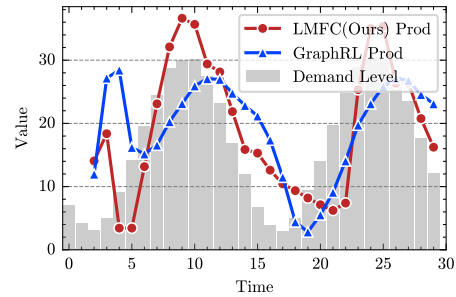
States, Demands, Actions, and Rewards. The system mean-field state μ_t captures the normalized inventory levels at nodes and in-transit flow on edges. At each time t , exogenous demand ξ_t arrives at store nodes $v \in \mathcal{V}_S$. If local inventory is available, demand is fulfilled with revenue p . Otherwise it remains pending and accumulates a delay penalty ϵ per time step. Each node select actions over feasible flow decisions. Factories control commodity production (conversion flow f_t^C) and outbound shipments (reposition flow f_t^R), incurring production cost m^P , transportation cost m^T , and delays t^P , t_{ij} . Stores passively fulfill demand (demand flow f_t^D) and accept incoming shipments. All nodes face storage constraints c_i , incur storage cost m^S , and are penalized by ϵ for overstocking. The policy $\pi_t(a_t|\mu_t, \xi_{1:t})$ maps currently observed inventory level μ_t and demand history $\xi_{1:t}$ to node-level flow actions. The reward function of the SCIM system is the total profit of the system, calculated as total revenue minus operational costs and penalties.

	DMFC	GRL	Improvement
1F2S	432 ± 193*	247 ± 110	+75%
1F3S	1671 ± 223*	875 ± 102	+91%
1F10S	4310 ± 358*	1244 ± 312	+246%

(a) Reward improvement of DMFC over GraphRL.

	DMFC	GRL	Reduction
1F2S	766*	1113	-31.2%
1F3S	676*	1227	-44.9%
1F10S	2278*	5080	-55.2%

(b) Storage cost reduction of DMFC over GraphRL.



(c) Production-demand comparison on 1F3S dataset: DMFC vs. GraphRL.

Figure 1: Performance comparison between DMFC algorithm and GraphRL baseline

Supply Chain Inventory Control. We apply the DMFC policy to the SCIM problem and compare it against the GraphRL baseline. All results are averaged over 30 runs, with detailed settings provided

in Appendix A.1.1. Tables 1a and 1b show that DMFC consistently outperforms GraphRL across all benchmark scenarios (1F2S, 1F3S, and 1F10S; F = factory, S = store). Specifically, DMFC improves total reward by 75–246% while reducing storage costs by 31–55%, without increasing penalties or degrading fulfillment performance. Figure 1c further illustrates DMFC’s advantage in production planning. On the 1F3S dataset, DMFC achieves better alignment between production and demand trends, resulting in smoother inventory turnover and lower storage costs. Results confirm that the mean-field modeling and lookahead scheduling enable more efficient resource utilization and improved responsiveness to exogenous demand fluctuations.

5.2 Vehicle Routing in Mobility-on-Demand Systems

Vehicle routing problems in mobility-on-demand systems (ECR) can also be formulated as mean-field Exo-NMDPs, where the goal is to optimize the repositioning of idle vehicles in a transportation network $\mathcal{G} = (\mathcal{V}, \mathcal{E})$ to maximize passenger demand fulfillment. Here, network nodes denote geographic regions and edges capture feasible inter-region routes. Each region is considered a central controller agent that governs the allocation and repositioning of idle vehicles within its area.

States, Demands, Actions, and Rewards. The mean-field state $\mu_t = (\mu_t^v(i), \mu_t^e(i, j))$ describes the normalized distribution of idle vehicles at each region $i \in \mathcal{V}$ and in-transit vehicles along each edge $(i, j) \in \mathcal{E}$. Exogenous demand is given by $\xi_t := (\lambda_t(v), \Phi_t(v, v'))$, where $\lambda_t(v)$ is the normalized arrival rate of passenger requests at region v , and $\Phi_t(v, v')$ represents the empirical distribution over destinations. When passengers arrive in region v , they are instantly matched to idle vehicles in the same region, generating a demand flow f^D . If supply falls short in that region, unmatched passengers leave the system. To reduce future imbalances, regions execute reposition actions, giving rise to the reposition flow f^R . Both f^D and f^R are subject to delays t_{ij} due to travel time. Upon arrival, vehicles automatically transition back to idle state, generating the conversion flow f^C . The reward is defined as total fare revenue minus fuel costs, and the objective is to optimize repositioning to maximize long-term returns.

Empty Car Coordination. We evaluate DMFC on four benchmarks (New York, Shenzhen, DiDi-9, and DiDi-20) and compare it against the GraphRL baseline. All results are averaged over 30 runs, with each episode consisting of 200 steps. The evaluation metric is the Order Response Rate (ORR) Lin et al. (2019), defined as the ratio of served requests to total requests. Detailed experimental settings can be found in Appendix A.2.1. Figure 2 illustrates the demand patterns for each dataset. To ensure consistent evaluation, all datasets have been adjusted to span 200 time slots. The red line represents the average passenger demand across regions, while the shaded area indicates the regional demand heterogeneity. Clear periodic trends are observed in the (extended) New York and Shenzhen datasets, whereas distinct diurnal patterns are evident in the DiDi-9 and DiDi-20 datasets. See Appendix A.2.2 for more information on the datasets.

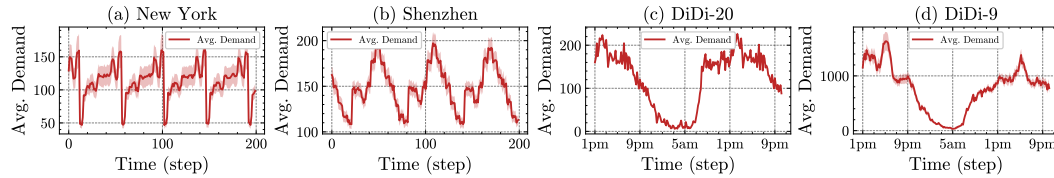


Figure 2: Demand pattern for each dataset in ECR environment

Figure 3a shows the ORR (Order Response Rate) performance on our (extended) New York and Shenzhen datasets with periodic demand. The GraphRL policy (GRL-Pretrain), which uses pre-trained weights directly from the original codebase without retraining, exhibits a clear performance degradation and limited adaptability across both cities. In contrast, DMFC maintains high and stable performance, effectively adapting to dynamic demand patterns in non-stationary environments. Figures 3b report the ORR under two pricing schemes: (i) fixed pricing (unit price, no fuel cost), and (ii) realistic pricing with or without fuel cost. DMFC consistently outperforms GraphRL across

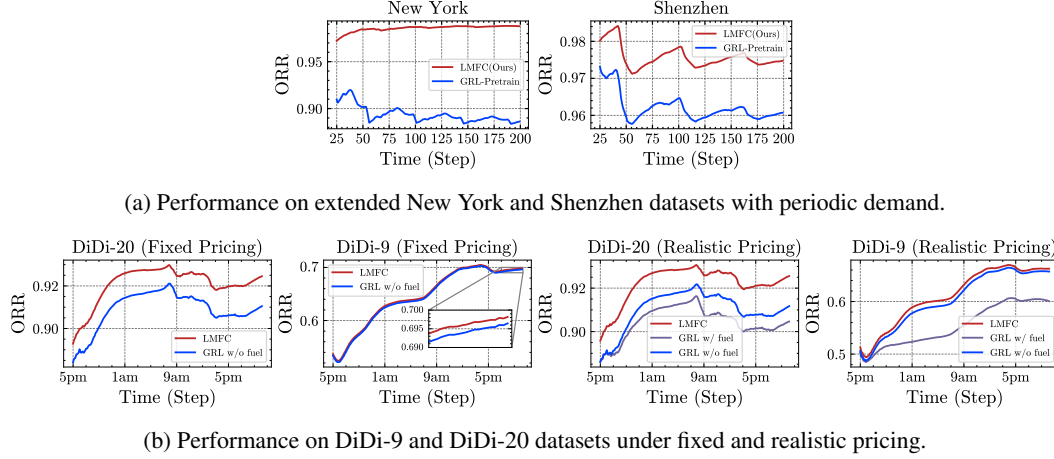


Figure 3: Order response rate (ORR) comparison across datasets and pricing schemes

all settings. Under fixed pricing, DMFC demonstrates a clear advantage on the DiDi-20 dataset and moderate improvements on DiDi-9. This is because DiDi-9 exhibits a more severe mismatch between supply and demand, where even marginal enhancements in routing yield noticeable benefits. The performance gap further widens under realistic pricing. While GraphRL tends to pursue short-term gains, often neglecting long-term system efficiency, DMFC emphasizes demand satisfaction, leading to more effective resource allocation and more stable performance under dynamic pricing and cost structures.

6 Conclusion

This paper proposes the Dynamic Mean-Field Control (DMFC) framework for optimizing resource allocation in network control problems with exogenous demands. Evaluations on supply chain inventory management and vehicle routing tasks demonstrate DMFC’s superiority over the GraphRL baseline, with enhanced adaptability to demand fluctuations and network complexity. Further, the framework achieves generalizability across diverse demand patterns (periodic, diurnal, sparse), robustness against supply-demand imbalances, and scalability to large networks via linear computational complexity.

References

- Ali Aalipour and Alireza Khani. Modeling, analysis, and control of autonomous mobility-on-demand systems: A discrete-time linear dynamical system and a model predictive control approach. *IEEE Transactions on Intelligent Transportation Systems*, 25(8):8615–8628, 2024. DOI: 10.1109/TITS.2024.3398562.
- Sajede Aminzadegan, Mohammad Tamannaie, and Morteza Rasti-Barzoki. Multi-agent supply chain scheduling problem by considering resource allocation and transportation. *Computers & Industrial Engineering*, 137:106003, 2019. ISSN 0360-8352. DOI: 10.1016/j.cie.2019.106003.
- Siddhartha Banerjee, Daniel Freund, and Thodoris Lykouris. Pricing and optimization in shared vehicle systems: An approximation framework. *Operations Research*, 70(3):1783–1805, 2022. DOI: 10.1287/opre.2021.2165.
- Forest Baskett, K. Mani Chandy, Richard R. Muntz, and Fernando G. Palacios. Open, closed, and mixed networks of queues with different classes of customers. *J. ACM*, 22(2):248–260, April 1975. ISSN 0004-5411. DOI: 10.1145/321879.321887.

- 306 Marcus Bellamy and Rahul Basole. Network analysis of supply chain systems: A systematic review
307 and future research. *Systems Engineering*, 16, 06 2013. DOI: 10.1002/sys.21238.
- 308 Francesco Borrelli, Alberto Bemporad, and Manfred Morari. *Predictive control for linear and hybrid*
309 *systems*. Cambridge University Press, 2017.
- 310 Anton Braverman, J. G. Dai, Xin Liu, and Lei Ying. Empty-car routing in ridesharing systems.
311 *Operations Research*, 67(5):1437–1452, 2019. DOI: 10.1287/opre.2018.1822.
- 312 Francesco Bullo, Emilio Frazzoli, Marco Pavone, Ketan Savla, and Stephen L. Smith. Dynamic
313 vehicle routing for robotic systems. *Proceedings of the IEEE*, 99(9):1482–1504, 2011. DOI:
314 10.1109/JPROC.2011.2158181.
- 315 Nicole Bäuerle. Mean field markov decision processes. *Applied Mathematics & Optimization*, 88,
316 04 2023. DOI: 10.1007/s00245-023-09985-1.
- 317 René Carmona, Mathieu Laurière, and Zongjun Tan. Model-free mean-field reinforcement learn-
318 ing: Mean-field MDP and mean-field Q-learning. *The Annals of Applied Probability*, 33(6B):
319 5334 – 5381, 2023. DOI: 10.1214/23-AAP1949. URL [https://doi.org/10.1214/](https://doi.org/10.1214/23-AAP1949)
320 [23-AAP1949](https://doi.org/10.1214/23-AAP1949).
- 321 DRI. http://research.xiaojukeji.com/index_en.html, 2016. Accessed: 2016-06-
322 30.
- 323 Daniele Gammelli, Kaidi Yang, James Harrison, Filipe Rodrigues, Francisco C. Pereira, and Marco
324 Pavone. Graph neural network reinforcement learning for autonomous mobility-on-demand sys-
325 tems. In *2021 60th IEEE Conference on Decision and Control (CDC)*, pp. 2996–3003, 2021.
326 DOI: 10.1109/CDC45484.2021.9683135.
- 327 Daniele Gammelli, James Harrison, Kaidi Yang, Marco Pavone, Filipe Rodrigues, and Francisco C.
328 Pereira. Graph reinforcement learning for network control via bi-level optimization. In Andreas
329 Krause, Emma Brunskill, Kyunghyun Cho, Barbara Engelhardt, Sivan Sabato, and Jonathan Scar-
330 lett (eds.), *Proceedings of the 40th International Conference on Machine Learning*, volume 202
331 of *Proceedings of Machine Learning Research*, pp. 10587–10610. PMLR, 23–29 Jul 2023.
- 332 Nicolas Gast, Bruno Gaujal, and Jean-Yves Le Boudec. Mean field for markov decision processes:
333 From discrete to continuous optimization. *IEEE Transactions on Automatic Control*, 57(9):2266–
334 2280, 2012. DOI: 10.1109/TAC.2012.2186176.
- 335 Mohammad Ghavamzadeh and Sridhar Mahadevan. Learning to communicate and act using hier-
336 archical reinforcement learning. In *Proceedings of the Third International Joint Conference on*
337 *Autonomous Agents and Multiagent Systems-Volume 3*, pp. 1114–1121. Citeseer, 2004.
- 338 Gurobi Optimization, LLC. Gurobi Optimizer Reference Manual, 2024. URL [https://www.](https://www.gurobi.com)
339 [gurobi.com](https://www.gurobi.com).
- 340 John Holler, Risto Vuorio, Zhiwei Qin, Xiaocheng Tang, Yan Jiao, Tiancheng Jin, Satinder Singh,
341 Chenxi Wang, and Jieping Ye. Deep reinforcement learning for multi-driver vehicle dispatching
342 and repositioning problem. In *2019 IEEE International Conference on Data Mining (ICDM)*, pp.
343 1090–1095, 2019. DOI: 10.1109/ICDM.2019.00129.
- 344 Ramon Iglesias, Federico Rossi, Rick Zhang, and Marco Pavone. A bcnp network approach to
345 modeling and controlling autonomous mobility-on-demand systems. *The International Journal*
346 *of Robotics Research*, 38, 07 2016. DOI: 10.1177/0278364918780335.
- 347 Ramon Iglesias, Federico Rossi, Kevin Wang, David Hallac, Jure Leskovec, and Marco Pavone.
348 Data-driven model predictive control of autonomous mobility-on-demand systems. In *2018 IEEE*
349 *International Conference on Robotics and Automation (ICRA)*, pp. 6019–6025, 2018. DOI: 10.
350 1109/ICRA.2018.8460966.

- 351 Matej Jusup, Barna Pásztor, Tadeusz Janik, Kenan Zhang, Francesco Corman, Andreas Krause, and
 352 Ilija Bogunovic. Safe model-based multi-agent mean-field reinforcement learning. In *Proceedings*
 353 *of the 23rd International Conference on Autonomous Agents and Multiagent Systems, AAMAS*
 354 '24, pp. 973–982, Richland, SC, 2024. International Foundation for Autonomous Agents and
 355 Multiagent Systems. ISBN 9798400704864.
- 356 Matej Jusup, Kenan Zhang, Zhiyuan Hu, Barna Pásztor, Andreas Krause, and Francesco Corman.
 357 Ride-sourcing vehicle rebalancing with service accessibility guarantees via constrained mean-
 358 field reinforcement learning. *CoRR*, abs/2503.24183, 2025. DOI: 10.48550/ARXIV.2503.24183.
 359 URL <https://doi.org/10.48550/arXiv.2503.24183>.
- 360 KDDCup. https://www.biendata.xyz/competition/kdd_didi/, 2020. Accessed:
 361 2020-07-12.
- 362 Basil Kouvaritakis and Mark Cannon. *Model Predictive Control: Classical, Robust and Stochastic*.
 363 Springer, 2016.
- 364 Pawel Ladosz, Lilian Weng, Minwoo Kim, and Hyondong Oh. Exploration in deep reinforcement
 365 learning: A survey. *Inf. Fusion*, 85:1–22, 2022. DOI: 10.1016/J.INFFUS.2022.03.003. URL
 366 <https://doi.org/10.1016/j.inffus.2022.03.003>.
- 367 Kaixiang Lin, Renyu Zhao, Zhe Xu, and Jiayu Zhou. Efficient collaborative multi-agent deep re-
 368 inforcement learning for large-scale fleet management, 2019. URL [https://arxiv.org/](https://arxiv.org/abs/1802.06444)
 369 [abs/1802.06444](https://arxiv.org/abs/1802.06444).
- 370 Chenxi Liu, Chao-Xiong Chen, and Chao Chen. Meta: A city-wide taxi repositioning framework
 371 based on multi-agent reinforcement learning. *IEEE Transactions on Intelligent Transportation*
 372 *Systems*, 23(8):13890–13895, 2022. DOI: 10.1109/TITS.2021.3096226.
- 373 Volodymyr Mnih, Koray Kavukcuoglu, David Silver, Alex Graves, Ioannis Antonoglou, Daan Wier-
 374 stra, and Martin Riedmiller. Playing atari with deep reinforcement learning, 2013.
- 375 Adam Paszke, Sam Gross, Francisco Massa, Adam Lerer, James Bradbury, Gregory Chanan, Trevor
 376 Killeen, Zeming Lin, Natalia Gimelshein, Luca Antiga, Alban Desmaison, Andreas Köpf, Ed-
 377 ward Z. Yang, Zachary DeVito, Martin Raison, Alykhan Tejani, Sasank Chilamkurthy, Benoit
 378 Steiner, Lu Fang, Junjie Bai, and Soumith Chintala. Pytorch: An imperative style, high-
 379 performance deep learning library. In Hanna M. Wallach, Hugo Larochelle, Alina Beygelzimer,
 380 Florence d’Alché-Buc, Emily B. Fox, and Roman Garnett (eds.), *Advances in Neural Informa-*
 381 *tion Processing Systems 32: Annual Conference on Neural Information Processing Systems 2019,*
 382 *NeurIPS 2019, December 8-14, 2019, Vancouver, BC, Canada*, pp. 8024–8035, 2019.
- 383 Barna Pásztor, Andreas Krause, and Ilija Bogunovic. Efficient model-based multi-agent mean-field
 384 reinforcement learning. *Trans. Mach. Learn. Res.*, 2023, 2023. URL [https://openreview.](https://openreview.net/forum?id=gvcDSDYUZx)
 385 [net/forum?id=gvcDSDYUZx](https://openreview.net/forum?id=gvcDSDYUZx).
- 386 Rick Zhang;Marco Pavone;. Control of robotic mobility-on-demand systems: A queueing-
 387 theoretical perspective. *The International Journal of Robotics Research*, 35(1-3):186–203, 2016.
 388 DOI: 10.1177/0278364915581863.
- 389 Zhiwei Tony Qin, Hongtu Zhu, and Jieping Ye. Reinforcement learning for ridesharing: An extended
 390 survey. *Transportation Research Part C: Emerging Technologies*, 144:103852, 2022.
- 391 Sheldon M Ross. *Applied probability models with optimization applications*. Dover Publications,
 392 1992.
- 393 Marco Pavone;Stephen L Smith;Emilio Frazzoli;Daniela Rus;. Robotic load balancing for mobility-
 394 on-demand systems. *The International Journal of Robotics Research*, 31(7):839–854, 2012. DOI:
 395 10.1177/0278364912444766.

- 396 John F Shortle, James M Thompson, Donald Gross, and Carl M Harris. *Fundamentals of queueing*
397 *theory*. John Wiley & Sons, 2018.
- 398 Sean R. Sinclair, Felipe Frujeri, Ching-An Cheng, Luke Marshall, Hugo Barbalho, Jingling Li,
399 Jennifer Neville, Ishai Menache, and Adith Swaminathan. Hindsight learning for mdps with
400 exogenous inputs, 2023. URL <https://arxiv.org/abs/2207.06272>.
- 401 Richard S. Sutton and Andrew G. Barto. *Reinforcement Learning: An Introduction*. A Bradford
402 Book, Cambridge, MA, USA, 2018. ISBN 0262039249.
- 403 Richard S. Sutton, Doina Precup, and Satinder Singh. Between mdps and semi-mdps: A framework
404 for temporal abstraction in reinforcement learning. *Artificial Intelligence*, 112(1):181–211, 1999.
405 ISSN 0004-3702. DOI: [https://doi.org/10.1016/S0004-3702\(99\)00052-1](https://doi.org/10.1016/S0004-3702(99)00052-1). URL <https://www.sciencedirect.com/science/article/pii/S0004370299000521>.
- 407 Matthew Tsao, Dejan Milojevic, Claudio Ruch, Mauro Salazar, Emilio Frazzoli, and Marco Pavone.
408 Model predictive control of ride-sharing autonomous mobility-on-demand systems. In *2019*
409 *International Conference on Robotics and Automation (ICRA)*, pp. 6665–6671, 2019. DOI:
410 10.1109/ICRA.2019.8794194.
- 411 Jingwei Wang, Qianyu Hao, Wenzhen Huang, Xiaochen Fan, Zhentao Tang, Bin Wang, Jianye
412 Hao, and Yong Li. Dyps: Dynamic parameter sharing in multi-agent reinforcement learning for
413 spatio-temporal resource allocation. In *Proceedings of the 30th ACM SIGKDD Conference on*
414 *Knowledge Discovery and Data Mining*, KDD ’24, pp. 3128–3139, New York, NY, USA, 2024.
415 Association for Computing Machinery. ISBN 9798400704901.
- 416 Honghao Wei, Zixian Yang, Xin Liu, Zhiwei Qin, Xiaocheng Tang, and Lei Ying. A reinforcement
417 learning and prediction-based lookahead policy for vehicle repositioning in online ride-hailing
418 systems. *IEEE Transactions on Intelligent Transportation Systems*, 25(2):1846–1856, 2024. DOI:
419 10.1109/TITS.2023.3312048.
- 420 Dacheng Wen, Yupeng Li, and Francis C. M. Lau. A survey of machine learning-based ride-hailing
421 planning. *IEEE Transactions on Intelligent Transportation Systems*, 25(6):4734–4753, 2024.
422 DOI: 10.1109/TITS.2023.3345174.
- 423 Zhi-Hai Zhang, Bin-Feng Li, Xiang Qian, and Lin-Ning Cai. An integrated supply chain network de-
424 sign problem for bidirectional flows. *Expert Systems with Applications*, 41(9):4298–4308, 2014.
425 ISSN 0957-4174. DOI: <https://doi.org/10.1016/j.eswa.2013.12.053>.

A Additional Experiment Details

This section provides further details about the experimental configuration and hyperparameters. All RL modules are taken from [Gammelli et al. \(2023\)](#) and implemented using Pytorch [Paszke et al. \(2019\)](#). The Gurobi Optimizer [Gurobi Optimization, LLC \(2024\)](#) is used for optimization problems. The code environment has been rewritten to facilitate more detailed and adaptable customization options. To ensure consistent behavior across environments in edge cases, some IBM CPLEX components from the original codebase of [Gammelli et al. \(2023\)](#) have been retained within our framework.

A.1 Supply Chain Inventory Management

A.1.1 Environment Details and Datasets

We follow the basic setting of [Gammelli et al. \(2023\)](#) and define the exogenous demand pattern λ_t as a co-sinusoidal function with a stochastic component:

$$\lambda_{i,t} = \left\lfloor \frac{\lambda_i^{\max}}{2} \left(1 + \cos \left(\frac{4\pi(2i+t)}{T} \right) + \mathcal{U}(0, \lambda_i^{\text{var}}) \right) \right\rfloor \quad (8)$$

where $\lfloor \cdot \rfloor$ is the floor function, λ_i^{\max} is the maximum demand value $\mathcal{U}(0, \lambda_i^{\text{var}})$ is a uniform distribution on the interval $[0, \lambda_i^{\text{var}}]$, and T is the episode length. Hyperparameters in simulation experiments are borrowed from [Gammelli et al. \(2023\)](#) and listed below (with minor modifications):

A.1.2 Mean-Field Control

In a supply chain inventory management problem, commodities are modeled as agents whose distributions evolve over time. At each time step t , the system state s_t comprises current inventory levels at factories and stores, along with pending demands at each store node that have accumulated due to prior stockouts. Given the system state and exogenous demand ξ_t , we define the mean-field state μ_t and the demand vector ξ_t , both normalized by the total number of commodities in the network.

The total commodity volume is not static, as items are continuously consumed (sold at stores) and replenished (produced by factories). To ensure a consistent mean-field representation, we incorporate demand-driven production into the scaling process. Specifically, we forecast the next-step demand, account for current inventories and backlog levels, and infer the appropriate production quantity to maintain system balance. This enables us to construct a normalized and dynamically adjusted mean-field state that reflects both supply and anticipated demand over time.

The mean-field control formulation of a SCIM problem is defined as follows:

$$\max_{f^R, f^C, \epsilon^V, \epsilon^S} \min_{i \in \mathcal{V}_S} \alpha_i - M \sum_{i \in \mathcal{V}} |\epsilon_i^V| - \sum_{i \in \mathcal{V}} |\epsilon_i^S| \quad (9)$$

$$\text{s.t.} \quad \mu_{i,t+1} = \mu_{i,t} - f_{i,t}^D + \sum_{k \in \mathcal{V}_F} f_{ki,t}^R, \quad \forall i \in \mathcal{V}_S \quad (10)$$

$$\mu_{i,t+1} = \mu_{i,t} - \sum_{j \in \mathcal{V}_S} f_{ij,t}^R + f_{i,t}^C, \quad \forall i \in \mathcal{V}_F \quad (11)$$

$$\hat{\mu}_{i,t+1}^* = \mu_{i,t+1} + \epsilon_i^S, \quad \forall i \in \mathcal{V} \quad (12)$$

$$\mu_{i,t+1} \leq V_i + \epsilon_i^V, \quad \forall i \in \mathcal{V} \quad (13)$$

The objective (9) adopts a minimax structure, where $\min_{i \in \mathcal{V}_S} \alpha_i$ denotes the worst-case demand fulfillment rate across all store nodes, analogous to the reward function $\mathcal{R}(\mu_t, \xi_t, a_t)$ in (6). The penalty terms $\sum_i |\epsilon_i^S|$ and $M \sum_i |\epsilon_i^V|$ serve two purposes: the former captures the deviation between the evolved mean-field state $\mu_t + 1$ and the target state $\hat{\mu}_{t+1}^*$ (as defined in (12)), while the latter penalizes inventory overflow beyond node capacity V_i , as constrained in (13).

Parameter	Explanation	Value
λ^{\max}	Maximum demand	[2, 16]
λ^{var}	Demand variance	[2, 2]
T	Episode length	30
t^P	Production time	1
t_{ij}	Travel time	[1, 1]
c	Storage capacity	[20, 9, 12]
m^P	Production cost	5
m^S	Storage cost	[3, 2, 1]
m^T	Transportation cost	[0.3, 0.6]
p	Price	15
ϵ	Penalty	21

(a) Parameters for the 1F2S environment

Parameter	Explanation	Value
λ^{\max}	Maximum demand	[1, 5, 24]
λ^{var}	Demand variance	[2, 2, 2]
T	Episode length	30
t^P	Production time	1
t_{ij}	Travel time	[1, 1, 1]
c	Storage capacity	[30, 15, 15, 15]
m^P	Production cost	5
m^S	Storage cost	[2, 1, 1, 1]
m^T	Transportation cost	[0.3, 0.3, 0.3]
p	Price	15
ϵ	Penalty	21

(b) Parameters for the 1F3S environment

Parameter	Explanation	Value
λ^{\max}	Maximum demand	[2, 2, 2, 2, 10, 10, 10, 18, 18, 18]
λ^{var}	Demand variance	[2] $_{i \in \mathcal{V}}$
T	Episode length	30
t^P	Production time	1
t_{ij}	Travel time	[1] $_{i \in \mathcal{V}}$
c	Storage capacity	[100, 15 $\forall i \in \mathcal{V} \setminus 0$]
m^P	Production cost	5
m^S	Storage cost	[1, 2 $\forall i \in \mathcal{V} \setminus 0$]
m^T	Transportation cost	[0.3] $_{i \in \mathcal{V}}$
p	Price	15
ϵ	Penalty	21

(c) Parameters for the 1F10S environment

Figure 4: Parameter settings for different SCIM environments

Constraints (10) and (11) refine the general flow conservation dynamics in (4) for the supply chain inventory management (SCIM) setting. Specifically, $f_{i,t}^D$ represents the quantity of demand fulfilled at store node i , $f_{ij,t}^R$ indicates the fraction of commodities transferred from factory i to store j , and $f_{i,t}^C$ denotes the amount of production at factory node i .

A.2 Vehicle Routing in Mobility-on-Demand Systems

A.2.1 Environment Details

We follow the experimental setup of Gammelli et al. (2023) and construct our environments accordingly. To ensure consistent and comprehensive evaluation, all experiments are conducted on episodes consisting of 200 time steps. We model vehicles within mobility-on-demand systems as agents. Exogenous demand is then modeled as a Poisson arrival process, with intensity specified for each time step. Agent repositioning and scheduling delays are determined by the datasets and may be either constant or time-varying. Note that the geographical regions in the environment are not necessarily fully connected.

A.2.2 Datasets

We conduct experiments on four benchmarks: the **New York** and **Shenzhen** datasets from Gammelli et al. (2023) (3-hour demand windows), the **DiDi-9** dataset from DRI (2016) used by Braverman et al. (2019) (21-day order records), and the **DiDi-20** dataset from KDDCup (2020) used by Wei et al. (2024) (a full-day order data). Table 1 gives an overview of each datasets. For preprocessing, we extended the New York and Shenzhen datasets to 200 timesteps. For the DiDi-9 and DiDi-20 datasets, we extracted contiguous 200-timestep segments spanning 1 PM to 10 PM across consecutive days. Figure 2 illustrates the preprocessed demand patterns, showing periodic trends in the New York/Shenzhen data and diurnal cycles in the DiDi-9/DiDi-20 datasets. The red line denotes average demand, with shading indicating regional variations.

	NYC	SZ	DiDi-9	DiDi-20
# of regions	14	17	9	20
# of cars	1200	1200	1500	500
minutes per step	4	3	10	10

Table 1: Configurations of ECR simulation experiments over all datasets

Each dataset contains comprehensive records of travel times, origin–destination pairs, pricing, and other relevant information. From these records, we extract key time-slot-specific parameters to model the dynamics of the traffic network, including the demand pattern ξ_t , inter-regional travel time t_{ij} , and the fuel cost coefficient β .

A.2.3 Mean-Field Control

When considering vehicle routing problems in a mobility-on-demand system, we consider the real-world as inter-connected geographical regions, and model each vehicle as an agent. At each time step t , the system mean-field state μ_t represents the density of agents across each region. The ideal mean-field state of the system is defined in proportional to anticipated next-step demand of each

491 region. We define the corresponding mean-field control problem as follows:

$$\min_{f^R} \sum_{(i,j) \in \mathcal{E}} t_{ij} \cdot f_{ij}^R + M \cdot \epsilon_i \quad (14)$$

$$\text{s.t.} \quad \mu_{i,t} + \sum_{k \neq i} f_{ki}^R - \sum_{j \neq i} f_{ij}^R + \epsilon_i \geq \hat{\mu}_{i,t+1}^*, \quad \forall i \in \mathcal{V} \quad (15)$$

$$\sum_{j \neq i} f_{ij}^R \leq \mu_{i,t} \quad \forall i \in \mathcal{V} \quad (16)$$

492 The objective in (14) seeks to minimize repositioning time, thereby reducing the fuel cost associated
 493 with vehicle routing. This objective is alternative to our overarching goal of maximizing long-term
 494 system revenue as defined in (1). The auxiliary variable ϵ captures deviations between the actual
 495 system mean-field state μ_{t+1} and the target mean-field state $\hat{\mu}_{t+1}^*$, with M denoting a large penalty
 496 coefficient. Constraint (15) ensures that the system evolves toward the desired target state, while
 497 constraint (16) enforces that the repositioning outflow from region i does not exceed the number of
 498 available empty vehicles.