

Safe Model-Based Multi-Agent Mean-Field Reinforcement Learning

Matej Jusup
ETH Zurich
Zurich, Switzerland
mhusup@ethz.ch

Kenan Zhang
EPFL Lausanne
Lausanne, Switzerland
kenan.zhang@epfl.ch

Barna Pásztor
ETH Zurich
Zurich, Switzerland
barna.pasztor@ai.ethz.ch

Francesco Corman
ETH Zurich
Zurich, Switzerland
corman@ethz.ch

Tadeusz Janik
ETH Zurich
Zurich, Switzerland
tjanik@student.ethz.ch

Andreas Krause
ETH Zurich
Zurich, Switzerland
krausea@ethz.ch

Ilija Bogunovic
University College London
London, United Kingdom
i.bogunovic@ucl.ac.uk

ABSTRACT

Many applications, e.g., in shared mobility, require coordinating a large number of agents. Mean-field reinforcement learning addresses the resulting scalability challenge by optimizing the policy of a representative agent interacting with the infinite population of identical agents instead of considering individual pairwise interactions. In this paper, we address an important generalization where there exist global constraints on the distribution of agents (e.g., requiring capacity constraints or minimum coverage requirements to be met). We propose SAFE-M³-UCRL, the first model-based mean-field reinforcement learning algorithm that attains safe policies even in the case of *unknown* transitions. As a key ingredient, it uses epistemic uncertainty in the transition model within a log-barrier approach to ensure pessimistic constraints satisfaction with high probability. Beyond the synthetic swarm motion benchmark, we showcase SAFE-M³-UCRL on the vehicle repositioning problem faced by many shared mobility operators and evaluate its performance through simulations built on vehicle trajectory data from a service provider in Shenzhen. Our algorithm effectively meets the demand in critical areas while ensuring service accessibility in regions with low demand.

KEYWORDS

Multi-agent reinforcement learning; Mean-field control; Global safety; Epistemic uncertainty; Probabilistic neural network ensemble; Shared mobility; Vehicle repositioning

ACM Reference Format:

Matej Jusup, Barna Pásztor, Tadeusz Janik, Kenan Zhang, Francesco Corman, Andreas Krause, and Ilija Bogunovic. 2024. Safe Model-Based Multi-Agent

Mean-Field Reinforcement Learning. In *Proc. of the 23rd International Conference on Autonomous Agents and Multiagent Systems (AAMAS 2024)*, Auckland, New Zealand, May 6 – 10, 2024, IFAAMAS, 10 pages.

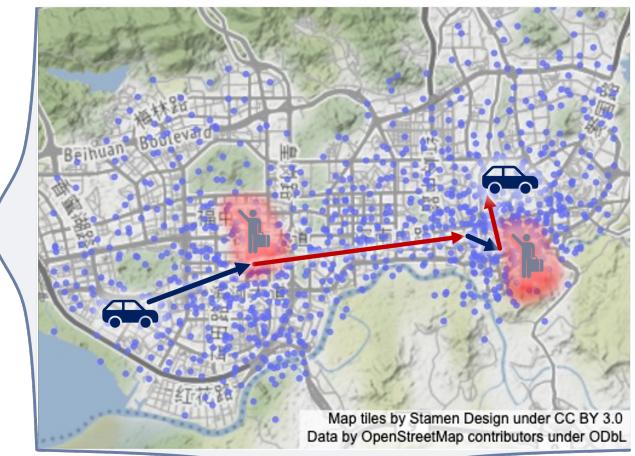


Figure 1: An illustration of vehicles' spatial distribution (light-blue scatters), mean-field marginal distributions (left and bottom curves), repositioning trips (blue arrows), and a trajectory of passenger trips (red arrows).

1 INTRODUCTION

Multi-Agent Reinforcement Learning (MARL) is a rapidly growing field that seeks to understand and optimize the behavior of multiple agents interacting in a shared environment. MARL has a wide range of potential applications, including vehicle repositioning in shared mobility services (e.g., moving idle vehicles from low-demand to high-demand areas [41]), swarm robotics (e.g., operating a swarm of drones [2]), and smart grids (e.g., operating a network of sensors in electric system [52]). The interactions between agents in these complex systems introduce several challenges, including non-stationarity, scalability, competing learning goals, and varying information structure. Mean-Field Control (MFC) addresses



This work is licensed under a Creative Commons Attribution International 4.0 License.

the scalability and non-stationarity hurdles associated with MARL by exploiting the insight that many relevant MARL problems involve a large number of very similar agents working towards the same goal. Instead of focusing on the individual agents and their interactions, MFC considers an asymptotically large population of identical cooperative agents and models them as a distribution on the state space. This approach circumvents the problem’s dependency on the population size, enabling the consideration of large populations. The solutions obtained by MFC are often sufficient for the finite-agent equivalent problem [13, 36, 48, 61] in spite of the introduced approximations. An example of such a system is a ride-hailing platform that serves on-demand trips with a fleet of vehicles. The platform needs to proactively reposition idle vehicles based on their current locations, the locations of the other vehicles in the fleet, and future demand patterns (see Figure 1) to maximize the number of fulfilled trips and minimize customer waiting times. Additionally, the platform may be obligated by external regulators to guarantee service accessibility across the entire service region. The problem quickly becomes intractable as the number of vehicles increases. A further difficulty lies in modeling the traffic flows. Due to numerous infrastructure, external, and driver behavioral factors, which are often region-specific, it is laborious and often difficult to determine transitions precisely [8, 18, 60].

In this paper, we focus on learning the *safe optimal policies* for a large multi-agent system when the underlying transitions are *unknown*. In most real-world systems, the transitions must be learned from the data obtained from repeated interactions with the environment. We assume that the cost of obtaining data from the environment is high and seek to design a model-based solution that efficiently uses the collected data. Existing works consider solving the MFC problem via model-free or model-based methods without safety guarantees. However, the proposed *Safe Model-Based Multi-Agent Mean-Field Upper-Confidence Reinforcement Learning* (SAFE-M³-UCRL) algorithm focuses on learning underlying transitions and deriving optimal policies for the mean-field setting while avoiding undesired or dangerous distributions of agents’ population.

Contributions. Section 3 extends the MFC setting with safety constraints and defines a novel comprehensive notion of global population-based safety. To address safety-constrained environments featuring a large population of agents, in Section 4, we propose a model-based mean-field reinforcement learning (MFRL) algorithm called SAFE-M³-UCRL. Our algorithm leverages epistemic uncertainty in the transition model, employing a log-barrier approach to guarantee pessimistic satisfaction of safety constraints and enables the derivation of safe policies. In Section 5, we conduct empirical testing on the synthetic swarm motion benchmark and real-world vehicle repositioning problem, a challenge commonly faced by shared mobility operators. Our results demonstrate that the policies derived using SAFE-M³-UCRL successfully fulfill demand in critical demand hotspots while ensuring service accessibility in areas with lower demand.

2 RELATED WORK

Our notion of safety for the mean-field problem extends the frameworks of *Mean-Field Games* (MFG) and *Mean-Field Control* (MFC) [30, 31, 38, 39]. For a summary of the progress focusing on MFGs, see [40] and references therein. We focus on MFCs in this work,

which assume cooperative agents in contrast to MFGs, which assume competition. [6, 24, 26, 27, 49] address the problem of solving MFCs under known transitions, i.e., planning, while [4, 5, 10–12, 58, 61, 62] consider model-free Q-learning and Policy Gradient methods in various settings. Closest to our approach, [50] introduces M³-UCRL, a model-based, on-policy algorithm, which is more sample efficient than other proposed approaches. Similarly to [15, 17, 33] for model-based single-agent RL and [53] for model-based MARL, M³-UCRL uses the epistemic uncertainty in the transition model to design optimistic policies that efficiently balance exploration and exploitation in the environment and maximize sample efficiency. This is also the setting of our interest. However, safety is not considered in any of these methods.

There are two main ways of handling *safety* in RL; assigning significantly lower rewards to unsafe states [46] and providing additional knowledge to the agents [56] or using the notion of controllability to avoid unsafe policies explicitly [25]. The following approaches combine the two methods; [7] uses Lyapunov functions to restrict the safe policy space, [14] projects unsafe policies to a safe set via a control barrier function, and [3] introduces shielding, i.e., correcting actions only if they lead to unsafe states. For comprehensive overviews on safe RL, we refer the reader to [23, 29]. As an alternative, [59] demonstrates that the general-purpose stochastic optimization methods can be used for constrained MDPs, i.e., safe RL formulations. Similar to our work, they use the log-barrier approach to turn constrained into unconstrained optimization. Nevertheless, the aforementioned works focus mainly on individual agents, while in large-scale multi-agent environments, maintaining individual safety becomes intractable, and the focus shifts towards global safety measures. For multi-agent problems, previous works focus on satisfying the individual constraints of the agents while learning in a multi-agent environment. For the cooperative problem, [28] proposes two model-free algorithms, MACPO and MAPPO-Lagrangian. MACPO is computationally expensive, while MAPPO-Lagrangian does not guarantee hard constraints for safety. Dec-PG solves the decentralized learning problem using a consensus network that shares weights between neighboring agents [44]. For the non-cooperative decentralized MARL problem with individual constraints, [55] adds a safety layer to multi-agent DDPG [43] similar to single-agent Safe DDPG [19] for continuous state-action spaces. Aggregated and population-based constraints have been addressed in the following works. CMIX [42] extends QMIX [51], which considers average and peak constraints defined over the whole population of agents in a centralized-learning decentralized-execution framework. Their formulation relies on the joint state and action spaces, making it infeasible for a large population of agents. [22] introduces an additional shielding layer that corrects unsafe actions. Their centralized approach suffers from scalability issues, while the factorized shielding method monitors only a subset of the state or action space. For mixed cooperative-competitive settings, [66] uses the notion of returnability to define a safe, decentralized, multi-agent version of Q-learning that ensures individual and joint constraints. However, their approach requires an estimation of other agents’ policies, which does not scale well for large systems. Works considering constraints on the whole population fail to overcome the exponential nature of multi-agent problems or require domain knowledge to factorize the problems into subsets.

Closest to our setting, [48] introduces constraints to the MFC by defining a cost function and a threshold that the discounted sum of costs can not exceed. We propose a different formulation that restricts the set of feasible mean-field distributions at every step, therefore, addressing the scalability issue and allowing for more specific control over constraints and safe population distributions.

3 PROBLEM STATEMENT

Formally, we consider the *episodic* setting, where episodes $n = 1, \dots, N$ each have $t = 0, \dots, T - 1$ discrete steps and the terminal step $t = T$. The state space $\mathcal{S} \subseteq \mathbb{R}^p$ and action space $\mathcal{A} \subseteq \mathbb{R}^q$ are the same for every agent. We use $s_{n,t}^{(i)} \in \mathcal{S}$ and $a_{n,t}^{(i)} \in \mathcal{A}$, to denote the state and action of agent $i \in \{1, \dots, m\}$ in episode n at step t . For every n and t , the *mean-field distribution* $\mu_{n,t} \in \mathcal{P}(\mathcal{S})$ describes the global state with m identical agents when $m \rightarrow +\infty$, i.e.,

$$\mu_{n,t}(ds) = \lim_{m \rightarrow \infty} \frac{1}{m} \sum_{i=1}^m \mathbb{I}(s_{n,t}^{(i)} \in ds),$$

where $\mathbb{I}(\cdot)$ is the indicator function, and $\mathcal{P}(\mathcal{S})$ is the set of probability measures over the state space \mathcal{S} .

We consider the MFC model to capture a collective behavior of a large number of *collaborative* agents operating in the shared *stochastic environment*. This model assumes the limiting regime of *infinitely* many agents and *homogeneity*. Namely, all agents are identical and indistinguishable, therefore, solving MFC amounts to finding an optimal policy for a single, so-called, *representative agent*. The representative agent interacts with the mean-field distribution of agents instead of focusing on individual interactions and optimizes a collective reward. Due to the homogeneity assumption, the representative agent's policy is used to control all the agents in the environment.

We posit that the reward $r : \mathcal{S} \times \mathcal{P}(\mathcal{S}) \times \mathcal{A} \rightarrow \mathbb{R}$ of the representative agent is known and that it depends on the states of the other agents through the mean-field distribution.¹ Before every episode n , the representative agent selects a non-stationary policy profile $\pi_n = (\pi_{n,0}, \dots, \pi_{n,T-1}) \in \Pi$ where individual policies are of the form $\pi_{n,t} : \mathcal{S} \times \mathcal{P}(\mathcal{S}) \rightarrow \mathcal{A}$ and Π is the set of admissible policy profiles. The policy profile is then shared with all the agents that choose their actions according to π_n during episode n .

We consider a general family of deterministic transitions $f : \mathcal{S} \times \mathcal{P}(\mathcal{S}) \times \mathcal{A} \rightarrow \mathcal{S}$. Given the current mean-field distribution $\mu_{n,t}$, the representative agent's state $s_{n,t}$ and its action $a_{n,t}$, the next representative agent's state $s_{n,t+1}$ is given by

$$s_{n,t+1} = f(s_{n,t}, \mu_{n,t}, a_{n,t}) + \varepsilon_{n,t}, \quad (1)$$

where $\varepsilon_{n,t}$ is a Gaussian noise with known variance. We assume that the transitions are *unknown* and are to be inferred from collected trajectories across episodes.

Mean-field transitions. State-to-state transition map in Equation (1) naturally extends to the *mean-field transitions* induced by a policy profile π_n and transitions f in episode n (see [50, Lemma 1])

$$\mu_{n,t+1}(ds') = \int_S \mathbb{P}[s_{n,t+1} \in ds'] \mu_{n,t}(ds), \quad (2)$$

where $s_{n,t+1}$ is the next representative agent state and $\mu_{n,t}(ds) = \mathbb{P}[s_{n,t} \in ds]$ under π_n for all t . To simplify the notation, we use $U(\cdot)$

¹Our framework easily extends to unknown reward by estimating its epistemic uncertainty and learning it similarly to learning unknown transitions (see [15]).

to denote the mean-field transition function from Equation (2), i.e., we have $\mu_{n,t+1} = U(\mu_{n,t}, \pi_{n,t}, f)$. We further introduce the notation $z_{n,t} \in \mathcal{Z} = \mathcal{S} \times \mathcal{P}(\mathcal{S}) \times \mathcal{A}$ to denote the tuple $(s_{n,t}, \mu_{n,t}, a_{n,t})$.

For a given policy profile π_n and mean-field distribution μ , the performance of the representative agent at step t is measured via its expected future reward for the rest of the episode, i.e.,

$$\mathbb{E} \left[\sum_{j=t}^{T-1} r(z_{n,j}) | \mu_{n,t} = \mu \right].$$

Here, the expectation is over the randomness in the transitions and over the sampling of the initial state, i.e., $s_{n,t} \sim \mu$.

3.1 Safe Mean-Field Reinforcement Learning

We extend the MFC with global safety constraints,² i.e., the constraints imposed on the mean-field distributions. We consider safety functions $h : \mathcal{P}(\mathcal{S}) \rightarrow \mathbb{R}$ over probability distributions. Given some hard safety threshold $C \in \mathbb{R}$, we consider a mean-field distribution μ as safe if it satisfies $h(\mu) \geq C$, or, equivalently

$$h_C(\mu) := h(\mu) - C \geq 0. \quad (3)$$

We denote the set of safe mean-field distributions for a safety constraint $h_C(\cdot)$ as $\zeta = \{\mu \in \mathcal{P}(\mathcal{S}) : h_C(\mu) \geq 0\}$. Hence, our focus is on the safety of the system as a whole rather than the safety of individual agents, as it becomes intractable to handle individual agents' states and interactions in the case of a large population.

For a given initial distribution μ_0 , we formally define the *Safe-MFC*³ problem as follows

$$\pi^* = \arg \max_{\pi \in \Pi} \mathbb{E} \left[\sum_{t=0}^{T-1} r(z_t) | \mu_0 \right] \quad (4a)$$

$$\text{subject to } a_t = \pi_t(s_t, \mu_t) \quad (4b)$$

$$s_{t+1} = f(z_t) + \varepsilon_t \quad (4c)$$

$$\mu_{t+1} = U(\mu_t, \pi_t, f) \quad (4d)$$

$$h_C(\mu_{t+1}) \geq 0, \quad (4e)$$

where we explicitly require induced mean-field distributions $\{\mu_t\}_{t=1}^T$ to reside in the safe set ζ by restricting the set of admissible policy profiles Π to policy profiles that induce safe distributions. To ensure complete safety, we note that the initial mean-field distribution μ_0 must be in the safe set ζ as our learning protocol does not induce it (see Algorithm 1).

We make the following assumptions about the environment using Wasserstein 1-distance defined by

$$W_1(\mu, \mu') := \inf_{\gamma \in \Gamma(\mu, \mu')} \mathbb{E}_{(x, y) \sim \gamma} \|x - y\|_1,$$

where $\Gamma(\mu, \mu')$ is the set of all couplings of μ and μ' (i.e., a joint probability distributions with marginals μ and μ'). We further define the distance between $z = (s, \mu, a)$ and $z' = (s', \mu', a')$ as

$$d(z, z') := \|s - s'\|_2 + \|a - a'\|_2 + W_1(\mu, \mu').$$

²For the exposition, we use a single constraint, however, our approach is directly extendable to multiple constraints.

³We refer to formulations under known transitions as control problems, while we reserve the term reinforcement learning for formulations under unknown transitions.

ASSUMPTION 1 (TRANSITIONS LIPSCHITZ CONTINUITY). *The transition function $f(\cdot)$ is L_f -Lipschitz-continuous, i.e.,*

$$\|f(z) - f(z')\|_2 \leq L_f d(z, z').$$

ASSUMPTION 2 (MEAN-FIELD POLICIES LIPSCHITZ CONTINUITY). *The individual policies π present in any admissible policy profile $\boldsymbol{\pi}$ in Π are L_π -Lipschitz-continuous, i.e.,*

$$\|\pi(s, \mu) - \pi(s', \mu')\|_2 \leq L_\pi (\|s - s'\|_2 + W_1(\mu, \mu'))$$

for all $\pi \in \boldsymbol{\pi} \in \Pi$.

ASSUMPTION 3 (REWARD LIPSCHITZ CONTINUITY). *The reward function $r(\cdot)$ is L_r -Lipschitz-continuous, i.e.,*

$$\|r(z) - r(z')\|_2 \leq L_r d(z, z').$$

These assumptions are considered standard in model-based learning [15, 33, 50, 53] and mild, as individual policies and rewards are typically designed such that they meet these smoothness requirements. For example, we use neural networks with Lipschitz-continuous activations to represent our policies.

3.2 Examples of Safety Constraints

We can model multiple classes of safety constraints $h_C(\cdot) \geq 0$ that naturally appear in real-world applications such as vehicle repositioning, traffic flow, congestion control, and others.

Entropic safety. Entropic constraints can be used in multi-agent systems to prevent overcrowding by promoting spatial diversity and avoiding excessive clustering. Incorporating an entropic term in the decision-making process encourages the controller to distribute the agents evenly within the state space. This might be particularly useful in applications that include crowd behavior, such as operating a swarm of drones or a fleet of vehicles. In such scenarios, we define safety by imposing a threshold $C \geq 0$ on the differential entropy

$$H(\mu) := - \int_S \log \mu(s) \mu(ds) \quad (5)$$

of the mean-field distribution μ , i.e.,

$$h_C(\mu) := H(\mu) - C.$$

Distribution similarity. Another way to define safety is by preventing μ from diverging from a prior distribution v_0 . The prior can be based on previous studies, expert opinions or regulatory requirements. We can use a penalty function that quantifies the allowed dissimilarity between the two distributions

$$h_C(\mu; v_0) := C - D(\mu, v_0),$$

with $C \geq 0$ and where the distance function between probability measures $D : \mathcal{P}(S) \times \mathcal{P}(S) \rightarrow \mathbb{R}_{\geq 0}$ depends on the problem at hand.

We provide further examples of safety functions in Appendix B [35] together with proofs that they satisfy Assumption 6.

3.3 Statistical Model and Safety Implications

The representative agent learns about unknown transitions by interacting with the environment. We take a model-based approach to achieve sample efficiency by sequentially updating and improving the transition model estimates based on the previously observed transitions. At the beginning of each episode n , the representative agent updates its model based on $\cup_{i=1}^{n-1} \mathcal{D}_i$ where $\mathcal{D}_i =$

$\{(z_{i,t}, s_{i,t+1})\}_{t=0}^{T-1}$ and $z_{i,t} = (s_{i,t}, \mu_{i,t}, a_{i,t})$ is the set of observations in episode i for $i = 1, \dots, n-1$, i.e., up until the beginning of episode n . We estimate the mean $\mathbf{m}_{n-1} : \mathcal{Z} \rightarrow \mathcal{S}$ and covariance $\Sigma_{n-1} : \mathcal{Z} \rightarrow \mathbb{R}^{p \times p}$ functions from the set of collected trajectories $\cup_{i=1}^{n-1} \mathcal{D}_i$, and denote model's confidence with $\sigma_{n-1}^2(z) = \text{diag}(\Sigma_{n-1}(z))$. We assume that the statistical model is calibrated, meaning that at the beginning of every episode, the agent has *high probability confidence bounds* around unknown transitions. The following assumptions are consistent with [15, 17, 50, 53, 57] and other literature which aims to exclude extreme functionals from consideration.

ASSUMPTION 4 (CALIBRATED MODEL). *Let $\mathbf{m}_{n-1}(\cdot)$ and $\Sigma_{n-1}(\cdot)$ be the mean and covariance functions of the statistical model of f conditioned on $n-1$ observed episodes. For the confidence function $\sigma_{n-1}(\cdot)$, there exists a non-decreasing, strictly positive sequence $\{\beta_n\}_{n \geq 0}$ such that for $\delta > 0$ with probability at least $1 - \delta$, we have jointly for all $n \geq 1$ and $z \in \mathcal{Z}$ that $|f(z) - \mathbf{m}_{n-1}(z)| \leq \beta_{n-1} \sigma_{n-1}(z)$ elementwise.*

ASSUMPTION 5 (ESTIMATED CONFIDENCE LIPSCHITZ CONTINUITY). *The confidence function $\sigma_n(\cdot)$ is L_σ -Lipschitz-continuous for all $n \geq 0$, i.e., $\|\sigma_n(z) - \sigma_n(z')\|_2 \leq L_\sigma d(z, z')$.*

Since the true transition model is unknown in Equation (4c) and Equation (4d), at the beginning of every episode n , the representative agent can only construct the confidence set of transitions \mathcal{F}_{n-1} with $\mathbf{m}_{n-1}(\cdot)$ and $\sigma_{n-1}(\cdot)$ estimated based on the observations up until the end of the previous episode $n-1$, i.e.,

$$\mathcal{F}_{n-1} = \left\{ \tilde{f} : \tilde{f} \text{ is calibrated w.r.t. } \mathbf{m}_{n-1}(\cdot) \text{ and } \sigma_{n-1}(\cdot) \right\} \quad (6)$$

The crucial challenge in *Safe-MFRL* is that the representative agent can only select transitions $\tilde{f} \in \mathcal{F}_{n-1}$ at the beginning of the episode n and use it instead of true transitions f when solving Equation (4) to find an optimal policy profile $\boldsymbol{\pi}_n^*$. The resulting mean-field distributions $\{\tilde{\mu}_t\}_{t=1}^T$ are then different from $\{\mu_t\}_{t=1}^T$ (i.e., the ones that correspond to the true transition model), and hence the constraint Equation (3) guarantees only the safety under the estimated transitions \tilde{f} , i.e., $h_C(\tilde{\mu}_t) \geq 0$. In contrast, the original environment constraint $h_C(\mu_t) \geq 0$ might be violated, resulting in unsafe mean-field distributions under true transitions f .

Next, we demonstrate how to modify the constraint Equation (4e) for the optimization problem Equation (4) when an estimated transition function \tilde{f} is used from the confidence set \mathcal{F}_{n-1} such that the mean-field distributions $\mu_{n,t}$ induced by the resulting policy profile $\boldsymbol{\pi}_n^*$ do not violate the original constraint under true transitions f . First, we require the following property for any safety function $h(\cdot)$.

ASSUMPTION 6 (SAFETY LIPSCHITZ CONTINUITY). *The safety function $h(\cdot)$ is L_h -Lipschitz-continuous, i.e., $|h(\mu) - h(\mu')| \leq L_h W_1(\mu, \mu')$.*

The following lemma shows that we can ensure safety under true transitions f by having tighter constraints under any estimated transitions \tilde{f} selected from \mathcal{F}_{n-1} .

LEMMA 1. *Given a fixed policy profile $\boldsymbol{\pi}_n$, a safety function $h(\cdot)$ satisfying Assumption 6 and a safety threshold $C \in \mathbb{R}$, we have in episode n for all steps*

$$|h_C(\tilde{\mu}_{n,t}) - h_C(\mu_{n,t})| \leq L_h C_{n,t},$$

where $C_{n,t}$ is an arbitrary constant that satisfies $C_{n,t} \geq W_1(\tilde{\mu}_{n,t}, \mu_{n,t})$.

PROOF. For arbitrary $\tilde{\mu}_{n,t}, \mu_{n,t} \in \mathcal{P}(\mathcal{S})$ we have

$$|h_C(\tilde{\mu}_{n,t}) - h_C(\mu_{n,t})| = |h(\tilde{\mu}_{n,t}) - h(\mu_{n,t})| \leq L_h W_1(\tilde{\mu}_{n,t}, \mu_{n,t}) \leq L_h C_{n,t},$$

where the first equality follows from the definition of $h_C(\cdot)$, the first inequality follows from Assumption 6 and the second inequality comes from $C_{n,t} \geq W_1(\tilde{\mu}_{n,t}, \mu_{n,t})$. \square

Crucially, using Lemma 1 we can formulate a safety constraint for the optimization under estimated transitions \tilde{f} that ensures that the constraint under true transitions f is satisfied with high probability.

COROLLARY 1. *For every episode n and step t , $h_C(\tilde{\mu}_{n,t}) \geq L_h C_{n,t}$ implies $h_C(\mu_{n,t}) \geq 0$ guaranteeing the safety of the original system.*

PROOF. The corollary follows directly from Lemma 1 and the triangle inequality, which are used in the third and the second inequality, respectively

$$\begin{aligned} L_h C_{n,t} &\leq h_C(\tilde{\mu}_{n,t}) \\ &\leq |h_C(\tilde{\mu}_{n,t}) - h_C(\mu_{n,t})| + h_C(\mu_{n,t}) \\ &\leq L_h C_{n,t} + h_C(\mu_{n,t}). \end{aligned}$$

The claim is obtained by subtracting the positive constant $L_h C_{n,t}$ from both sides. \square

Then, $C_{n,t}$ for $t = 1, \dots, T$ become parameters of the optimization problem (as defined in Section 4) that the representative agent faces at the beginning of episode n . However, choosing the appropriate values that comply with the condition $C_{n,t} \geq W_1(\tilde{\mu}_{n,t}, \mu_{n,t})$ is not trivial since $\mu_{n,t}$ depends on unknown true transitions of the system. Note that computing $C_{n,0}$ at the initial step $t = 0$ is not necessary because the inequality is always guaranteed due to the initialization $\tilde{\mu}_{n,0} = \mu_{n,0}$ for every episode n . In Appendix A [35], we demonstrate how to efficiently upper bound $W_1(\tilde{\mu}_{n,t}, \mu_{n,t})$ and obtain $C_{n,t}$ using the Lipschitz constants of the system and the statistical model's epistemic uncertainty. In particular, $C_{n,t}$ approaches zero, and $h_C(\tilde{\mu}_{n,t}) \geq L_h C_{n,t}$ reduces to the constraint Equation (4e) as the estimated confidence $\sigma_{n-1}(\cdot)$ shrinks due to the increasing number of observations available to estimate true transitions.

4 SAFE-M³-UCRL

In this section, we introduce a model-based approach for the *Safe-MFRL* problem that combines the safety guarantees in Corollary 1 with upper-bound confidence interval optimization. At the beginning of each episode n , the representative agent constructs the confidence set of transitions \mathcal{F}_{n-1} (see Equation (6)) given the calibrated statistical model and previously observed data and selects a safe *optimistic* policy profile π_n^* to obtain the highest value function within \mathcal{F}_{n-1} while satisfying the safety constraint derived in Corollary 1. In particular, the optimal policy profile π^* from Equation (4) is approximated at the episode n by

$$\pi_n^* = \arg \max_{\pi_n \in \Pi} \max_{\tilde{f}_{n-1} \in \mathcal{F}_{n-1}} \mathbb{E} \left[\sum_{t=0}^{T-1} r(\tilde{z}_{n,t}) \middle| \tilde{\mu}_{n,0} = \mu_0 \right] \quad (7a)$$

$$\text{subject to } \tilde{a}_{n,t} = \pi_{n,t}(\tilde{s}_{n,t}, \tilde{\mu}_{n,t}) \quad (7b)$$

$$\tilde{s}_{n,t+1} = \tilde{f}_{n-1}(\tilde{z}_{n,t}) + \varepsilon_{n,t} \quad (7c)$$

$$\tilde{\mu}_{n,t+1} = U(\tilde{\mu}_{n,t}, \pi_{n,t}, \tilde{f}_{n-1}) \quad (7d)$$

$$h_C(\tilde{\mu}_{n,t+1}) \geq L_h C_{n,t+1}, \quad (7e)$$

with $\tilde{z}_{n,t} = (\tilde{s}_{n,t}, \tilde{\mu}_{n,t}, \tilde{a}_{n,t})$. Equation (7) optimizes over the function space \mathcal{F}_{n-1} which is usually intractable even in bandit settings [20]. Additionally, it must comply with the safety constraint Equation (7e), further complicating the optimization. We utilize the *hallucinated control* reparametrization and the *log-barrier* method to alleviate these issues. After the reformulation of the problem, model-free or model-based mean-field optimization algorithms can be applied to find policy profile π_n^* at the beginning of episode n .

We use an established approach known as *Hallucinated Upper Confidence Reinforcement Learning* (H-UCRL) [17, 47, 50] and introduce an auxiliary function $\eta : \mathcal{Z} \rightarrow [-1, 1]^p$, where p is the dimensionality of the state space \mathcal{S} , to define hallucinated transitions

$$\tilde{f}_{n-1}(z) = \mathbf{m}_{n-1}(z) + \beta_{n-1} \Sigma_{n-1}(z) \eta(z). \quad (8)$$

Notice that \tilde{f}_{n-1} is calibrated for any $\eta(\cdot)$ under Assumption 4, i.e., $\tilde{f}_{n-1} \in \mathcal{F}_{n-1}$. Assumption 4 further guarantees that every function \tilde{f}_{n-1} can be expressed in the auxiliary form Equation (8)

$$\forall \tilde{f}_{n-1} \in \mathcal{F}_{n-1} \exists \eta : \mathcal{Z} \rightarrow [-1, 1]^p \text{ such that}$$

$$\tilde{f}_{n-1}(z) = \mathbf{m}_{n-1}(z) + \beta_{n-1} \Sigma_{n-1}(z) \eta(z), \forall z \in \mathcal{Z}.$$

Thus, the intractable optimization over the function space \mathcal{F}_{n-1} in Equation (7) can be expressed as an optimization over the set of admissible policy profiles Π and auxiliary function $\eta(\cdot)$. Note that $\eta(z) = \eta(s, \mu, \pi(s, \mu)) = \eta(s, \mu)$ for a fixed individual policy π . This turns $\eta(\cdot)$ into a policy that exerts *hallucinated control* over the epistemic uncertainty of the confidence set of transitions \mathcal{F}_{n-1} [17]. Furthermore, Equation (8) allows us to optimize over parametrizable functions (e.g., *neural networks*) π and $\eta(\cdot)$ using gradient ascent.

We introduce the safety constraint to the objective using the *log-barrier method* [64]. This restricts the domain on which the objective function is defined only to values that satisfy the constraint Equation (7e), hence, turning Equation (7) to an unconstrained optimization problem. Combining these two methods yields the following optimization problem

$$\pi_n^* = \arg \max_{\pi_n \in \Pi} \max_{\eta(\cdot) \in [-1, 1]^p} \mathbb{E} \left[\sum_{t=0}^{T-1} r(\tilde{z}_{n,t}) + \lambda \log(h_C(\tilde{\mu}_{n,t+1}) - L_h C_{n,t+1}) \middle| \tilde{\mu}_{n,0} = \mu_0 \right] \quad (9a)$$

$$\text{subject to } \tilde{a}_{n,t} = \pi_{n,t}(\tilde{s}_{n,t}, \tilde{\mu}_{n,t}) \quad (9b)$$

$$\tilde{f}_{n-1}(\tilde{z}_{n,t}) = \mathbf{m}_{n-1}(\tilde{z}_{n,t}) + \beta_{n-1} \Sigma_{n-1}(\tilde{z}_{n,t}) \eta(\tilde{z}_{n,t}) \quad (9c)$$

$$\tilde{s}_{n,t+1} = \tilde{f}_{n-1}(\tilde{z}_{n,t}) + \varepsilon_{n,t} \quad (9d)$$

$$\tilde{\mu}_{n,t+1} = U(\tilde{\mu}_{n,t}, \pi_{n,t}, \tilde{f}_{n-1}), \quad (9e)$$

with $\tilde{z}_{n,t} = (\tilde{s}_{n,t}, \tilde{\mu}_{n,t}, \tilde{a}_{n,t})$ and $\lambda > 0$ being a tuneable hyperparameter used to balance between the reward and the safety constraint. Provided that the set of safe mean-field distributions (assuming the safe initial distribution μ_0 is given) is not empty, π_n^* is guaranteed to satisfy the safety constraint during the policy rollout in episode n .

REMARK 1. *Note that Equation (9) can also be used under known transitions by setting $\mathbf{m}_{n-1}(\cdot) = f(\cdot)$, $\Sigma_{n-1}(\cdot) = 0$ and $L_h C_{n,t} = 0$, hence, recovering the original constraint $h_C(\tilde{\mu}_{n,t}) \geq 0$ from Equation (3). In Section 5, we utilize this useful property to construct the upper bound for the reward obtained under unknown transitions.*

Algorithm 1 Model-Based Learning Protocol in SAFE-M³-UCRL

Input: Set of admissible policy profiles Π , safety constraint $h_C(\cdot)$, calibrated statistical model represented by $\mathbf{m}_{n-1}(\cdot)$ and $\Sigma_{n-1}(\cdot)$, initial mean-field distribution μ_0 , known reward $r(\cdot)$, safety Lipschitz constant L_h , hyperparameter λ , number of episodes N , number of steps T

- 1: **for** $n = 1, \dots, N$ **do**
- 2: Compute $C_{n,t}$ for $t = 1, \dots, T$ as described in Appendix A [35]
- 3: Optimize the objective in Equation (9) over the admissible policy profiles Π and hallucinated transitions Equation (8)
- 4: Execute the obtained policy profile π_n^* and collect the trajectories $\mathcal{D}_n = \{(z_{n,t}, s_{n,t+1})\}_{t=0}^{T-1}$ from the representative agent
- 5: Update the confidence set of transitions \mathcal{F}_{n-1} with the collected data to obtain \mathcal{F}_n for the next episode
- 6: **end for**

Return $\pi_N^* = (\pi_{N,0}^*, \dots, \pi_{N,T-1}^*)$

We summarize the model-based learning protocol used by SAFE-M³-UCRL in Algorithm 1. The first step computes constants $C_{n,t}$ introduced in Lemma 1. The second step optimizes the objective in Equation (9). The third and fourth steps collect trajectories from the representative agent and update the calibrated model. While the learning protocol is model-based, the subroutine in Line 3 can use either model-based or model-free algorithms proposed for the MFC due to our reformulation in Equation (9). In Appendix C.3 [35], we introduce modifications of well-known algorithms for optimizing the mean-field setting.

5 EXPERIMENTS

In this section, we demonstrate the performance of SAFE-M³-UCRL on the swarm motion benchmark and showcase that it can tackle the real-world large-scale vehicle repositioning problem faced by ride-hailing platforms.

5.1 Swarm Motion

Due to the infancy of MFRL as a research topic, one of the rare benchmarks used by multiple authors is the swarm motion. [12, 50] view it as MFRL problem, while [21] uses it in the context of MFGs. In this setting, an infinite population of agents is moving around toroidal state space with the aim of maximizing a location-dependent reward function while avoiding congested areas [1].

Modeling. We model the state space \mathcal{S} as the unit torus on the interval $[0, 1]$, and the action space is the interval $\mathcal{A} = [-7, 7]$. We approximate the continuous-time swarm motion by partitioning unit time into $T = 100$ equal steps of length $\Delta t = 1/T$. The next state $s_{n,t+1} = f(z_{n,t}) + \varepsilon_{n,t}$ is induced by the unknown transitions $f(z_{n,t}) = s_{n,t} + a_{n,t}\Delta t$ with $\varepsilon_{n,t} \sim N(0, \Delta t)$. The reward function is defined by $r(z_{n,t}) = \phi(s_{n,t}) - \frac{1}{2}a_{n,t}^2 - \log(\mu_{n,t})$, where the first term $\phi(s) = 2\pi^2(\sin(2\pi s) - \cos^2(2\pi s)) + 2\sin(2\pi s)$ determines the positional reward received at the state s , the second term defines the kinetic energy penalizing large actions, and the last term penalizes overcrowding. Note that the optimal solution for continuous time setting, $\Delta t \rightarrow 0$ can be obtained analytically [1] and used as a benchmark. [12, 50] show that MFRL discrete-time, $\Delta t > 0$, methods can learn good approximations of the optimal solution. The disadvantage of these methods is that they can influence the skewness of the mean-field distribution only via overcrowding penalty. Therefore, to control skewness, their only option is to introduce a hyperparameter to the reward that regulates the level of overcrowding penalization. On the other hand, SAFE-M³-UCRL controls skewness without trial-and-error reward shaping by imposing the

entropic safety constraint $h_C(\mu_{n,t}) = H(\mu) - C \geq 0$, with $H(\cdot)$ defined in Equation (5), instead of having the overcrowding penalty term $\log(\mu_{n,t})$. Since higher entropy translates into less overcrowding, we can upfront determine and upper-bound the acceptable level of overcrowding by setting a desirable threshold C .

We use a neural network to parametrize the policy profile $\pi_n = (\pi_{n,0}, \dots, \pi_{n,T-1})$ during the optimization of Equation (9). The optimization is done by *Mean-Field Back-Propagation Through Time* (MF-BPTT) (see Appendix C.3.1 [35]). In our experiments, a single neural network shows enough predictive power to represent the whole policy profile, but using T networks, one for each individual policy $\pi_{n,t}$, $t = 0, \dots, T-1$, is a natural extension. We use a *Probabilistic Neural Network Ensemble* [16, 37] to represent a statistical model of transitions, which we elaborate in Appendix C.1 [35]. We represent the mean-field distribution by discretizing the state space uniformly and assigning the probability of the representative agent residing within each interval. We set the safety threshold C as a proportion $p \in [0, 1]$ of the maximum entropy, i.e., $C = p \max_{\mathcal{P}(\mathcal{S})} H(\mu)$. Note that SAFE-M³-UCRL guarantees safe mean-field distributions only if the initial mean-field distributions $\mu_{n,0}$ at time $t = 0$ are safe for every episode n for given threshold C . A generic way for safe initialization is setting $\mu_{n,0}$ as the maximum entropy distribution among all safe distributions ζ

$$\mu_{n,0} = \arg \max_{\mu \in \zeta} H(\mu). \quad (10)$$

Note that, in general, the safe initial distribution might not exist.

Results. In Figure 2a, we observe the learning curve of SAFE-M³-UCRL for $p = 0.95$ for 10 randomly initialized runs. The learning process is volatile in the initial phase due to the high epistemic uncertainty, but after 50 episodes, all policies converge toward the solution as if the transitions were known. In Figure 2b, we use various thresholds C to show that the entropic constraint effectively influences the degree of agents' greediness to collect the highest positional reward. By increasing p towards 1, we force agents to put increasingly more emphasis on global welfare rather than on individual rewards. We see that for $p = 0.5$ we obtain a distribution that matches the distribution obtained by the unconstrained M³-UCRL [50] that relies on the overcrowding penalty, while for $p = 0.95$ we significantly surpass the effect that the penalty has on the distribution's skewness. We also show that the discrete-time solutions with low p serve as a good approximation of the continuous-time optimal distribution μ^* . Furthermore, we observe that the policies under unknown transitions overlap with the solutions learned under known transitions. In Figure 2e, we observe

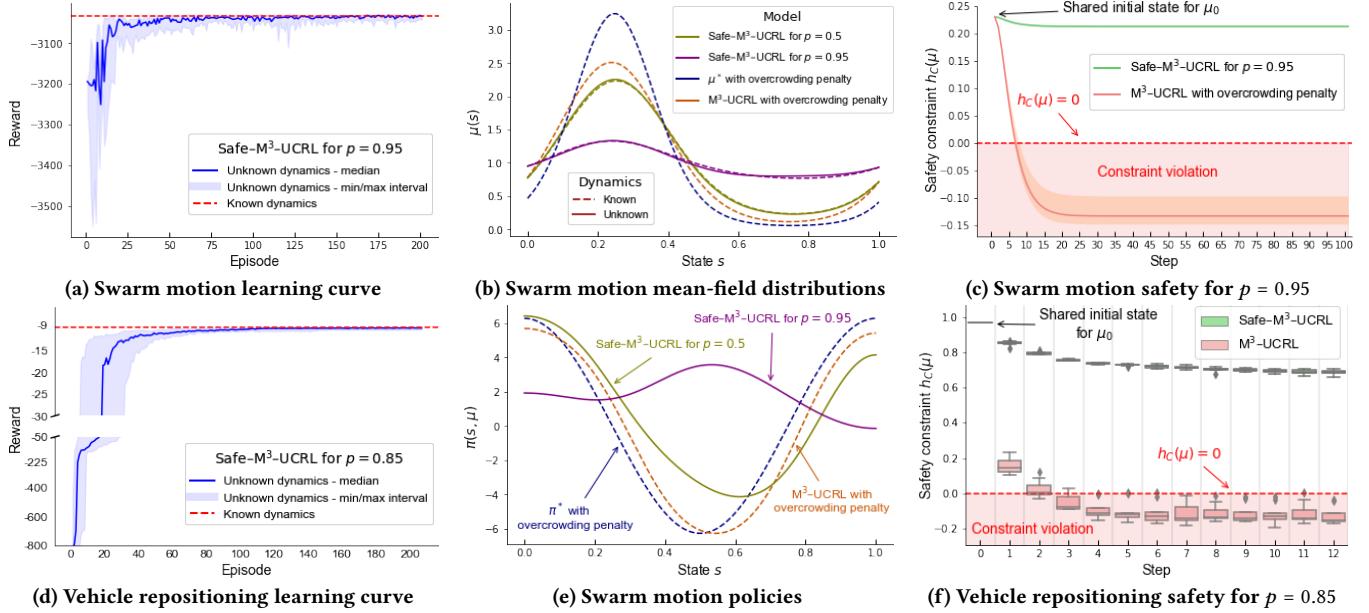


Figure 2: Performance analysis of **SAFE-M³-UCRL** for swarm motion and vehicle repositioning. The policy and statistical model were trained on 10 randomly initialized neural networks for each hyperparameter p .

that unconstrained policies and policies with low p push agents towards high individual rewards. Note that for states close to 1, the algorithms learn it requires less kinetic energy to push agents over the border due to the toroidal shape of the state space. For high p 's learned policies push agents always in the same direction to maintain uniformity of the system. Importantly, Figure 2c shows that **SAFE-M³-UCRL** for $p = 0.95$ keeps the mean-field distribution safe throughout the entire execution, unlike the solutions that rely on the overcrowding penalty term.

5.2 Vehicle Repositioning Problem

Since ride-hailing services, such as Uber, Lyft, and Bolt, gained market share, vehicle repositioning has been a long-standing challenge for these platforms, i.e., moving idle vehicles to areas with high demand. A similar challenge is present in bike-sharing services in many cities and, as of more recently, dockless electric scooter-sharing services such as Bird and Lime. The operator significantly increases profit in the competitive environment by successfully repositioning idle vehicles to high-demand areas. Nevertheless, there might exist regulations imposed by the countries or cities that enforce service providers to either guarantee fair service accessibility or restrict the number of vehicles in districts with high traffic density. Such restrictions prevent operators from greedily maximizing the profit and can be encapsulated by some dispersion metric such as entropy. Solving this problem helps prevent prolonged vehicle cruising and extensive passenger waiting times in the demand hotspots, increasing the service provider's efficiency and reducing its carbon footprint. Existing approaches to vehicle repositioning range from static optimization over a queuing network [9, 65], model predictive control [32], to RL [41, 45, 63]. The main advantage of **SAFE-M³-UCRL** is the capability of controlling a large fleet of homogeneous vehicles and enforcing efficient coordination to match the spatiotemporal demand distribution. Additionally, the

safety constraint introduced into the model guarantees service accessibility by ensuring idle vehicles are spreading over the study region. Although accessibility has not been widely discussed in the literature on vehicle repositioning, it is expected to be an important fairness constraint when shared mobility services become a prevailing travel mode [54].

Modeling. Ride-hailing operations can be modeled as sequential decision-making, which consists of passenger trips followed by repositioning trips operated by a central controller as illustrated in Figure 1. We assume that the controller has access to the locations of vehicles in its fleet and communicates the real-time repositioning actions to the drivers via electronic devices. Nevertheless, since the fleet is operating in a noisy traffic environment, repositioning usually cannot be executed perfectly. We assume that vehicles can move freely within the area of our interest, which is represented by a two-dimensional unit square, i.e., the state-space $\mathcal{S} = [0, 1]^2$, and repositioning actions are taken from $\mathcal{A} = [-1, 1]^2$. The objective of our model is to satisfy the demand in the central district of Shenzhen while providing service accessibility in the wider city center. We restrict our modeling horizon to three evening peak hours, which are discretized in fifteen-minute operational intervals, i.e., $T = 12$, and each episode n represents one day. We model service providers goal of maximizing the coverage of the demand by the negative of the Kullback-Leibler divergence between the vehicles' distribution $\mu_{n,t}$ and demand for service denoted as ρ_0 , i.e., $r(z_{n,t}) = -D_{KL}(\rho_0 || \mu_{n,t})$. In particular, the demand distribution $\rho_0 \in \mathcal{P}(\mathcal{S})$ represents a probability of a trip originating in the infinitesimal neighborhood of state $s \in \mathcal{S}$ during peak hours (see Figure 3a). We estimate a stationary demand distribution ρ_0 from the vehicle trajectories collected in Shenzhen, China, in 2016. If the passenger's trip originates at $s \in \mathcal{S}$ the likelihood of its destinations is defined by the mapping $\Phi : \mathcal{S} \rightarrow \mathcal{P}(\mathcal{S})$, which we fit from the trip trajectories. We use $\Phi(\cdot)$ to define sequential transitions by

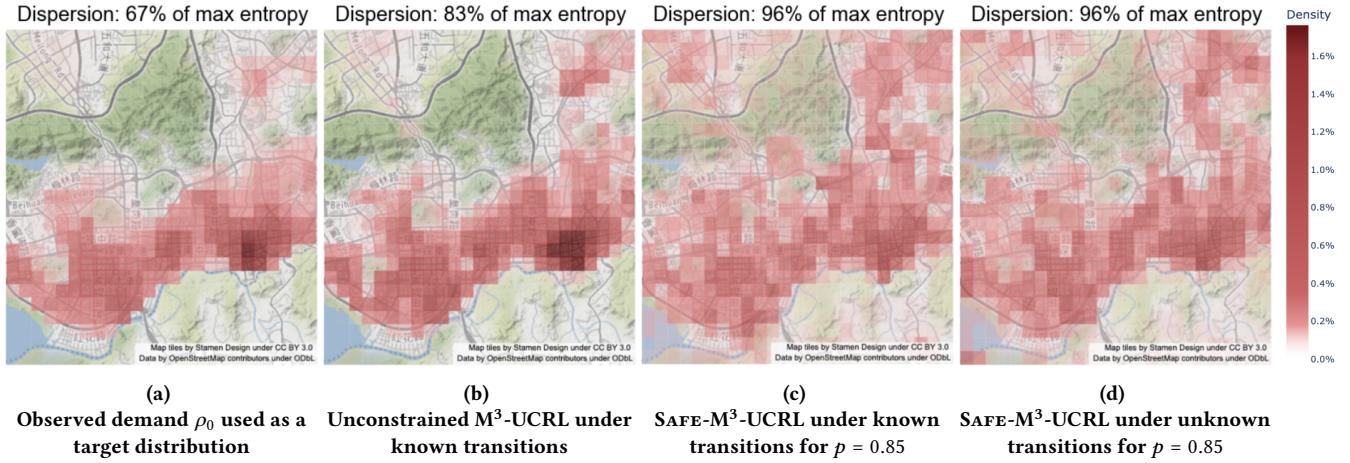


Figure 3: SAFE-M³-UCRL guided vehicle distribution in Shenzhen in the evening peak hours.

first executing passenger trips followed by vehicle repositioning. Formally, the next state $s_{n,t+1} = f(z_{n,t}) + \varepsilon_{n,t}$ is induced by the unknown transitions $f(z_{n,t}) = \text{clip}(s_{n,t}^\Phi + a_{n,t}, 0, 1)$, where $s_{n,t}^\Phi \sim \Phi(s_{n,t})$ and $\varepsilon_{n,t} \sim \text{TN}(0, \sigma^2 I_2)$ is a Gaussian with a known variance σ^2 truncated at the borders of \mathcal{S} and I_2 is the 2×2 unit matrix. Notice that the controller determines repositioning actions given intermediate states $s_{n,t}^\Phi$ obtained after executing passenger trips. We use entropic safety constraint $h_C(\mu_{n,t}) = H(\mu) - C \geq 0$ to enforce the service accessibility across all residential areas (see Figures 3c to 3d). Consequently, the optimization objective in Equation (9) trades off between greedily satisfying the demand ρ_0 and adhering to accessibility constraint imposed by $h_C(\cdot)$. Identically to the swarm motion experiment, we use a neural network to parametrize the policy profile π_n , which we optimize by MF-BPTT. A statistical model of the transitions is represented by a Probabilistic Neural Network Ensemble, while $\mu_{n,0}$ is initialized using Equation (10). We represent the mean-field distribution by discretizing the state space into the uniform grid.

Results. The entropy of the target distribution, ρ_0 in Figure 3a, already achieves $p = 0.67$ of the maximum due to a wide horizontal spread. To enforce vertical spread, we require an additional 18 percentage points of entropy as a safety constraint. Concretely, we use $p = 0.85$ to set the threshold as the proportion of maximum entropy and proceed by optimizing the policy profile in Equation (9). Due to the lack of an analytical solution for Equation (9), we use a policy profile trained under known transitions as a benchmark. We observe that the learned policy profile π_n^* converges to the policy profile under known transitions in $n = 80$ episodes. Figure 2d shows two phases of the learning process. During the first 60 episodes, the performance is volatile, but once the epistemic uncertainty around true transitions is tight, the model exploits it rapidly by episode 80.

In Figure 2f, we empirically show that SAFE-M³-UCRL satisfies safety constraints during the entire execution. In Figure 3, we use a city map of Shenzhen to show that SAFE-M³-UCRL improves service accessibility in low-demand areas. Figure 3b shows that M³-UCRL under known transitions learns how to satisfy the demand ρ_0 effectively at the cost of violating safety constraint (see Figure 2f). Figure 3c shows that SAFE-M³-UCRL under known transitions improves safety by distributing vehicles to residential areas

in the northwest and northeast. Finally, Figure 3d emphasizes the capability of SAFE-M³-UCRL to learn transitions while the policy profile π_n^* simultaneously converges towards the results achieved under known transitions with the number of episodes n passed.

Note that the results in this section are generated assuming the infinite regime. In Appendix D.5 [35], we show that in the finite regime the policy profile π_n^* can be successfully applied to millions of individual agents in real-time, which might be of particular importance to real-world practitioners. The code we use to train and evaluate SAFE-M³-UCRL is available in our GitHub repository [34].

6 CONCLUSION

We present a novel formulation of the mean-field model-based reinforcement learning problem incorporating safety constraints. SAFE-M³-UCRL addresses this problem by leveraging epistemic uncertainty under an unknown transition model and employing a log-barrier approach to ensure conservative satisfaction of the constraints. Beyond the synthetic swarm motion experiment, we showcase the potential of our algorithm for real-world applications by effectively matching the demand distribution in a shared mobility service while consistently upholding service accessibility. In the future, we believe that integrating safety in intelligent multi-agent systems will have a crucial impact on various applications, such as autonomous ride-hailing, firefighting robots and drone/robot search-and-rescue operations in complex and confined spaces.

ACKNOWLEDGMENTS

The authors would like to thank Mojmír Mutný for the fruitful discussions and the Shenzhen Urban Transport Planning Center for collecting and sharing the data used in this work. Matej Jusup acknowledges support from the Swiss National Science Foundation (SNSF) under the research project DADA/181210. Barna Pasztor acknowledges a doctoral fellowship at ETH AI Center. Francesco Corman acknowledges Grant 2023-FS-331 for Research in the area of Public Transport. Andreas Krause acknowledges support by the European Research Council under the European Union’s Horizon 2020 research and innovation program grant agreement no. 815943 and the SNSF under NCCR Automation, grant agreement 51NF40 180545. Ilija Bogunovic acknowledges support from the EPSRC New Investigator Award EP/X03917X/1.

REFERENCES

- [1] Noha Almulla, Rita Ferreira, and Diogo Gomes. 2017. Two numerical approaches to stationary mean-field games. *Dynamic Games and Applications* 7 (2017), 657–682.
- [2] Yoav Alon and Huiyu Zhou. 2020. Multi-agent reinforcement learning for unmanned aerial vehicle coordination by multi-critic policy gradient optimization. *arXiv preprint arXiv:2012.15472* (2020).
- [3] Mohammed Alshiekh, Roderick Bloem, Rüdiger Ehlers, Bettina Könighofer, Scott Niekuim, and Ufuk Topcu. 2018. Safe reinforcement learning via shielding. In *Proceedings of the AAAI Conference on Artificial Intelligence*, Vol. 32.
- [4] Andrea Angiuli, Jean-Pierre Fouque, and Mathieu Laurière. 2021. Reinforcement Learning for Mean Field Games, with Applications to Economics. *arXiv preprint arXiv:2106.13755* (2021).
- [5] Andrea Angiuli, Jean-Pierre Fouque, and Mathieu Laurière. 2022. Unified reinforcement Q-learning for mean field game and control problems. *Mathematics of Control, Signals, and Systems* (2022), 1–55.
- [6] Nicole Bauerle. 2021. Mean Field Markov Decision Processes. *arXiv preprint arXiv:2106.08755* (2021).
- [7] Felix Berkenkamp, Matteo Turchetta, Angela Schoellig, and Andreas Krause. 2017. Safe model-based reinforcement learning with stability guarantees. *Advances in neural information processing systems* 30 (2017).
- [8] Michiel CJ Bliemer and Mark PH Raadsen. 2020. Static traffic assignment with residual queues and spillback. *Transportation Research Part B: Methodological* 132 (2020), 303–319.
- [9] Anton Braverman, Jim G Dai, Xin Liu, and Lei Ying. 2019. Empty-car routing in ridesharing systems. *Operations Research* 67, 5 (2019), 1437–1452.
- [10] René Carmona, Kenza Hamidouche, Mathieu Laurière, and Zongjun Tan. 2021. Linear-quadratic zero-sum mean-field type games: Optimality conditions and policy optimization. *Journal of Dynamics and Games*. 2021, Volume 8, Pages 403–443 8, 4 (2021), 403.
- [11] René Carmona, Mathieu Laurière, and Zongjun Tan. 2019. Linear-quadratic mean-field reinforcement learning: convergence of policy gradient methods. *arXiv preprint arXiv:1910.04295* (2019).
- [12] René Carmona, Mathieu Laurière, and Zongjun Tan. 2019. Model-free mean-field reinforcement learning: mean-field MDP and mean-field Q-learning. *arXiv preprint arXiv:1910.12802* (2019).
- [13] Minshuo Chen, Yan Li, Ethan Wang, Zhuoran Yang, Zhaoran Wang, and Tuo Zhao. 2021. Pessimism meets invariance: Provably efficient offline mean-field multi-agent RL. *Advances in Neural Information Processing Systems* 34 (2021), 17913–17926.
- [14] Richard Cheng, Gábor Orosz, Richard M Murray, and Joel W Burdick. 2019. End-to-end safe reinforcement learning through barrier functions for safety-critical continuous control tasks. In *Proceedings of the AAAI Conference on Artificial Intelligence*, Vol. 33, 3387–3395.
- [15] Sayak Ray Chowdhury and Aditya Gopalan. 2019. Online learning in kernelized markov decision processes. In *The 22nd International Conference on Artificial Intelligence and Statistics*. PMLR, 3197–3205.
- [16] Kurtland Chua, Roberto Calandra, Rowan McAllister, and Sergey Levine. 2018. Deep reinforcement learning in a handful of trials using probabilistic dynamics models. *Advances in neural information processing systems* 31 (2018).
- [17] Sebastian Curi, Felix Berkenkamp, and Andreas Krause. 2020. Efficient model-based reinforcement learning through optimistic policy search and planning. *Advances in Neural Information Processing Systems* 33 (2020), 14156–14170.
- [18] Carlos F Daganzo. 1994. The cell transmission model: A dynamic representation of highway traffic consistent with the hydrodynamic theory. *Transportation research part B: methodological* 28, 4 (1994), 269–287.
- [19] Gal Dalal, Krishnamurthy Dvijotham, Matej Vecerik, Todd Hester, Cosmin Paduraru, and Yuval Tassa. 2018. Safe exploration in continuous action spaces. *arXiv preprint arXiv:1801.08757* (2018).
- [20] Varsha Dani, Thomas P Hayes, and Sham M Kakade. 2008. Stochastic linear optimization under bandit feedback. (2008).
- [21] Romuald Elie, Julien Perolat, Mathieu Laurière, Matthieu Geist, and Olivier Pietquin. 2020. On the convergence of model free learning in mean field games. In *Proceedings of the AAAI Conference on Artificial Intelligence*, Vol. 34, 7143–7150.
- [22] Ingry El Sayed-Aly, Suda Bharadwaj, Christopher Amato, Rüdiger Ehlers, Ufuk Topcu, and Lu Feng. 2021. Safe Multi-Agent Reinforcement Learning via Shielding. In *Proceedings of the 20th International Conference on Autonomous Agents and MultiAgent Systems*. 483–491.
- [23] Javier Garcia and Fernando Fernández. 2015. A comprehensive survey on safe reinforcement learning. *Journal of Machine Learning Research* 16, 1 (2015), 1437–1480.
- [24] Nicolas Gast, Bruno Gaujal, and Jean-Yves Le Boudec. 2012. Mean field for Markov decision processes: from discrete to continuous optimization. *IEEE Trans. Automat. Control* (2012), 2266–2280.
- [25] Clement Gehring and Doina Precup. 2013. Smart exploration in reinforcement learning using absolute temporal difference errors. In *Proceedings of the 2013 international conference on Autonomous agents and multi-agent systems*. 1037–1044.
- [26] Haotian Gu, Xin Guo, Xiaoli Wei, and Renyuan Xu. 2020. Dynamic Programming Principles for Mean-Field Controls with Learning. *arXiv preprint arXiv:1911.07314* (2020).
- [27] Haotian Gu, Xin Guo, Xiaoli Wei, and Renyuan Xu. 2021. Mean-Field Controls with Q-Learning for Cooperative MARL: Convergence and Complexity Analysis. *SIAM Journal on Mathematics of Data Science* 3, 4 (2021), 1168–1196.
- [28] Shangding Gu, Jakub Grudzień Kuba, Munning Wen, Ruiqing Chen, Ziyuan Wang, Zheng Tian, Jun Wang, Alois Knoll, and Yaodong Yang. 2021. Multi-agent constrained policy optimisation. *arXiv preprint arXiv:2110.02793* (2021).
- [29] Shangding Gu, Long Yang, Yali Du, Guang Chen, Florian Walter, Jun Wang, Yaodong Yang, and Alois Knoll. 2022. A review of safe reinforcement learning: Methods, theory and applications. *arXiv preprint arXiv:2205.10330* (2022).
- [30] Minyi Huang, Peter E Caines, and Roland P Malhamé. 2007. Large-population cost-coupled LQG problems with nonuniform agents: individual-mass behavior and decentralized ϵ -equilibria. *IEEE Trans. Automat. Control* 52, 9 (2007), 1560–1571.
- [31] Minyi Huang, Roland P Malhamé, Peter E Caines, et al. 2006. Large population stochastic dynamic games: closed-loop McKean-Vlasov systems and the Nash certainty equivalence principle. *Communications in Information & Systems* 6, 3 (2006), 221–252.
- [32] Ramon Iglesias, Federico Rossi, Kevin Wang, David Hallac, Jure Leskovec, and Marco Pavone. 2018. Data-driven model predictive control of autonomous mobility-on-demand systems. In *2018 IEEE international conference on robotics and automation (ICRA)*. IEEE, 6019–6025.
- [33] Thomas Jaksch, Ronald Ortner, and Peter Auer. 2010. Near-optimal Regret Bounds for Reinforcement Learning. *Journal of Machine Learning Research* (2010).
- [34] Matej Jusup and Tadeusz Janik. 2023. *Safe Model-Based Multi-Agent Mean-Field Reinforcement Learning*. <https://doi.org/10.5281/zenodo.10431636>
- [35] Matej Jusup, Barna Pásztor, Tadeusz Janik, Kenan Zhang, Francesco Corman, Andreas Krause, and Ilija Bogunovic. 2023. Safe Model-Based Multi-Agent Mean-Field Reinforcement Learning. *arXiv preprint arXiv:2306.17052* (2023).
- [36] Daniel Lackner. 2017. Limit theory for controlled McKean–Vlasov dynamics. *SIAM Journal on Control and Optimization* 55, 3 (2017), 1641–1672.
- [37] Balaji Lakshminarayanan, Alexander Pritzel, and Charles Blundell. 2017. Simple and scalable predictive uncertainty estimation using deep ensembles. *Advances in neural information processing systems* 30 (2017).
- [38] Jean-Michel Lasry and Pierre-Louis Lions. 2006. Jeux à champ moyen. i–le cas stationnaire. *Comptes Rendus Mathématique* 343, 9 (2006), 619–625.
- [39] Jean-Michel Lasry and Pierre-Louis Lions. 2006. Jeux à champ moyen. ii–horizon fini et contrôle optimal. *Comptes Rendus Mathématique* 343, 10 (2006), 679–684.
- [40] Mathieu Laurière, Sarah Perrin, Matthieu Geist, and Olivier Pietquin. 2022. Learning Mean Field Games: A Survey. *arXiv preprint arXiv:2205.12944v2* (2022).
- [41] Kaixiang Lin, Renyu Zhao, Zhe Xu, and Jiayu Zhou. 2018. Efficient large-scale fleet management via multi-agent deep reinforcement learning. In *Proceedings of the 24th ACM SIGKDD international conference on knowledge discovery & data mining*. 1774–1783.
- [42] Chenyi Liu, Nan Geng, Vaneet Aggarwal, Tian Lan, Yuan Yang, and Mingwei Xu. 2021. Cmix: Deep multi-agent reinforcement learning with peak and average constraints. In *Machine Learning and Knowledge Discovery in Databases. Research Track: European Conference, ECML PKDD 2021, Bilbao, Spain, September 13–17, 2021, Proceedings, Part I* 21. Springer, 157–173.
- [43] Ryan Lowe, Yi I Wu, Aviratamar, Jean Harb, OpenAI Pieter Abbeel, and Igor Mordatch. 2017. Multi-agent actor-critic for mixed cooperative-competitive environments. *Advances in neural information processing systems* 30 (2017).
- [44] Songtao Lu, Kaiqing Zhang, Tianyi Chen, Tamer Başar, and Lior Horesh. 2021. Decentralized policy gradient descent ascent for safe multi-agent reinforcement learning. In *Proceedings of the AAAI Conference on Artificial Intelligence*, Vol. 35, 8767–8775.
- [45] Chao Mao, Yulin Liu, and Zuo-Jun Max Shen. 2020. Dispatch of autonomous vehicles for taxi services: A deep reinforcement learning approach. *Transportation Research Part C: Emerging Technologies* 115 (2020), 102626.
- [46] Teodor Mihai Moldovan and Pieter Abbeel. 2012. Safe exploration in markov decision processes. *arXiv preprint arXiv:1205.4810* (2012).
- [47] Teodor Mihai Moldovan, Sergey Levine, Michael Jordan, and Pieter Abbeel. 2015. Optimism-driven exploration for nonlinear systems. In *2015 IEEE International Conference on Robotics and Automation (ICRA)*. IEEE, 3239–3246.
- [48] Washim Uddin Mondal, Vaneet Aggarwal, and Satish V Ukkusuri. 2022. Mean-Field Approximation of Cooperative Constrained Multi-Agent Reinforcement Learning (CMARL). *arXiv preprint arXiv:2209.07437* (2022).
- [49] Méderic Motte and Huyê Pham. 2019. Mean-field Markov decision processes with common noise and open-loop controls. *arXiv preprint arXiv:1912.07883* (2019).
- [50] Barna Pásztor, Andreas Krause, and Ilija Bogunovic. 2023. Efficient Model-Based Multi-Agent Mean-Field Reinforcement Learning. *Transactions on Machine Learning Research* (2023).
- [51] Tabish Rashid, Mikayel Samvelyan, Christian Schroeder De Witt, Gregory Farquhar, Jakob Foerster, and Shimon Whiteson. 2020. Monotonic value function factorisation for deep multi-agent reinforcement learning. *The Journal of Machine*

- Learning Research* 21, 1 (2020), 7234–7284.
- [52] Martin Roesch, Christian Linder, Roland Zimmermann, Andreas Rudolf, Andrea Hohmann, and Gunther Reinhart. 2020. Smart Grid for Industry Using Multi-Agent Reinforcement Learning. *Applied Sciences* 10, 19 (2020). <https://doi.org/10.3390/app10196900>
 - [53] Pier Giuseppe Sessa, Maryam Kamgarpour, and Andreas Krause. 2022. Efficient Model-based Multi-agent Reinforcement Learning via Optimistic Equilibrium Computation. , 19580–19597 pages.
 - [54] Susan Shaheen, Corwin Bell, Adam Cohen, Balaji Yelchuru, Booz Allen Hamilton, et al. 2017. *Travel behavior: Shared mobility and transportation equity*. Technical Report. United States. Federal Highway Administration. Office of Policy
 - [55] Ziyad Sheebaelhamd, Konstantinos Zisis, Athina Nisioti, Dimitris Gkouletsos, Dario Pavllo, and Jonas Kohler. 2021. Safe Deep Reinforcement Learning for Multi-Agent Systems with Continuous Action Spaces. *arXiv preprint arXiv:2108.03952* (2021).
 - [56] Yong Song, Yi-bin Li, Cai-hong Li, and Gui-fang Zhang. 2012. An efficient initialization approach of Q-learning for mobile robots. *International Journal of Control, Automation and Systems* 10, 1 (2012), 166–172.
 - [57] Nirajan Srinivas, Andreas Krause, Sham Kakade, and Matthias Seeger. 2010. Gaussian Process Optimization in the Bandit Setting: No Regret and Experimental Design. In *International Conference on Machine Learning*, 1015–1022.
 - [58] Jayakumar Subramanian and Aditya Mahajan. 2019. Reinforcement Learning in Stationary Mean-Field Games. In *International Conference on Autonomous Agents and MultiAgent Systems*. 251–259.
 - [59] Ilnura Usmanova, Yarden As, Maryam Kamgarpour, and Andreas Krause. 2022. Log barriers for safe black-box optimization with application to safe reinforcement learning. *arXiv preprint arXiv:2207.10415* (2022).
 - [60] Jeroen PT van der Gun, Adam J Pel, and Bart Van Arem. 2018. The link transmission model with variable fundamental diagrams and initial conditions. *Transportmetrica B: Transport Dynamics* (2018).
 - [61] Lingxiao Wang, Zhuoran Yang, and Zhaoran Wang. 2020. Breaking the Curse of Many Agents: Provably Mean Embedding Q-Iteration for Mean-Field Reinforcement Learning. In *International Conference on Machine Learning*. 10092–10103.
 - [62] Weichen Wang, Jiequn Han, Zhuoran Yang, and Zhaoran Wang. 2021. Global Convergence of Policy Gradient for Linear-Quadratic Mean-Field Control/Game in Continuous Time. , 10772–10782 pages.
 - [63] Jian Wen, Jinhua Zhao, and Patrick Jaillet. 2017. Rebalancing shared mobility-on-demand systems: A reinforcement learning approach. In *2017 IEEE 20th international conference on intelligent transportation systems (ITSC)*. Ieee, 220–225.
 - [64] Margaret H Wright. 1992. Interior methods for constrained optimization. *Acta numerica* 1 (1992), 341–407.
 - [65] Rick Zhang and Marco Pavone. 2016. Control of robotic mobility-on-demand systems: a queueing-theoretical perspective. *The International Journal of Robotics Research* 35, 1-3 (2016), 186–203.
 - [66] Zheqing Zhu, Erdem Biyik, and Dorsa Sadigh. 2020. Multi-agent safe planning with gaussian processes. In *2020 IEEE/RSJ International Conference on Intelligent Robots and Systems (IROS)*. IEEE, 6260–6267.