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

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

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

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

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

Francesco Corman
ETH Zurich
Zurich, Switzerland
corman@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-M$^3$-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-M$^3$-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

**Figure 1: An illustration of vehicles' spatial distribution (light-blue scatters), 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 [47]), swarm robotics (e.g., operating a swarm of drones [3]), and smart grids (e.g., operating a network of sensors in electric system [62]). 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 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 [14, 41, 55, 71] 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 [9, 20, 70].

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-$M^3$-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 algorithm called Safe-$M^3$-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-$M^3$-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) [35, 36, 43, 44]. For a summary of the progress focusing on MFGs, see [45] and references therein. We focus on MFCs in this work,

which assume cooperative agents in contrast to MFGs, which assume competition. [7, 27, 30, 31, 56] address the problem of solving MFCs under known transitions, i.e., planning, while [5, 6, 11–13, 68, 71, 72] consider model-free Q-learning and Policy Gradient methods in various settings. Closest to our approach, [58] introduces $M^3$-UCRL, a model-based, on-policy algorithm, which is more sample efficient than other proposed approaches. Similarly to [16, 19, 39] for model-based single-agent RL and [63] for model-based MARL, $M^3$-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.

In terms of *safety*, there are two main ways of handling it in RL; assigning significantly lower rewards to unsafe states [53] and providing additional knowledge to the agents [66] or using the notion of controllability to avoid unsafe policies explicitly [28]. Furthermore, the following approaches combine the two methods; [8] uses Lyapunov functions to restrict the safe policy space, [15] projects unsafe policies to a safe set via a control barrier function, and [4] introduces shielding, i.e., correcting actions only if they lead to unsafe states. For comprehensive overviews on safe RL, we refer the reader to [26, 33]. As an alternative, [69] 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, [32] 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 [50]. For the non-cooperative decentralized MARL problem with individual constraints, [65] adds a safety layer to multi-agent DDPG [49] similar to single-agent Safe DDPG [21] for continuous state-action spaces. Aggregated and population-based constraints have been addressed in the following works. CMIX [48] extends QMIX [61], 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. [25] 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, [76] 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, [55] 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 $n = 1, \ldots, N$ each have $t = 0, \ldots, T-1$ 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 $i \in \{1, \ldots, m\}$ 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 $m \to +\infty$, i.e.,

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

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 $r : \mathcal{S} \times \mathcal{P}(\mathcal{S}) \times \mathcal{A} \to \mathbb{R}$ 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 $\boldsymbol{\pi}_n = (\pi_{n,0}, \ldots, \pi_{n,T-1}) \in \Pi$ where individual policies are of the form $\pi_{n,t} : \mathcal{S} \times \mathcal{P}(\mathcal{S}) \to \mathcal{A}$ and $\Pi$ is the set of admissible policy profiles. The policy profile is then shared with all the agents that choose their actions according to $\boldsymbol{\pi}_n$ during episode $n$.

We consider a general family of deterministic transitions $f : \mathcal{S} \times \mathcal{P}(\mathcal{S}) \times \mathcal{A} \to \mathcal{S}$. 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

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

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 $\boldsymbol{\pi}_n$ and transitions $f$ in episode $n$ (see [58, Lemma 1])

$$\mu_{n,t+1}(ds') = \int_{\mathcal{S}} \mathbb{P}[s_{n,t+1} \in ds']\mu_{n,t}(ds), \tag{2}$$

where $s_{n,t+1}$ is the next representative agent state and $\mu_{n,t}(ds) = \mathbb{P}[s_{n,t} \in ds]$ under $\boldsymbol{\pi}_n$ for all $t$. To simplify the notation, we use $U(\cdot)$ 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 $\boldsymbol{\pi}_n$ and mean-field distribution $\mu$, 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 $h : \mathcal{P}(\mathcal{S}) \to \mathbb{R}$ over probability distributions. Given some hard safety threshold $C \in \mathbb{R}$, we consider a mean-field distribution $\mu$ as safe if it satisfies $h(\mu) \geq C$, or, equivalently

$$h_C(\mu) := h(\mu) - C \geq 0. \tag{3}$$

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 $\mu_0$, we formally define the *Safe-MFC*[3] problem as follows

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

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

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

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

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

where we explicitly require induced mean-field distributions $\{\mu_t\}_{t=1}^{T}$ to reside in the safe set $\zeta$ by restricting the set of admissible policy profiles $\Pi$ to policy profiles that induce safe distributions. To ensure complete safety, we note that the initial mean-field distribution $\mu_0$ *must be in the safe set $\zeta$ 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 $\mu$ and $\mu'$ (i.e., a joint probability distributions with marginals $\mu$ and $\mu'$). We further define the distance between $z = (s, \mu, a)$ and $z' = (s', \mu', a')$ as

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

---

[1]Our framework easily extends to unknown reward by estimating its epistemic uncertainty and learning it similarly to learning unknown transitions (see [16]).

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

[3]We 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 $\pi$ present in any admissible policy profile $\boldsymbol{\pi}$ in $\Pi$ are $L_\pi$-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 [16, 39, 58, 63] 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 (see Appendix D.3).

## 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

$$H(\mu) := - \int_S \log \mu(s) \mu(ds) \tag{5}$$

of the mean-field distribution $\mu$, i.e.,

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

**Distribution similarity.** Another way to define safety is by preventing $\mu$ from diverging from a prior distribution $\nu_0$. 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

$$h_C(\mu; \nu_0) := C - D(\mu, \nu_0),$$

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

We provide further examples of safety functions in Appendix B 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 $i = 1, ..., n - 1$, i.e., up until the beginning of episode $n$. We estimate the mean $\boldsymbol{m}_{n-1} : \mathcal{Z} \to \mathcal{S}$ and covariance $\Sigma_{n-1} : \mathcal{Z} \to \mathbb{R}^{p \times p}$ functions from the set of collected trajectories $\cup_{i=1}^{n-1} \mathcal{D}_i$, and denote model's confidence with $\boldsymbol{\sigma}_{n-1}^2(z) = \text{diag}(\Sigma_{n-1}(z))$. 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 [16, 19, 58, 63, 67] and other literature which aims to exclude extreme functionals from consideration.

ASSUMPTION 4 (CALIBRATED MODEL). *Let $\boldsymbol{m}_{n-1}(\cdot)$ 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 $\boldsymbol{\sigma}_{n-1}(\cdot)$, 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 $|f(z) - \boldsymbol{m}_{n-1}(z)| \leq \beta_{n-1} \boldsymbol{\sigma}_{n-1}(z)$ elementwise.*

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

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 $\boldsymbol{m}_{n-1}(\cdot)$ and $\boldsymbol{\sigma}_{n-1}(\cdot)$ estimated based on the observations up until the end of the previous episode $n - 1$, i.e.,

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

The crucial challenge in *Safe Mean-Field Reinforcement Learning* 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 $t$*

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

*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

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

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 $C_{n,t} \ge W_1(\tilde{\mu}_{n,t}, \mu_{n,t})$. □

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 $t$, $h_C(\tilde{\mu}_{n,t}) \ge L_h C_{n,t}$ implies $h_C(\mu_{n,t}) \ge 0$ 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

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

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

Then, $C_{n,t}$ for $t = 1, \ldots, T$ 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 $C_{n,t} \ge W_1(\tilde{\mu}_{n,t}, \mu_{n,t})$ 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, 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 $h_C(\tilde{\mu}_{n,t}) \ge L_h C_{n,t}$ 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 Mean-Field Reinforcement Learning* 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 $\pi_n^*$ 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 $\pi^*$ from Equation (4) is approximated at the episode $n$ by

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

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

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

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

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

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 [22]. 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 $\pi_n^*$ at the beginning of episode $n$.

We use an established approach known as *Hallucinated Upper Confidence Reinforcement Learning* (H-UCRL) [19, 54, 58] and introduce an auxiliary function $\eta : \mathcal{Z} \to [-1, 1]^p$, where $p$ is the dimensionality of the state space $\mathcal{S}$, to define hallucinated transitions

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

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)

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

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

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 $\Pi$ and auxiliary function $\eta(\cdot)$ (see Appendix C.2 for further details). Note that $\eta(z) = \eta(s, \mu, \pi(s, \mu)) = \eta(s, \mu)$ for a fixed individual policy $\pi$. This turns $\eta(\cdot)$ into a policy that exerts *hallucinated control* over the epistemic uncertainty of the confidence set of transitions $\mathcal{F}_{n-1}$ [19]. Furthermore, Equation (8) allows us to optimize over parametrizable functions (e.g., *neural networks*) $\pi$ and $\eta(\cdot)$ using gradient ascent.

We introduce the safety constraint to the objective using the *log-barrier method* [74]. 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

$$\pi_n^* = \underset{\pi_n \in \Pi}{\arg\max} \; \underset{\eta(\cdot) \in [-1,1]^p}{\max}$$

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

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

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

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

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

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 $\mu_0$ is given) is not empty, $\pi_n^*$ 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 $\boldsymbol{m}_{n-1}(\cdot) = f(\cdot)$, $\Sigma_{n-1}(\cdot) = 0$ and $L_h C_{n,t} = 0$, hence, recovering the original constraint $h_C(\tilde{\mu}_{n,t}) \ge 0$ 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 $\Pi$, safety constraint $h_C(\cdot)$, calibrated statistical model represented by $\boldsymbol{m}_{n-1}(\cdot)$ and $\Sigma_{n-1}(\cdot)$, initial mean-field distribution $\mu_0$, known reward $r(\cdot)$, safety Lipschitz constant $L_h$, hyperparameter $\lambda$, number of episodes $N$, number of steps $T$

1: **for** $n = 1, \ldots N$ **do**
2:     Compute $C_{n,t}$ for $t = 1, \ldots, T$ as described in Appendix A
3:     Optimize the objective in Equation (9) over the admissible policy profiles $\Pi$ and hallucinated transitions Equation (8)
4:     Execute the obtained policy profile $\boldsymbol{\pi}_n^*$ 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** $\boldsymbol{\pi}_N^* = (\pi_{N,0}^*, \ldots, \pi_{N,T-1}^*)$

---

We summarize the model-based learning protocol used by SAFE-M³-UCRL in Algorithm 1. The first step computes constants $C_{n,t}$ (see Appendix A) 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, 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 Mean-Field RL as a research topic, one of the rare benchmarks used by multiple authors is the swarm motion. [13, 58] view it as Mean-Field RL problem, while [24] 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 [2].

**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 $\varepsilon_{n,t} \sim N(0, \Delta t)$ for all episodes $n$ and steps $t$. 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 $s$ (see Appendix E.3), 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, $\Delta t \to 0$ can be obtained analytically [2] (see Appendix E.1) and used as a benchmark. [13, 58] show that Mean-Field RL 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 $\boldsymbol{\pi}_n = (\pi_{n,0}, \ldots, \pi_{n,T-1})$, for every episode $n$, during the optimization of Equation (9). The optimization is done by *Mean-Field Back-Propagation Through Time* (MF-BPTT) (see Appendix C.3.1). 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 $\pi_{n,t}$, $t = 0, \ldots, T - 1$, is a natural extension. We use a *Probabilistic Neural Network Ensemble* [17, 42] to represent a statistical model of transitions, which we elaborate in Appendix C.1. 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 $\zeta$

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

Note that, in general, the safe initial distribution might not exist. We provide the implementation details in Appendix D.2.

**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 [58] that relies on the overcrowding penalty, while for $p = 0.95$ we significantly surpass the effect that the penalty has on

(a) Swarm motion learning curve     (b) Swarm motion mean-field distributions     (c) Swarm motion safety for $p$ = 0.95

(d) Vehicle repositioning learning curve     (e) Swarm motion policies     (f) Vehicle repositioning safety for $p$ = 0.85
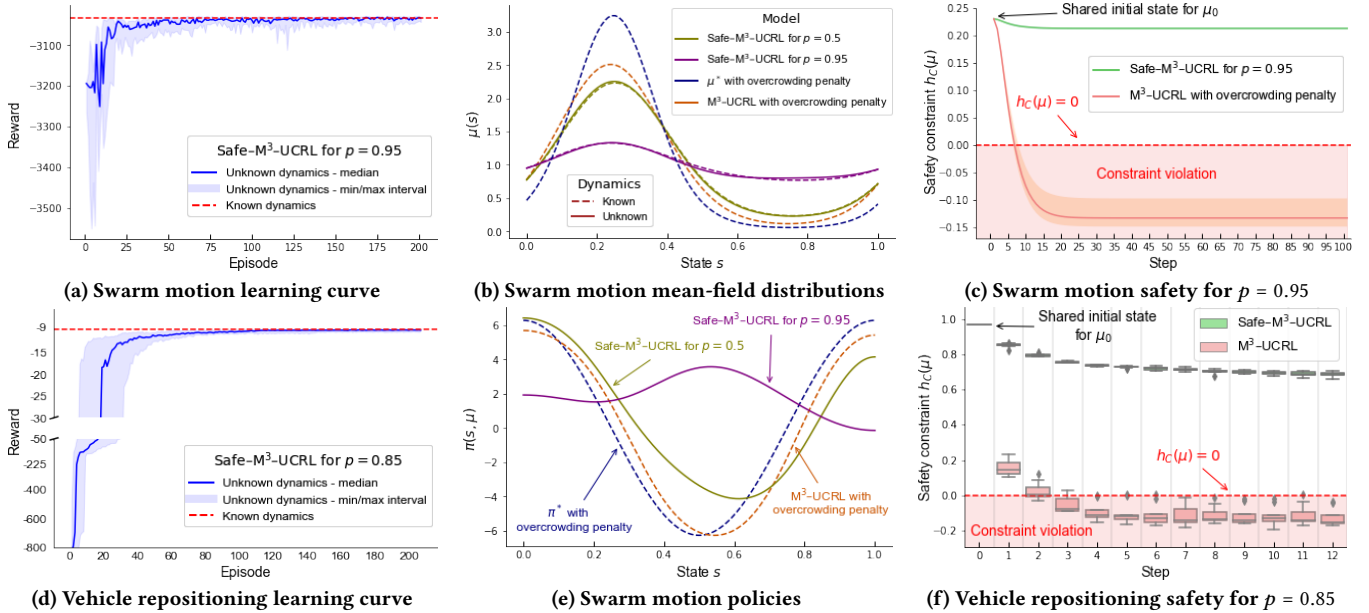
**Figure 2: Performance analysis of SAFE-M³-UCRL for swarm motion and vehicle repositioning. The policy and statistical model were trained on 10 randomly initialized neural networks for each hyperparameter $p$.**

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 $\mu^*$. Furthermore, we observe that the policies under unknown transitions overlap with the solutions learned under known transitions. In Figure 2e, we observe 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. For further details, see Appendix E.3.

## 5.2 Vehicle Repositioning Problem

Since ride-hailing services, such as Uber, Lyft, and Bolt, gained popularity and market share, vehicle repositioning has been a long-standing challenge for these platforms, i.e., moving idle vehicles to areas with high-demand potential. A similar challenge is present in bike-sharing services accessible in many cities and, as of more recently, dockless electric scooter-sharing services such as Bird and Lime. In the competitive environment, the operator significantly increases profit 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 [10, 75], model predictive control [37], to RL [47, 52, 73]. 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 [64].

**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 (for more details, see Appendix D). 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
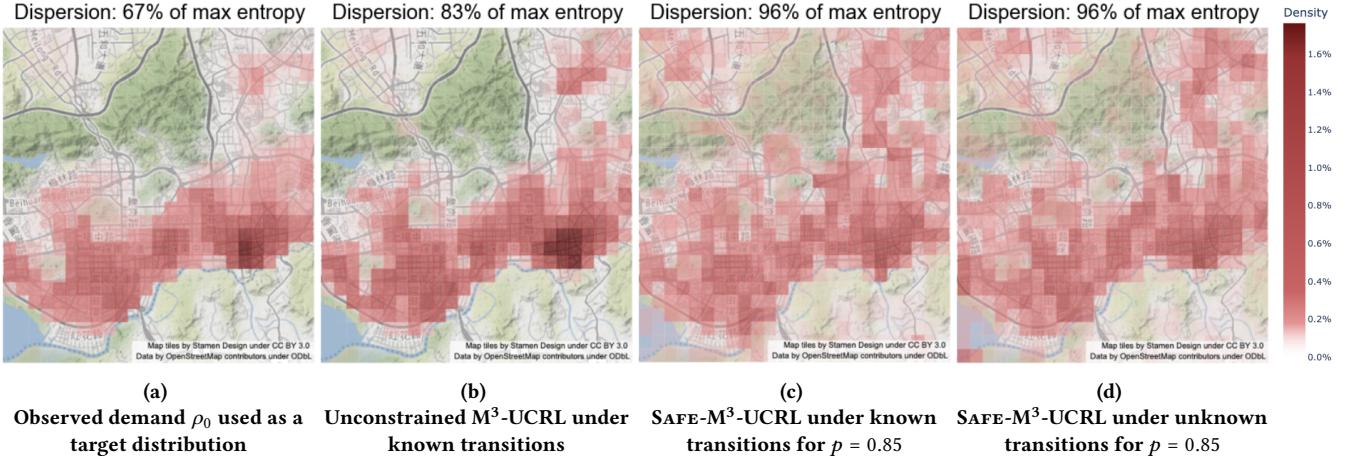
| | | | Density |
|---|---|---|---|
| Dispersion: 67% of max entropy | Dispersion: 83% of max entropy | Dispersion: 96% of max entropy | Dispersion: 96% of max entropy |

**(a)** Observed demand $\rho_0$ used as a target distribution

**(b)** Unconstrained $M^3$-UCRL under known transitions

**(c)** Safe-$M^3$-UCRL under known transitions for $p = 0.85$

**(d)** Safe-$M^3$-UCRL under unknown transitions for $p = 0.85$

Figure 3: Safe-$M^3$-UCRL guided vehicle distribution in Shenzhen in the evening peak hours.

service denoted as $\rho_0$, 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 $\rho_0$ 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 (see Appendix D.1). We use $\Phi(\cdot)$ to define sequential transitions by 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 $\sigma^2$ truncated at the borders of $\mathcal{S}$ and $I_2$ is the $2 \times 2$ 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). Therefore, the optimization objective in Equation (9) trades off between greedily satisfying the demand $\rho_0$ 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 $\pi_n$, 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 as elaborated in Appendix D.1.

**Results.** The entropy of the target distribution, $\rho_0$ in Figure 3a, already achieves $p = 0.67$ of the maximum due to a wide horizontal spread. To achieve 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 $\pi_n^*$ 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-$M^3$-UCRL satisfies safety constraints during the entire execution. In Figure 3, we use a city map of Shenzhen to show that Safe-$M^3$-UCRL improves service accessibility in low-demand areas. Figure 3b shows that $M^3$-UCRL under known transitions learns how to satisfy the demand $\rho_0$ effectively at the cost of violating safety constraint (see Figure 2f). Figure 3c shows that Safe-$M^3$-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-$M^3$-UCRL to learn complex transitions while the policy profile $\pi_n^*$ simultaneously converges towards the results achieved under known transitions with the number of episodes $n$ passed.

In Appendix D.3, we provide the details on the parameters used during the training and exhaustive performance analysis. The results in this section are generated assuming the infinite regime. At the same time, in Appendix D.5, we showcase that in the finite regime the policy profile $\pi_n^*$ 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-$M^3$-UCRL is available in our GitHub repository [40].

## 6 Conclusion

We present a novel formulation of the mean-field model-based reinforcement learning problem incorporating safety constraints. Safe-$M^3$-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 considerations 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.

## REFERENCES

[1] Yasin Abbasi-Yadkori and Csaba Szepesvári. 2011. Regret bounds for the adaptive control of Linear Quadratic systems. *Journal of Machine Learning Research* 19 (2011), 1–26.

[2] Noha Almulla, Rita Ferreira, and Diogo Gomes. 2017. Two numerical approaches to stationary mean-field games. *Dynamic Games and Applications* 7 (2017), 657–682.

[3] 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).

[4] Mohammed Alshiekh, Roderick Bloem, Rüdiger Ehlers, Bettina Könighofer, Scott Niekum, and Ufuk Topcu. 2018. Safe reinforcement learning via shielding. In *Proceedings of the AAAI Conference on Artificial Intelligence*, Vol. 32.

[5] 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).

[6] 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.

[7] Nicole Bäuerle. 2021. Mean Field Markov Decision Processes. *arXiv preprint arXiv:2106.08755* (2021).

[8] 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).

[9] 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.

[10] Anton Braverman, Jim G Dai, Xin Liu, and Lei Ying. 2019. Empty-car routing in ridesharing systems. *Operations Research* 67, 5 (2019), 1437–1452.

[11] 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.

[12] 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).

[13] 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).

[14] 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.

[15] 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.

[16] 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.

[17] 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).

[18] Philippe Clement and Wolfgang Desch. 2008. An elementary proof of the triangle inequality for the Wasserstein metric. *Proc. Amer. Math. Soc.* 136, 1 (2008), 333–339.

[19] 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.

[20] 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.

[21] 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).

[22] Varsha Dani, Thomas P Hayes, and Sham M Kakade. 2008. Stochastic linear optimization under bandit feedback. (2008).

[23] Richard M Dudley. 2018. *Real analysis and probability*. CRC Press.

[24] 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.

[25] Ingy ElSayed-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.

[26] Javier Garcıa and Fernando Fernández. 2015. A comprehensive survey on safe reinforcement learning. *Journal of Machine Learning Research* 16, 1 (2015), 1437–1480.

[27] 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.

[28] 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.

[29] Xavier Glorot and Yoshua Bengio. 2010. Understanding the difficulty of training deep feedforward neural networks. In *Proceedings of the thirteenth international conference on artificial intelligence and statistics*. JMLR Workshop and Conference Proceedings, 249–256.

[30] 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).

[31] 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.

[32] Shangding Gu, Jakub Grudzien Kuba, Munning Wen, Ruiqing Chen, Ziyan Wang, Zheng Tian, Jun Wang, Alois Knoll, and Yaodong Yang. 2021. Multi-agent constrained policy optimisation. *arXiv preprint arXiv:2110.02793* (2021).

[33] 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).

[34] Charles R Harris, K Jarrod Millman, Stéfan J Van Der Walt, Ralf Gommers, Pauli Virtanen, David Cournapeau, Eric Wieser, Julian Taylor, Sebastian Berg, Nathaniel J Smith, et al. 2020. Array programming with NumPy. *Nature* 585, 7825 (2020), 357–362.

[35] Minyi Huang, Peter E Caines, and Roland P Malhamé. 2007. Large-population cost-coupled LQG problems with nonuniform agents: individual-mass behavior and decentralized $\varepsilon$-equilibria. *IEEE Trans. Automat. Control* 52, 9 (2007), 1560–1571.

[36] 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.

[37] 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.

[38] Sergey Ioffe and Christian Szegedy. 2015. Batch normalization: Accelerating deep network training by reducing internal covariate shift. In *International conference on machine learning*. pmlr, 448–456.

[39] Thomas Jaksch, Ronald Ortner, and Peter Auer. 2010. Near-optimal Regret Bounds for Reinforcement Learning. *Journal of Machine Learning Research* (2010).

[40] Matej Jusup and Tadeusz Janik. 2023. *Safe Model-Based Multi-Agent Mean-Field Reinforcement Learning*. https://doi.org/10.5281/zenodo.10431636

[41] Daniel Lacker. 2017. Limit theory for controlled McKean–Vlasov dynamics. *SIAM Journal on Control and Optimization* 55, 3 (2017), 1641–1672.

[42] 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).

[43] Jean-Michel Lasry and Pierre-Louis Lions. 2006. Jeux à champ moyen. i–le cas stationnaire. *Comptes Rendus Mathématique* 343, 9 (2006), 619–625.

[44] 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.

[45] Mathieu Laurière, Sarah Perrin, Matthieu Geist, and Olivier Pietquin. 2022. Learning Mean Field Games: A Survey. *arXiv preprint arXiv:2205.12944v2* (2022).

[46] Timothy P Lillicrap, Jonathan J Hunt, Alexander Pritzel, Nicolas Heess, Tom Erez, Yuval Tassa, David Silver, and Daan Wierstra. 2015. Continuous control with deep reinforcement learning. *arXiv preprint arXiv:1509.02971* (2015).

[47] 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.

[48] 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.

[49] Ryan Lowe, Yi I Wu, Aviv Tamar, 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).

[50] 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.

[51] Ali Malik, Volodymyr Kuleshov, Jiaming Song, Danny Nemer, Harlan Seymour, and Stefano Ermon. 2019. Calibrated model-based deep reinforcement learning. In *International Conference on Machine Learning*. PMLR, 4314–4323.

[52] 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.

[53] Teodor Mihai Moldovan and Pieter Abbeel. 2012. Safe exploration in markov decision processes. *arXiv preprint arXiv:1205.4810* (2012).

[54] Teodor Mihai Moldovan, Sergey Levine, Michael I Jordan, and Pieter Abbeel. 2015. Optimism-driven exploration for nonlinear systems. In *2015 IEEE International Conference on Robotics and Automation (ICRA)*. IEEE, 3239–3246.

[55] 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).

[56] Médéric Motte and Huyên Pham. 2019. Mean-field Markov decision processes with common noise and open-loop controls. *arXiv preprint arXiv:1912.07883* (2019).

[57] Adam Paszke, Sam Gross, Soumith Chintala, Gregory Chanan, Edward Yang, Zachary DeVito, Zeming Lin, Alban Desmaison, Luca Antiga, and Adam Lerer. 2017. Automatic differentiation in pytorch. (2017).

[58] Barna Pásztor, Andreas Krause, and Ilija Bogunovic. 2023. Efficient Model-Based Multi-Agent Mean-Field Reinforcement Learning. *Transactions on Machine Learning Research* (2023).

[59] Matthias Plappert, Rein Houthooft, Prafulla Dhariwal, Szymon Sidor, Richard Y Chen, Xi Chen, Tamim Asfour, Pieter Abbeel, and Marcin Andrychowicz. 2017. Parameter space noise for exploration. *arXiv preprint arXiv:1706.01905* (2017).

[60] Yury Polyanskiy and Yihong Wu. 2016. Wasserstein continuity of entropy and outer bounds for interference channels. *IEEE Transactions on Information Theory* 62, 7 (2016), 3992–4002.

[61] 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.

[62] 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

[63] Pier Giuseppe Sessa, Maryam Kamgarpour, and Andreas Krause. 2022. Efficient Model-based Multi-agent Reinforcement Learning via Optimistic Equilibrium Computation. , 19580–19597 pages.

[64] 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 ….

[65] 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).

[66] 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.

[67] Niranjan 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.

[68] Jayakumar Subramanian and Aditya Mahajan. 2019. Reinforcement Learning in Stationary Mean-Field Games. In *International Conference on Autonomous Agents and MultiAgent Systems*. 251–259.

[69] 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).

[70] 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).

[71] Lingxiao Wang, Zhuoran Yang, and Zhaoran Wang. 2020. Breaking the Curse of Many Agents: Provable Mean Embedding Q-Iteration for Mean-Field Reinforcement Learning. In *International Conference on Machine Learning*. 10092–10103.

[72] 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.

[73] 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.

[74] Margaret H Wright. 1992. Interior methods for constrained optimization. *Acta numerica* 1 (1992), 341–407.

[75] 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.

[76] Zheqing Zhu, Erdem Bıyık, 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.

# Appendix

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

### A  Relationship Between Mean-Field Distributions under True and Estimated Transitions

In this section, we describe the procedure for computing constant $C_{n,t}$ (defined in Section 4) at every episode $n = 1, \ldots, N$ and step $t = 1, \ldots, T$. This is necessary for establishing the connection between the mean-field distribution $\mu_{n,t}$ in the original system under known transitions $f$ (see Equation (4)) and the mean-field distribution $\tilde{\mu}_{n,t}$ in the system induced by a calibrated statistical model $\tilde{f}_{n-1}$ (see Equation (7), Section 3.3 and Assumption 4). In particular, Corollary 1 shows that the safety in the original system $h_C(\mu_{n,t})$ is guaranteed with high probability if $h_C(\tilde{\mu}_{n,t}) \geq L_h C_{n,t}$ for a safety constant $C_{n,t} \geq W_1(\tilde{\mu}_{n,t}, \mu_{n,t})$ and Lipschitz constant $L_h$ introduced in Assumption 6. We use the following result from [58, Lemma 5] to make a connection between the mean-field distributions under true and estimated transitions.

LEMMA 2. *Under Assumptions 1 to 5 and assuming that the event in Assumption 4 holds true, for episodes $n = 1, \ldots, N$, steps $t = 1, \ldots, T$ and fixed policy profile $\boldsymbol{\pi}_n = (\pi_{n,0}, \ldots, \pi_{n,T-1})$, we have:*

$$W_1(\tilde{\mu}_{n,t}, \mu_{n,t}) \leq 2\beta_{n-1}\overline{L}_{n-1}^{t-1} \sum_{i=0}^{t-1} I_{n,t},$$

*where* $I_{n,t} = \int_{\mathcal{S}} \|\sigma_{n-1}(s, \mu_{n,i}, \pi_{n,t}(s, \mu_{n,i}))\|_2 \mu_{n,i}(ds)$ *and* $\overline{L}_{n-1} = 1 + 2(1 + L_\pi)(L_f + 2\beta_{n-1}L_\sigma)$.

Let $K_{n,t} := 2\beta_{n-1}\overline{L}_{n-1}^{t-1}$ and $z = (s, \mu, \pi_{n,t}(s, \mu))$. Then we have:

$$W_1(\tilde{\mu}_{n,t}, \mu_{n,t}) \leq K_{n,t} \sum_{i=0}^{t-1} \int_{\mathcal{S}} \|\sigma_{n-1}(s, \mu_{n,i}, \pi_{n,t}(s, \mu_{n,i}))\|_2 \mu_{n,i}(ds)$$
$$\leq t K_{n,t} \max_{z \in \mathcal{Z}} \|\sigma_{n-1}(z)\|_2,$$

for $t = 1, \ldots, T$, where the first inequality is due to Lemma 2 while the second one follows since maximum upper bounds expectation. Hence, we can set $C_{n,t}$ as $t K_{n,t} \max_{z \in \mathcal{Z}} \|\sigma_{n-1}(z)\|_2$, and calculate $\max_{z \in \mathcal{Z}} \|\sigma_{n-1}(z)\|_2$ once at the beginning of each episode $n$. Notice that our learning protocol (see Algorithm 1) does not require computing $C_{n,0}$ at the initial step $t = 0$ since the mean-field distributions share the initial state, i.e., $\tilde{\mu}_{n,0} = \mu_{n,0}$. It is worth noting that in certain models, such as Gaussian Processes, this upper bound provides a meaningful interpretation. Specifically, $\sigma_{n-1}(z)$ represents the epistemic uncertainty of the model, which tends to decrease monotonically as more data is observed.

### B  Examples of Safety Constraints

In this section, we show some important classes of safety constraints $h_C(\cdot) \geq 0$ satisfying Assumption 6.

### B.1  Entropic Safety

Entropic safety serves the purpose of controlling the dispersion of the mean-field distribution, which is useful in many applications, such as vehicle repositioning (see Section 5). A natural way of defining entropic safety is via differential entropy, but in general, the differential entropy is not Lipschitz continuous due to the unboundedness of the natural logarithm. Nevertheless, the issue can be easily circumvented by considering $\varepsilon$-smoothed differential entropy $H^\varepsilon : \mathcal{P}(\mathcal{S}) \to \mathbb{R}_{\geq 0}$. For $\varepsilon > 0$ and $C \geq 0$ we define $\varepsilon$-smoothed differential entropy and associated entropic safety constraint $H_C^\varepsilon(\cdot)$ as

$$H^\varepsilon(\mu) := -\int_{\mathcal{S}} \log(\mu(s) + \varepsilon)\mu(ds)$$
$$H_C^\varepsilon(\mu) := H^\varepsilon(\mu) - C$$

To show that $H_C^\varepsilon(\cdot)$ satisfies Assumption 6 let $h_C(\cdot) := H_C^\varepsilon(\cdot)$ and assume that $\mathcal{S} \subset \mathbb{R}^p$ is a compact set. First, note that $f(x) = \log(x + \varepsilon)$ is $\frac{1}{\varepsilon}$-Lipschitz continuous for $\varepsilon > 0$, i.e., $\varepsilon f(x)$ is 1-Lipschitz. Second, for every $S \subseteq \mathcal{S}$ and $L$-Lipschitz function $f : \mathcal{S} \to \mathbb{R}$, a function $g(S) = \int_S f(s)\mu(ds)$ is $L$-Lipschitz due to the boundedness of $f$. The following derivation shows that $H_C^\varepsilon(\cdot)$ is $\frac{1}{\varepsilon}$-Lipschitz continuous.

$$|h_C(\mu) - h_C(\mu')| = \left|H_C^\varepsilon(\mu) - H_C^\varepsilon(\mu')\right|$$

$$= \left|\int_{\mathcal{S}} \log(\mu'(s) + \varepsilon)\mu'(ds) - \int_{\mathcal{S}} \log(\mu(s) + \varepsilon)\mu(ds)\right|$$

$$= \left|\int_{\mathcal{S}} \log(\mu'(s) + \varepsilon)(\mu' - \mu)(ds) - \int_{\mathcal{S}} \log(\mu(s) + \varepsilon)\mu(ds)\right.$$
$$\left. + \int_{\mathcal{S}} \log(\mu'(s) + \varepsilon)\mu(ds)\right|$$

$$\leq \left|\int_{\mathcal{S}} \log(\mu'(s) + \varepsilon)(\mu' - \mu)(ds)\right|$$
$$+ \left|\int_{\mathcal{S}} \log(\mu(s) + \varepsilon)\mu(ds) - \int_{\mathcal{S}} \log(\mu'(s) + \varepsilon)\mu(ds)\right|$$

$$\leq \left|\int_{\mathcal{S}} \log(\mu'(s) + \varepsilon)(\mu' - \mu)(ds)\right| + \frac{1}{\varepsilon}W_1(\mu, \mu')$$

$$\leq M\left|\int_{\mathcal{S}} (\mu' - \mu)(ds)\right| + \frac{1}{\varepsilon}W_1(\mu, \mu')$$

$$= \underbrace{M(\mu' - \mu)(\mathcal{S})}_{\overline{M}} + \frac{1}{\varepsilon}W_1(\mu, \mu')$$

$$= \overline{M} + \frac{1}{\varepsilon}W_1(\mu, \mu'),$$

where the first inequality comes from the triangle inequality, the second inequality comes from the Lipschitz continuity of the $\varepsilon$-smoothed logarithm and integral, the third comes from the upper-boundedness of logarithm on a compact set, and the last equality comes from the fact that the measure of a compact set is finite.

REMARK 2. *In Section 5, Appendix D and Appendix E, we use the discrete equivalent of entropic safety, i.e., Shannon entropy, for our experiments because of the discrete representation of the mean-field distribution. [60, Proposition 8] shows that Shannon entropy is Lipschitz continuous with respect to the scaled Wasserstein 1-distance, i.e., with respect to $\frac{1}{n} W_1(\cdot)$ known as Ornstein's distance.*

## B.2 Distribution Similarity

We can define safety by preventing the mean-field distribution $\mu$ from diverging from a prior distribution $\nu_0$ by quantifying the allowed dissimilarity between the two distributions

$$h_C(\mu; \nu_0) := C - D(\mu, \nu_0),$$

where $C \geq 0$ and $D : \mathcal{P}(\mathcal{S}) \times \mathcal{P}(\mathcal{S}) \to \mathbb{R}_{\geq 0}$ is the distance function between probability measures. Commonly used distances are Wasserstein p-distance for $p \geq 1$ and $f$-divergences such as KL-divergence, Hellinger distance, and total variation distance.

A concrete example of a distance function $D(\cdot)$ that satisfies Assumption 6 is Wasserstein 1-distance

$$W_1^C(\mu, \nu_0) := C - W_1(\mu, \nu_0)$$

$$h_C(\mu) := W_1^C(\mu, \nu_0).$$

[18] shows the triangle inequality of Wasserstein p-distance for probability measures on separable metric spaces. [23] shows that Wasserstein 1-distance induces a metric space $(\mathcal{P}(\mathcal{S}), W_1)$ over probability measures. The result now trivially follows from the reverse triangle inequality

$$|h_C(\mu) - h_C(\mu')| = |W_1^C(\mu, \nu_0) - W_1^C(\mu', \nu_0)|$$
$$= |W_1(\mu, \nu_0) - W_1(\mu', \nu_0)|$$
$$\leq W_1(\mu, \mu').$$

REMARK 3. *Avoiding risky distributions can be modeled by setting $W_1^C(\mu, \nu_0) = W_1(\mu, \nu_0) - C$ with the proof of Assumption 6 being equivalent to the above.*

**Weighted safety constraints.** In applications that require emphasis on certain regions of the state space, we can generalize the above safety constraints by introducing the weight function $w : \mathcal{S} \to \mathbb{R}_{\geq 0}$. We extend Appendix B.1 to weighted-differential entropy by defining

$$H^{w,\varepsilon}(\mu) := - \int_{\mathcal{S}} w(s) \log(\mu(s) + \varepsilon) \mu(ds)$$

and Appendix B.2 by considering weighted Wasserstein 1-distance

$$W_1^w(\mu, \nu_0) := \sup_{f : \text{Lip}(f) \leq 1} \int_{\mathcal{S}} w(s) f(s)(\mu - \nu_0)(ds).$$

Here, we use Kantorovic-Rubinstein dual definition of Wasserstein 1-distance

$$W_1(\mu, \nu_0) := \sup_{f : \text{Lip}(f) \leq 1} \int_{\mathcal{S}} f(s)(\mu - \nu_0)(ds),$$

where $f : \mathcal{S} \to \mathbb{R}$ is a continuous function and $\text{Lip}(f)$ denotes the minimal Lipschitz constant for $f$.

## C Implementation Details

In this section, we provide additional details on the practical implementation of SAFE-M$^3$-UCRL. In particular, we discuss *Probabilistic Neural Network Ensemble* model [17, 42] to implement the statistical model from Section 3.3 in Appendix C.1, the hallucinated control reparametrization from Section 4 in more detail in Appendix C.2, and optimization methods to solve Equation (9) in Appendix C.3

### C.1 Probabilistic Neural Network Ensemble Model of Transitions

As discussed in Section 3.3, we take a model-based approach to handling unknown transitions $f$. The representative agent learns the *Statistical Model* (see Section 3.3) of the transitions from the observed transitions $\cup_{i=1}^{n-1} \mathcal{D}_i$ at the beginning of each episode $n$, where $\mathcal{D}_i = \{(z_{i,t}, s_{i,t+1})\}_{t=0}^{T-1}$ with $z_{i,t} = (s_{i,t}, \mu_{i,t}, a_{i,t})$. We use *Probabilistic Neural Network Ensemble* [17, 42] that consists of $K$ neural networks parametrized by $\theta_k$ for $k \in \{1, \ldots, K\}$ (the episode index should be clear from the context so we omit it for the notation simplicity). Each neural network $f_{\theta_k}$ returns a mean vector, $\boldsymbol{m}_{\theta_k}(z) \in \mathcal{S} \subseteq \mathbb{R}^p$, and a covariance function $\Sigma_{\theta_k}(z) \in \mathbb{R}^{p \times p}$, that represents the aleatoric uncertainty. We further assume diagonal covariance functions. These outputs then form Gaussian distributions from which new states are sampled, i.e., $s_{t+1} \sim \mathcal{N}(\boldsymbol{m}_{\theta_k}(z_t), \Sigma_{\theta_k}(z_t))$ for $t = 0, \ldots, T - 1$. The models are trained with the negative log-likelihood loss function (NLL), $L(\theta) = -\sum_{\mathcal{D}_{1:n-1}} \log \mathbb{P}(s_{t+1} | z_t)$ as described in [42]. The ensemble means, and the aleatoric and epistemic uncertainties are then estimated as follows

$$\boldsymbol{m}_{n-1}(\cdot) = \frac{1}{K} \sum_{k=1}^{K} \boldsymbol{m}_{\theta_k}(\cdot)$$

$$\Sigma_{n-1}^e(\cdot) = \frac{1}{K-1} \sum_{k=1}^{K} (\boldsymbol{m}_{\theta_k}(\cdot) - \boldsymbol{m}_{n-1}(\cdot))(\boldsymbol{m}_{\theta_k}(\cdot) - \boldsymbol{m}_{n-1}(\cdot))^\top$$

$$\Sigma_{n-1}^a(\cdot) = \frac{1}{K} \sum_{k=1}^{K} \Sigma_{\theta_k}(\cdot),$$

where $\boldsymbol{m}_{n-1}(\cdot)$ is the ensemble prediction for the mean, while $\Sigma_{n-1}^a(\cdot)$ and $\Sigma_{n-1}^e(\cdot)$ denote the aleatoric and epistemic uncertainty estimates, respectively. Note that the epistemic uncertainty estimate $\Sigma_{n-1}^e(\cdot)$ is used to construct the calibrated model (see Section 3.3 and Assumption 4).

Even though Gaussian Processes (GPs) are proven to be calibrated under certain regularity assumptions [1, 67], we chose probabilistic neural network ensemble model due to their better practical performance, i.e., scalability to higher dimensions and larger datasets. The disadvantage of the probabilistic neural network ensemble is that, unlike GPs, it does not guarantee the calibrated model (see Assumption 4). Nevertheless, it can be recalibrated using one-step ahead predictions and temperature scaling as shown in [51]. We note that in our experiment (see Section 5), such recalibration was not needed, and the above-defined mean and epistemic uncertainty outputs were sufficiently accurate for training SAFE-M$^3$-UCRL (see Appendix E.3).

**Algorithm 2** Mean-Field Back-Propagation-Through-Time

---

**Input:** Safety constraint $h_C(\cdot)$, calibrated transitions $\tilde{f}_{n-1}$ represented by $\boldsymbol{m}_{n-1}(\cdot)$ and $\Sigma_{n-1}(\cdot)$, initial mean-field distribution $\mu_0$, known reward $r(\cdot)$, constants $C_{n,t}$, safety Lipschitz constant $L_h$; hyperparameter $\lambda$, number of epochs $K$, number of rollout steps $T$

1: Initialize $\boldsymbol{\pi}^\psi = (\pi_0^\psi, \ldots, \pi_{T-1}^\psi)$
2: **for** $k = 1, \ldots, K$ **do**
3:      Initialize $\tilde{\mu}_0 \leftarrow \mu_0$, $\tilde{s}_0 \sim \mu_0$
4:      Initialize $r \leftarrow 0$
5:      **for** $t = 0, \ldots, T-1$ **do**
6:          $\tilde{a}_t \leftarrow \pi_t^\psi(\tilde{s}_t, \tilde{\mu}_t)$
7:          $\tilde{s}_{t+1} \leftarrow \tilde{f}_{n-1}(\tilde{s}_t, \tilde{\mu}_t, \tilde{a}_t) + \varepsilon_t$
8:          $\tilde{\mu}_{t+1} \leftarrow U(\tilde{\mu}_t, \pi_t^\psi, \tilde{f}_{n-1})$
9:      **end for**
10:      Update $\psi$ with gradient ascent
11:      $\nabla_\psi \sum_{t=0}^{T-1} r(\tilde{s}_t, \tilde{\mu}_t, \tilde{a}_t) + \lambda \log(h_C(\tilde{\mu}_{t+1}) - L_h C_{n,t+1})$
12: **end for**
**Return** $\boldsymbol{\pi}_n^\psi \leftarrow \boldsymbol{\pi}^\psi$

---

## C.2 Hallucinated Control Implementation Trick

In Section 4, we introduced Safe-M$^3$-UCRL, a model that optimizes over the confidence set of transitions $\mathcal{F}_{n-1}$ and admissible policy profiles $\Pi$ (see Equation (7)). Unfortunately, optimizing directly over the function space is usually intractable since $\mathcal{F}_{n-1}$ is not convex, in general, [22]. Thus, to make the optimization tractable, we describe a *hallucinated control* trick, which leads to a practical reformulation (see Equation (9)). The structure in $\mathcal{F}_{n-1}$ allows us to parametrize the problem and use gradient-based optimization to find a policy profile $\boldsymbol{\pi}_n^*$ at every episode $n$. Namely, we use the mean-field variant of an established approach known as *Hallucinated Upper Confidence Reinforcement Learning* (H-UCRL) [19, 54, 58]. We introduce an auxiliary function $\eta : \mathcal{Z} \to [-1,1]^p$, where $p$ is the dimensionality of the state space $\mathcal{S}$, to define hallucinated transitions

$$\tilde{f}_{n-1}(z) = \boldsymbol{m}_{n-1}(z) + \beta_{n-1}\Sigma_{n-1}(z)\eta(z), \tag{11}$$

where $\boldsymbol{m}_{n-1}(\cdot)$ and $\Sigma_{n-1}(\cdot)$ are estimated from the past observations collected until the end of the previous episode $n-1$. 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 in Equation (11)

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

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

Furthermore, note that, for a fixed individual policy $\pi$, the auxiliary function $\eta(z) = \eta(s, \mu, \pi(s, \mu)) = \eta(s, \mu)$ has the same functional form as the policy $\pi$. This turns $\eta(\cdot)$ into a policy that exerts *hallucinated control* over the epistemic uncertainty of the confidence set of transitions $\mathcal{F}_{n-1}$ [19]. The reformulation of the optimization problem in Equation (9) allows us to optimize over parametrizable functions (e.g., *neural networks*) $\boldsymbol{\pi}$ and $\eta(\cdot)$ using gradient-based methods on the functions' parameters. Notice that the shared functional form of $\boldsymbol{\pi}$ and $\eta(\cdot)$ allows us to conveniently represent them by a single neural network. Further, note that the parametrization of $\eta(\cdot)$ must

be sufficiently flexible not to restrict the space of $\tilde{f}_{n-1}$. In Appendix C.3, we provide several algorithms that can be used to solve this optimization problem.

## C.3 Optimization Methods – MF-BPTT and MF-DDPG

In this section, we describe two algorithms to solve the optimization problem in Equation (9). Namely, in Appendix C.3.1, we outline the key steps to apply the mean-field variant of the *Back-Propagation-Through-Time* (BPTT) when a differentiable simulator is available and, in Appendix C.3.2, we describe the mean-field variant of the *Deep Deterministic Policy Gradient* (DDPG) [46] algorithm appropriate for non-differentiable simulators.

*C.3.1 Mean-Field Back-Propagation-Through-Time (MF-BPTT)*
Mean-Field Back-Propagation-Through-Time (MF-BPTT) assumes access to a differentiable simulator that returns a policy rollout given transition and policy functions. In each episode $n$, the representative agent initializes the policy profile $\boldsymbol{\pi}_n$, the auxiliary function $\eta(\cdot)$, and the estimated transitions $\tilde{f}_{n-1}$ using the mean $\boldsymbol{m}_{n-1}(\cdot)$ and covariance $\Sigma_{n-1}(\cdot)$ functions according to Equation (8). Then, the representative agent repeatedly calls the simulator with inputs $\boldsymbol{\pi}_n$ and $\tilde{f}_{n-1}$ that returns the episode reward defined in Equation (9a) as a differentiable object. After each policy rollout, a gradient ascent step is carried out on the parameters of $\boldsymbol{\pi}_n$ and $\eta(\cdot)$ before calling the simulator again. To simplify the notation, we overload the notation $\boldsymbol{\pi}_n^\psi$ with the combination of the two policy functions $\boldsymbol{\pi}_n$ and $\eta(\cdot)$, i.e., $\boldsymbol{\pi}_n^\psi = (\boldsymbol{\pi}_n^\psi, \eta^\psi)$ where the superscript $\psi$ represent the parameters of both functions. We outline the described steps in Algorithm 2. During the optimization phase (Line 3 in Algorithm 1), we optimize both functions, $\boldsymbol{\pi}_n$ and $\eta(\cdot)$, jointly. However, during the execution (Line 4 in Algorithm 1), we only use the outputs corresponding to the policy profile $\boldsymbol{\pi}_n$. The main distinction in Algorithm 2 compared to traditional BPTT lies inside the simulator that has to simulate the mean-field distribution $\tilde{\mu}_{n,t}$ for $t = 1, \ldots, T$ for each parameter update and calculate gradients with respect to this time dependency as well. More details on the implementation used for our experiments reported in Section 5 are provided in Appendix D.3 and Appendix E.2.

*C.3.2 Mean-Field Deep Deterministic Policy Gradient (MF-DDPG)*
In this section, we adopt the *Deep Deterministic Policy Gradient* (DDPG) algorithm [46] to the MFC. DDPG is a model-free actor-critic algorithm, hence, it can optimize Equation (9) without the assumption of a differentiable simulator. However, it can not be applied directly to the problem because the Q-value for $\tilde{z}_{n,t} = (\tilde{s}_{n,t}, \tilde{\mu}_{n,t}, \tilde{a}_{n,t})$ is ambiguous. The important insight here is that the value of a certain state $\tilde{s}_{n,t}$ of the environment reflects the whole population, i.e., the expected reward over the remainder of an episode for a given $\tilde{\mu}_{n,t}$ if every agent in every state $\tilde{s}_{n,t}$ chooses actions following $\pi_{n,t}$. In essence, the Q-value is a function of $\tilde{\mu}_{n,t}$ and $\pi_{n,t}$ and not $\tilde{z}_{n,t}$.

To overcome this issue, we introduce the *lifted mean-field Markov decision process* (MF-MDP) similarly to [13, 30, 31, 56, 58]. First, we rewrite the reward function as a function of the mean-field distribution and the policy, i.e.,

$$\tilde{r}(\tilde{\mu}_{n,t}, \pi_{n,t}) = \int_\mathcal{S} r(s, \tilde{\mu}_{n,t}, \pi_{n,t}(s, \tilde{\mu}_{n,t}))\tilde{\mu}_{n,t}(ds).$$

**Algorithm 3** Mean-Field Deep Deterministic Policy Gradient

---

**Input:** Safety constraint $h_C(\cdot)$, calibrated transitions $\tilde{f}_{n-1}$ represented by $\boldsymbol{m}_{n-1}(\cdot)$ and $\Sigma_{n-1}(\cdot)$, initial mean-field distribution $\mu_0$, known expected reward $\hat{r}(\cdot)$, constants $C_{n,t}$; safety Lipschitz constant $L_h$, hyperparameter $\lambda$, number of epochs $K$, number of rollout steps $T$, mini-batch size $B$, learning rate $\alpha$

1: Initialize $\boldsymbol{\pi}^\psi = (\pi_0^\psi, \ldots, \pi_{T-1}^\psi)$
2: Initialize $Q^\theta = (Q_0^\theta, \ldots, Q_{T-1}^\theta)$
3: Initialize $\theta' \leftarrow \theta$, $\psi' \leftarrow \psi$
4: Initialize replay buffer $R \leftarrow \emptyset$
5: **for** $k = 1, \ldots, K$ **do**
6:     Initialize $\tilde{\mu}_0 \leftarrow \mu_0$
7:     **for** $t = 0, \ldots, T - 1$ **do**
8:         $\tilde{\mu}_{t+1} \leftarrow U(\tilde{\mu}_t, \pi_t^\psi, \tilde{f}_{n-1})$
9:         $c_t \leftarrow \hat{r}(\tilde{\mu}_t, \pi_t^\psi) + \lambda \log(h_C(\tilde{\mu}_{t+1}) - L_h C_{n,t+1})$
10:        $R \leftarrow R \cup \{(\tilde{\mu}_t, c_t, \tilde{\mu}_{t+1})\}$
11:        Sample a mini-batch of $B$ random transitions $\{(\tilde{\mu}_i, c_i, \tilde{\mu}_{i+1})\}_{i=1}^B$ from $R$
12:        **for** $i = 1, \ldots, B$ **do**
13:            $q_i \leftarrow c_i + \gamma Q_t^{\theta'}(\tilde{\mu}_{i+1}, \boldsymbol{\pi}^{\psi'}(\cdot, \tilde{\mu}_{i+1}))$
14:        **end for**
15:        Update $\theta$ with gradient descent $\nabla_\theta \frac{1}{B} \sum_i (q_i - Q_t^\theta(\tilde{\mu}_i, \boldsymbol{\pi}^\psi(\cdot, \tilde{\mu}_i)))^2$
16:        Update $\psi$ with gradient ascent $\nabla_\psi \frac{1}{B} \sum_i Q_t^\theta(\tilde{\mu}_i, \boldsymbol{\pi}^\psi(\cdot, \tilde{\mu}_i))$
17:        $\theta' \leftarrow \alpha\theta + (1-\alpha)\theta'$
18:        $\psi' \leftarrow \alpha\psi + (1-\alpha)\psi'$
19:     **end for**
20: **end for**
Return $\boldsymbol{\pi}_n^\psi \leftarrow \boldsymbol{\pi}^{\psi'}$

---

To simplify the notation, we overload the notation $\boldsymbol{\pi}_n^\psi = (\boldsymbol{\pi}_n^\psi, \eta^\psi)$ as described in Appendix C.3.1. Then, we restate Equation (9) as

$$\psi^* = \arg\max_\psi \sum_{t=0}^{T-1} \tilde{r}(\tilde{\mu}_{n,t}, \pi_{n,t}^\psi) + \lambda \log(h_C(\tilde{\mu}_{n,t+1}) - L_h C_{n,t+1})$$

$$(12a)$$

$$\text{subject to} \quad \tilde{f}_{n-1}(\tilde{z}) = \boldsymbol{m}_{n-1}(\tilde{z}) + \beta_{n-1}\Sigma_{n-1}(\tilde{z})\eta^\psi(\tilde{z}) \quad (12b)$$

$$\tilde{\mu}_{n,t+1} = U(\tilde{\mu}_{n,t}, \pi_{n,t}^\psi, \tilde{f}_{n-1}), \quad (12c)$$

where $\tilde{\mu}_{n,0} = \mu_0$ for every $n$. The MF-MDP formulation in Equation (12) turns the MFC in Equation (9) into a Markov Decision Process on the state space of $\mathcal{P}(\mathcal{S})$ and action space $\{\pi : \mathcal{S} \times \mathcal{P}(\mathcal{S}) \to \mathcal{A}\}$ with deterministic transition function $U(\cdot)$. We define the Q-value as follows

$$Q_{n,t}(\tilde{\mu}_{n,t}, \boldsymbol{\pi}_n^\psi) = \sum_{j=t}^{T-1} \tilde{r}(\tilde{\mu}_{n,t}, \pi_{n,t}^\psi).$$

The Q-function above is parameterized by $\theta$ and denoted as $Q_{n,t}^\theta$ for episode $n$. The DDPG algorithm can now be stated for the lifted MF-MDP problem as outlined in Algorithm 3. The main learning loop consists of alternating updates to the policy $\boldsymbol{\pi}_n^\psi$ and the critic $Q_n^\theta$. The most recent version of the policy is executed in the environment to collect more transitions into the replay buffer.

Notice that Algorithm 3 uses a model-free approach to optimize the objective, which makes the exploration of particular importance. In comparison to traditional MDPs, MF-MDPs usually have a very constrained set of highly rewarding distributions, and most distributions offer poor rewards, which makes the exploration even more important. This is further complicated by what randomized actions imply in this scenario. In a traditional MDP, we can often assume, for instance, that executing random actions for a fixed number of initial steps would help in finding diversity in the reward space. This is not the case in MF-MDPs – depending on the granularity of the discretization that we use to represent probability distributions, we can expect the mean-field distribution to stay fairly stable. Similarly to [13], we might alleviate this issue by Gaussian mean-field initialization in each episode. Note, however, that this may not necessarily be appropriate in safety-constrained settings, as the initial mean-field distribution is expected to be safe. On top of that, we can add exploration noise to the actions obtained via the policy. Alternatively, we could add noise to the parameters of the policy network for a more consistent approach as in [59].

# D Experiments – Vehicle Repositioning

In this section, we provide further analysis, the motivation behind our modeling decisions, and details for making our experiments easily replicable. We use a private cluster with GPUs to run our experiments. SAFE-M³-UCRL and M³-UCRL under known transitions each used 15 minutes of one Intel Xeon Gold 511 CPU core, 32 GB of RAM and one Nvidia GeForce RTX 3090 GPU. Training SAFE-M³-UCRL and M³-UCRL under unknown transitions to produce results in Figure 2d, Figure 2f and Figure 8b had 24 hours of access to fifty Intel Xeon Gold 511 CPU cores, 64 GB of RAM and fifty Nvidia GeForce RTX 3090 GPUs during the batch execution necessary for training. The evaluation, i.e., generating the results for, e.g., Figure 9b and Figure 10 took around 1 hour of one Intel Xeon Gold 511 CPU core, 64 GB of RAM and one Nvidia GeForce RTX 3090 GPU. The only computationally intensive evaluation task was for Figure 9a for more than 1 million agents. We had access to Xeon Gold 511 CPU core, 128-256 GB of RAM, for around 8 hours. The implementation was predominantly done in Python packages PyTorch [57] and NumPy [34].

Our experimental workflow has the following structure:

(1) Input data preprocessing
- Estimating the demand distribution for the service $\rho_0$
- Estimating passenger's trip destinations' likelihood mapping $\Phi(\cdot)$

(2) Modeling assumptions, model parameters, and distributions' representation
- State-space, action space, noise, reward, safety constraint
- Mean-field distribution representation
- Mean-field transition $U(\cdot)$

(3) Executing model-based learning protocol (Algorithm 1)
- Optimizing Equation (9)
- Learning unknown transitions using probabilistic neural network ensemble, i.e., statistical estimators $\boldsymbol{m}_{n-1}(\cdot)$ and $\Sigma_{n-1}(\cdot)$

(4) Performance evaluation in the infinite regime

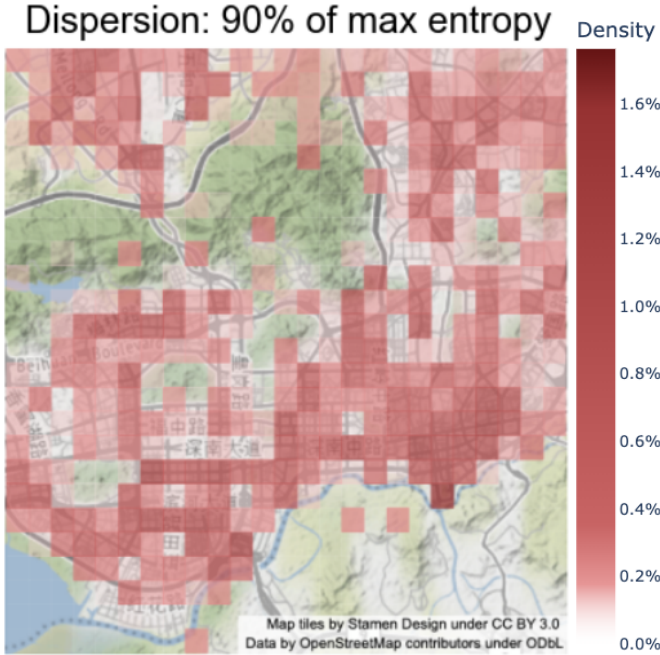(5) Performance evaluation in the finite regime

**Figure 4: Passenger's trip destinations' likelihood during evening peak hours given an arbitrary trip origin (see Appendix D.1).** We observe that the trips to most residential areas are almost equally likely and dispersed across the entire study region, which is consistent with the intuition that some residents commute back home (outside of the city center) after a work day, while others commute towards the city center for, e.g., leisure activities.

## D.1 Input Data Preprocessing

We consider vehicle trajectories collected in Shenzhen's extended city center with the geographical area spanning from 114.015 to 114.14 degrees longitude and from 22.5 to 22.625 degrees latitude. We have access to the trajectories of five full working weeks (Monday to Sunday) collected between $18^{th}$ January 2016 and $25^{th}$ September 2016. We restrict ourselves to evening peak hours between 19:00 and 22:00. We represent probability distributions by discretizing the space into $k \times k$ unit grid with $k = 25$ where each cell, $C_{ij} = [\frac{i}{k}, \frac{i+1}{k}\rangle \times [\frac{j}{k}, \frac{j+1}{k}\rangle$ for $i, j \in \{0, \ldots, k-1\}$, represents a square neighborhood of around $550 \times 550$ meters on the city map. We represent the service demand distribution $\rho_0 \in \mathcal{P}(\mathcal{S})$ as a $k \times k$ matrix where entries $[\rho_0]_{ij}, i, j \in \{0, \ldots, k-1\}$ represent a probability of a trip originating in the neighborhood $C_{ij}$. The probabilities are estimated as an average over the considered period and kept constant during each step $t$ of the learning protocol (see Appendix D.3). To smooth out the demand distribution and remove possible noise in the raw data, we apply 2-dimensional median smoothing with a window equal to 3 and show the output in Figure 3a. If the passenger's trip originates at the state $s \in \mathcal{S}$ the likelihood of its destinations is defined by the mapping $\Phi : \mathcal{S} \rightarrow \mathcal{P}(\mathcal{S})$ which we use in Appendix D.2 to define sequential transitions by first executing passenger trips followed by vehicle repositioning (see Figure 1). We flatten the $k \times k$ space grid into $k^2$-dimensional vector and

represent $\Phi(\cdot)$ as a $k^2 \times k^2$ probability matrix, i.e., rows summing to one represent outgoing mass for each cell. The entries $[\Phi]_{ij}$ with $i, j \in \{0, \ldots, k^2 - 1\}$ represent the likelihood of a passenger's trip that originated in the neighborhood $i$ ending in the neighborhood $j$. The likelihoods are estimated as an average over the considered period and kept constant. Figure 4 shows $\Phi$ column average, i.e., the trip destinations' likelihood given an arbitrary trip origin.

## D.2 Modeling Assumptions and Model Parameters

We represent our area of interest as a two-dimensional unit square, i.e., the state-space $\mathcal{S} = [0, 1]^2$, and assume that vehicles can move freely using repositioning actions taken from $\mathcal{A} = [-1, 1]^2$. If the action takes a vehicle outside of the state-space borders, we project it back onto the border. Since the fleet is operating in a noisy traffic environment, repositioning usually cannot be executed perfectly. We model the noise $\varepsilon_{n,t} \sim \text{TN}(0, \sigma^2 I_2)$ by a Gaussian with a known variance $\sigma^2$ truncated at the borders of $\mathcal{S}$ and $I_2$ is the $2 \times 2$ unit matrix. We use a standard deviation $\sigma = 0.0175$ to represent that the vehicle will be repositioned with a 68% probability within a circle with a 240-meter radius of the desired location or with 95% probability within a circle with a 480-meter radius. For simplicity, we assume that every passenger ride lasts fifteen minutes and that the repositioning is executed instantaneously. We assume that the representative agent (or multiple representative agents) reports the trajectories obtained during the interaction with the environment to the global controller at the end of the day. Therefore, we discretize our modeling horizon in fifteen-minute operational intervals, i.e., $T = 12$, and each episode $n$ represents one day. The goal of the service provider is maximizing the profit, which correlates with the amount of satisfied demand. Therefore, we model the service provider 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 $\rho_0$, i.e., $r(z_{n,t}) = -D_{KL}(\rho_0 \| \mu_{n,t})$. In other words, $r(\cdot)$ measures the similarity between vehicles' and demand distributions; the closer the distributions, the higher the profit. In practice, greedy profit maximization is often prevented by imposing service accessibility requirements by regulatory bodies. Due to the discrete representation of the mean-field distribution (discussed in the next paragraph), we use Shannon entropy to define the safety constraint

$$h_C(\mu_{n,t}) = -\sum_{i,j} \log([\mu_{n,t}]_{ij})[\mu_{n,t}]_{ij} - C \quad (13)$$

to enforce the service accessibility across all residential areas. Therefore, the optimization objective in Equation (9) trades off between maximizing the profit $r(\cdot)$ and adhering to accessibility requirements imposed by $h_C(\cdot)$. We use matrix $\Phi$ introduced Appendix D.1 to model sequential transitions by 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 $s_{n,t}^\Phi \sim \Phi(s_{n,t})$. Firstly, we find origin cell $i \in k^2$ such that $s_{n,t}$ resides in it and sample destination cell $j \in k^2$ given likelihoods defined in row $i$ of the probability matrix $\Phi$. Secondly, we determine the destination state $s_{n,t}^\Phi$ by uniform sampling from the destination cell $j$, which is a simplified model

of the passenger preferences of the final destination. Notice that the controller determines repositioning actions given intermediate states $s_{n,t}^{\Phi}$ obtained after executing passenger trips.

**Mean-field transitions.** During the learning/training phase, we assume that the number of agents $m \to \infty$ and that these agents induce the mean-field distribution $\mu_{n,t}$ for episode $n$ and step $t$. But, one of the major practical challenges is implementing the mean-field transition function $U(\cdot)$ (see Equation (2)). The main difficulties are representing the mean-field distribution $\mu_{n,t}$ and computing the integral in Equation (2). We use discretization to represent the mean-field distribution even though other representations, such as a mixture of Gaussians, are possible. Concretely, we represent the mean-field distribution as $\mu_{n,t} = [\mu_{n,t}]_{ij}$ with $i, j \in \{0, \dots, k-1\}$ with $k = 25$ by associating the probability $[\mu_{n,t}]_{ij} = \mathbb{P}[s_{n,t} \in C_{ij}]$ of the representative agent residing within each cell $C_{ij}$ during episode $n$ at step $t$. Note that the discrete representation of the mean-field distribution does not affect the state and action spaces, which remain continuous. The initial mean-field distributions, $\mu_{n,0}$ for every $n$, follow the uniform distribution which maximizes the Shannon entropy, ensuring the safety at the beginning of every episode $n$ at $t = 0$, i.e., $h_C(\mu_{n,0}) \geq 0$. In vehicle repositioning, the mean-field transitions $U(\cdot)$ consist of two sequential steps induced by transitions $f$. Namely, first, the demand shifts the mean-field distribution, followed by the transition induced by the controller. Formally, the mean-field demand transition is computed as $\mu_{n,t}^{\Phi} = (\mu_{n,t} \cdot p) \times \Phi + \mu_{n,t} \cdot (1-p)$, where $\cdot$ denotes elementwise multiplication, $\times$ denotes matrix multiplication and $p = \min(1, \frac{\rho_0}{\mu_{n,t}})$ represents elementwise proportion of occupied vehicles. The mean-field controller transition requires computing the integral in Equation (2) for which we use the discrete approximation given the points $c_{ij}$ uniformly chosen from cells $C_{ij}$ for $i, j \in \{0, \dots, k-1\}$

$$[\mu_{n,t+1}]_{ij} = \sum_{k,l} \mathbb{P}[f(c_{kl}, \mu_{n,t}^{\Phi}, \pi_{n,t}(c_{kl}, \mu_{n,t}^{\Phi})) + \varepsilon_{n,t} \in C_{ij}][\mu_{n,t}^{\Phi}]_{kl},$$

(14)

for the episode $n$ and step $t$. We assume that the noise term $\varepsilon_{n,t}$ is independent across episodes and steps as well as along the two dimensions while the truncation parameters are adjusted relative to the state space borders. Thus, we have the following

$$\mathbb{P}\left[f(c_{kl}, \mu_{n,t}^{\Phi}, \pi_{n,t}(c_{kl}, \mu_{n,t}^{\Phi})) + \varepsilon_{n,t} \in C_{ij}\right]$$

$$= \mathbb{P}\left[f(c_{kl}, \mu_{n,t}^{\Phi}, \pi_{n,t}(c_{kl}, \mu_{n,t}^{\Phi}))_x + \varepsilon_{n,t,x} \in \left[\frac{i}{k}, \frac{i+1}{k}\right)\right]$$

$$\times \mathbb{P}\left[f(c_{kl}, \mu_{n,t}^{\Phi}, \pi_{n,t}(c_{kl}, \mu_{n,t}^{\Phi}))_y + \varepsilon_{n,t,y} \in \left[\frac{j}{k}, \frac{j+1}{k}\right)\right]$$

$$= \left[\phi\left(\frac{i+1}{k} - f(c_{kl}, \mu_{n,t}^{\Phi}, \pi_{n,t}(c_{kl}, \mu_{n,t}^{\Phi}))_x\right)\right.$$

$$\left. - \phi\left(\frac{i}{k} - f(c_{kl}, \mu_{n,t}^{\Phi}, \pi_{n,t}(c_{kl}, \mu_{n,t}^{\Phi}))_x\right)\right]$$

$$\cdot \left[\phi\left(\frac{j+1}{k} - f(c_{kl}, \mu_{n,t}^{\Phi}, \pi_{n,t}(c_{kl}, \mu_{n,t}^{\Phi}))_y\right)\right.$$

$$\left. - \phi\left(\frac{j}{k} - f(c_{kl}, \mu_{n,t}^{\Phi}, \pi_{n,t}(c_{kl}, \mu_{n,t}^{\Phi}))_y\right)\right],$$

(15)

where $\phi(\cdot)$ is the cumulative distribution function of truncated Gaussian $\text{TN}(0, \sigma^2)$.

**Mean-field transitions in the finite regime** In Appendix D.5, we instantiate a finite number of vehicles $m < \infty$ to evaluate the policy performance in a realistic setting. We keep track of vehicles' states $s_t^{(l)}$ for every vehicle $l \in \{1, \dots, m\}$ at steps $t = 0, \dots, T$. In this setting, we approximate the mean-field transition $U(\cdot)$ with the normalized two-dimensional histogram $[\mu_{t+1}]_{ij}$ with bins defined by the cells $C_{ij}$ for $i, j \in \{0, \dots, k-1\}$ given vehicles' next states $s_{t+1}^{(l)} = f(s_t^{(l)}, \mu_t^{\Phi}, \pi(s_t^{(l)}, \mu_t^{\Phi})) + \varepsilon_t$.

## D.3 Model-Based Learning Protocol

We use the learning protocol introduced in Algorithm 1 to train SAFE-M³-UCRL. For hyperparameters of the model, see Table 1.

To optimize the objective in Equation (9) in the subroutine in Line 3, we use MF-BPTT (see Appendix C.3.1). We parametrize the policy via a fully-connected neural network with two hidden layers of 256 neurons and Leaky-ReLU activations. The output layer returns the agents' actions using Tanh activation. We use Xavier uniform initialization [29] to randomly initialize weights while we set bias terms to zero. We use 20,000 training epochs with the early stopping if the policy does not improve at least 0.5% within 500 epochs. To prevent gradient explosion, we use L2-norm gradient clipping with max-norm set to 1. Note that in our experiments, a single neural network had enough predictive power to represent the whole policy profile $\boldsymbol{\pi} = (\pi_0, \dots, \pi_{T-1})$, but using $T$ networks, one for each individual policy $\pi_t$ is a natural extension. For further details about hyperparameters, see Table 2. Additionally, note that we parametrize policy by a neural network with Lipschitz continuous activations (Leakly-ReLU and Tanh). In the case of the bounded neural network's weights, policy satisfies the Lipschitz continuity assumption (see Assumption 2). In practice, we bound the weights via L2-regularization.

We estimate the confidence set of transitions $\mathcal{F}_{n-1}$ in the subroutine in Line 5 using a probabilistic neural network ensemble (see Appendix C.1). We use an ensemble of 10 fully-connected neural networks with two hidden layers of 16 neurons and Leaky-ReLU activations. The output layer returns the mean vector and variance vector (because of the covariance matrix diagonality assumption) of the confidence set; the mean uses linear activation, and the variance uses Softplus activation. We minimize the negative log-likelihood (NLL) for each neural network under the assumption of heteroscedastic Gaussian noise as described in [42]. We randomly split the data from the replay buffer into a training set (90%) and a validation set (10%). We use 10,000 training epochs with the early stopping if the performance on the validation set does not improve for at least 0.5% within 100 consecutive epochs. For further details about hyperparameters, see Table 3.

## D.4 Performance Evaluation in the Infinite Regime

In this section, we extend the experiments presented in Section 5.2. We assume that the number of vehicles $m \to +\infty$ and show two important results. First, we show a conservative behavior of SAFE-M³-UCRL by training model with Shannon entropy safety constraint (see Remark 2) with the safety threshold $C$ set to $p = 0.50$ of the maximum Shannon entropy, i.e., $C = p \log(k^2)$. We then evaluate its performance against a much higher safety threshold, i.e., against

the safety threshold induced by $p = 0.85$. Figure 8a shows that the model never violates stricter safety constraint regardless of its training in a weaker setting. Second, we show the data efficiency of the learning protocol (Algorithm 1) by training models with access to one, five, and ten representative agents for data collection. Figure 8b shows that the model under unknown transitions converges to the model trained under known transitions almost 6 times faster when using ten representative agents instead of one representative agent. It is a very useful result that can be utilized in many applications. For example, in most transportation applications, using tens or even hundreds of representative agents often does not cause cost-related issues. On the contrary, it might be more cost-effective to have ten representative agents for a month than one representative agent for a year.

## D.5 Performance Evaluation in the Finite Regime

In this section, we assume a finite number of vehicles $m < +\infty$ and approximate the mean-field distribution

$$\mu_t(s) = \lim_{m \to \infty} \frac{1}{m} \sum_{i=1}^{m} \mathbb{I}(s_t^{(i)} = s)$$

with the empirical distribution as explained in Appendix D.5. The policy profile $\pi_n^*$ trained in the infinite regime is used to reposition each of $m$ individual vehicles in the fleet. In Figure 9a, we show the relationship between model performance and the number of vehicles in the system. As expected, we observe that increasing the number of vehicles leads to better performance due to the increased precision of mean-field distribution $\mu_t$ approximation at step $t$. By performing 100 randomly initialized runs for each of various fleet sizes, we see that SAFE-M$^3$-UCRL learned under unknown transitions and applied in the finite regime converges to the solution achieved in the infinite regime under known transitions. Furthermore, the model performs very well already for a fleet of 10,000 vehicles, which is in the order of magnitude of the fleet size that operates in Shenzhen. Concretely, in 2016, the fleet had around 17,000 vehicles, with an increasing trend. We also observe that in the finite regime, the difference in performance between SAFE-M$^3$-UCRL trained under known and unknown transitions becomes insignificant. To showcase the practical usefulness of the algorithm, in Figure 9b, we instantiate 10,000 vehicles and display their positions after the final repositioning at step $T = 12$. We observe that the majority of vehicles are repositioned to areas of high demand. At the same time, some of them are sent to residential zones in the northwest and northeast to enforce accessibility. It is important to note that some vehicles are repositioned to aquatic areas and areas without infrastructure due to two reasons. First, our model guarantees global safety without explicit guarantees for individual/local safety. Notice that undesirable areas might be avoided by safety constraint shaping, e.g., by setting the weight function for these areas to zero as discussed in Appendix B.2. Second, the model loses some of its accuracy due to the finite regime approximation errors. In Figure 10, we observe only a slight decrease in dispersion when SAFE-M$^3$-UCRL finite regime approximation is compared to the infinite regime performance. To conclude, we showcase the potential of SAFE-M$^3$-UCRL for vehicle repositioning in the finite regime, which might be a positive signal for real-world practitioners.

## E  Experiments – Swarm Motion

In this section, we extend the swarm motion experiments discussed in Section 5.1 and complement the vehicle repositioning experiments elaborated in Section 5.2 and Appendix D. For this experiment, we use the same private cluster and approximately the same amount of computational resources as reported in Appendix D.

### E.1  Modeling Assumptions and Model Parameters

We model the state space $\mathcal{S}$ as the unit torus on the interval $[0, 1]$ and set the action space as the interval $\mathcal{A} = [-7, 7]$ due to the knowledge of the range of actions from the continuous-time analytical solution [2]. 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 $\varepsilon_{n,t} \sim \mathrm{N}(0, \Delta t)$ for all episodes $n$ and steps $t$. 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 $s$ (see Figure 5), 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, $\Delta t \to 0$, can be obtained analytically. Namely, for the infinite time horizon, i.e., $T \to \infty$, we have

$$\pi^*(s, \mu) = 2\pi \cos(2\pi s)$$
$$\mu^*(s) = \frac{e^{2\sin(2\pi s)}}{\int_{\mathcal{S}} e^{2\sin(2\pi s')}ds'}, \tag{16}$$

where $\pi^*$ and $\mu^*$ form an ergodic solution satisfying $\mu^* = U(\mu^*, \pi^*, f)$. We use $\mu^*$ as a benchmark but note that it might no longer be an optimal solution in the discrete-time setting. Therefore, discrete-time solutions obtained under known transitions serve as a good benchmark for SAFE-M$^3$-UCRL performance under unknown transitions. To control overcrowding, we use the Shannon-entropic constraint introduced in Equation (13) instead of having the overcrowding penalty term $\log(\mu_{n,t})$ in the reward. As discussed in Appendix D.2, Shannon entropy is used because of the discrete mean-field distribution representation. Since higher entropy translates into less overcrowding, we can upfront determine and upper-bound the acceptable level of overcrowding by setting a desirable threshold $C$. Similarly to the discussion in Appendix D.2, we represent the mean-field distribution by discretizing the state space into $k = 100$ uniform intervals and assigning the probability of the representative agent residing within each of them. To compute the mean-field transitions $U(\cdot)$, we use one-dimensional equivalent of Equation (14) and Equation (15). We set the safety threshold $C$ as a proportion $p \in [0, 1]$ of the maximum Shannon entropy Equation (13), i.e., $C = p\log(k)$. We initialize safe mean-field distributions $\mu_{n,0}$ as uniform distributions since they maximize Shannon entropy, which makes them safe for every threshold $C$.

### E.2  Model-Based Learning Protocol

We follow the same procedure as described in Appendix D.3 with the hyperparameters from Tables 4 to 6. The only difference compared
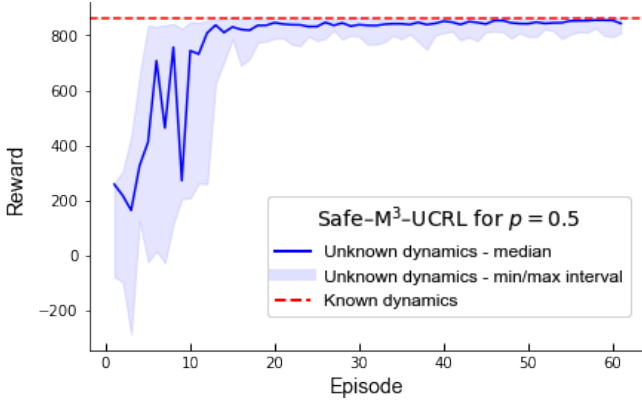
Figure 6: SAFE-M³-UCRL learning curve for $p = 0.5$ for swarm motion.

to Appendix D.3 is the increase in the complexity of the computational graph because of the high number of steps $T$. Therefore, we use batch normalization [38] to prevent vanishing gradients.

### E.3 Performance Evaluation

In this section, we complement the results shown in Section 5.1. We first visualize the positional reward $\phi(\cdot)$ in Figure 5 for the ease of interpretation of the obtained results. We see that the reward has two local maxima, but due to a significant difference in their value, unconstrained benchmarks ignore the lower maxima. On the other hand, SAFE-M³-UCRL for $p = 0.95$, and to a certain extent for $p = 0.5$, take advantage of it by reducing the kinetic energy in the neighborhood of lower maxima as shown in Figure 2e. In Figure 12,
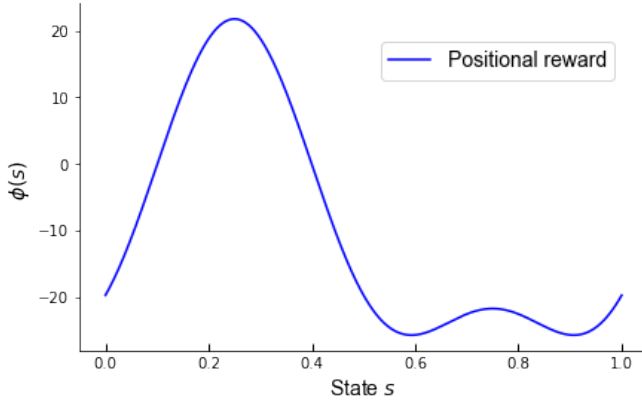


Figure 5: Swarm motion positional reward $\phi(\cdot)$.

we show the mean-field distributions progression over time guided

by policies learned by SAFE-M³-UCRL for $p = 0.5$ and $p = 0.95$. We observe that for $p = 0.5$, we reach near-stationary distribution after only 10 steps (see Figure 12a), i.e., the distribution remains the same until the algorithm terminates at $T = 100$. For $p = 0.95$, we reach stationarity even faster, as shown in Figure 12b. The learning process presented in Figure 2a is explained by the reduction of the epistemic uncertainty in the estimated transitions $\tilde{f}_{n-1}$. Before the first episode $n = 1$, the statistical model is estimated only from trajectories collected by randomly initialized policy $\pi_0$. Due to the high epistemic uncertainty in regions that random policy did not explore, upper-confidence hallucinated transitions Equation (8) do not approximate well true transitions (see Figure 11a). By episode $n = 5$, the model already has a good approximation of the transitions (see Figure 11b), while at episode $n = 50$, the transitions are known with near-certainty (see Figure 11c). These results coincide with the observation in Figure 2a where around episode $n = 50$, SAFE-M³-UCRL starts obtaining the results as if the transitions were known. Note that we implement the toroidal state space on $\mathcal{S} = [0, 1]$ by assuming a sufficiently large extension, e.g., $[-1, 2]$, of the interval over its borders such that a new state resulting from any possible action is captured with high probability. The new state is then mapped back to interval $[0, 1]$ using the modulo operation. For completeness of the analysis, in Figure 6, we show that SAFE-M³-UCRL for $p = 0.5$ converges to the result obtained under known transitions. We further validate the observation presented in Figure 2b that the constraint for $p = 0.5$ results in similar overcrowding as the reward penalty term $-\log(\mu)$. Namely, in Figure 7, we see that SAFE-M³-UCRL for $p = 0.5$ and M³-UCRL with overcrowding penalty term satisfy the safety constraint $h_C(\mu) = 0.5 \log(k)$ with a similar margin.
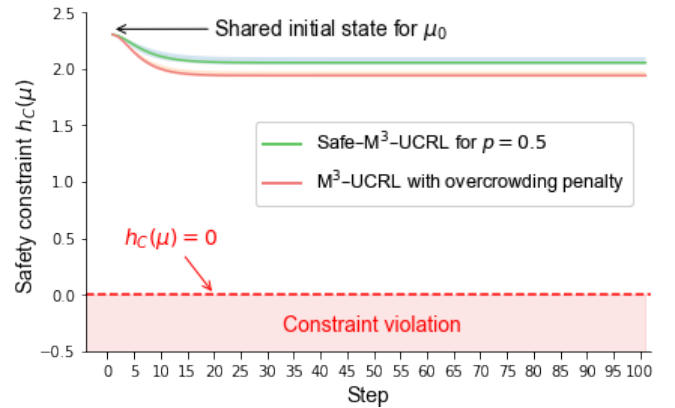


Figure 7: Swarm motion safety for $p = 0.5$.

(a) Safe-M³-UCRL conservative behavior



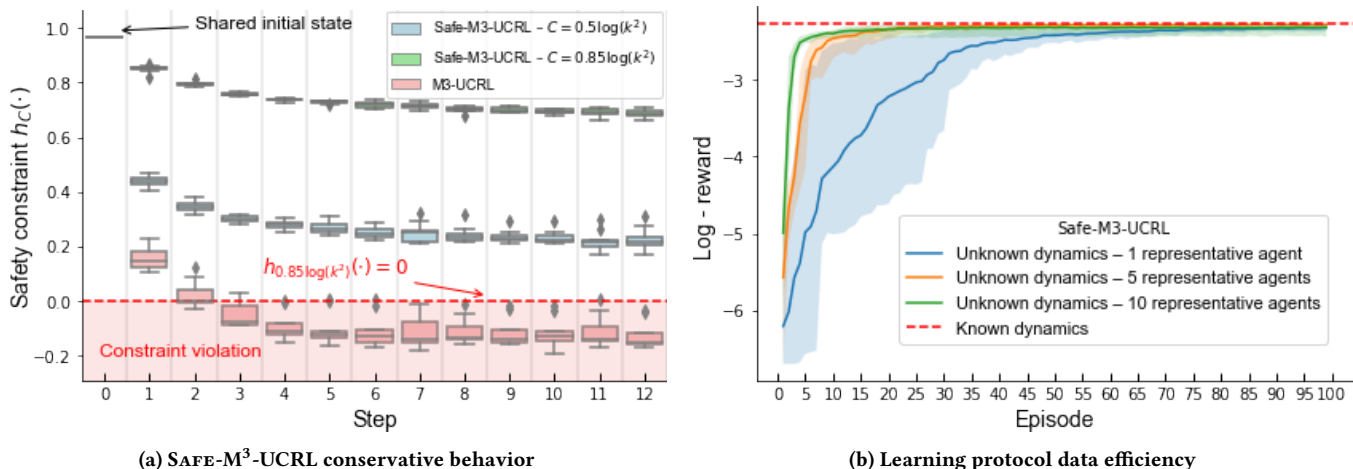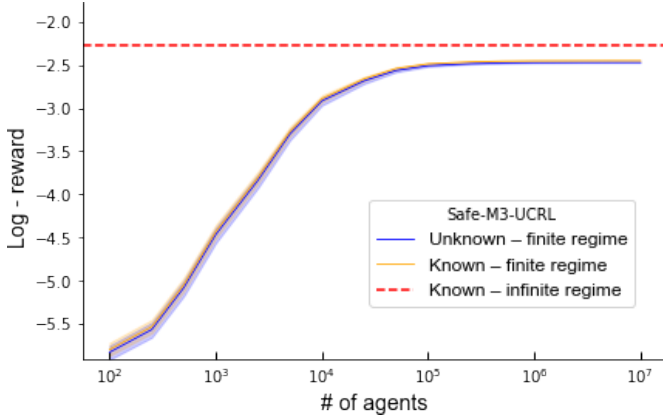(b) Learning protocol data efficiency

Figure 8: We showcase Safe-M³-UCRL conservative behavior and data efficiency by training 10 randomly initialized policy profiles and statistical models where each setup uses the entropic safety constraint $h_C(\cdot) \geq 0$. We set the safety threshold $C$ as a proportion $p$ of the maximum Shannon entropy, i.e., $C = p \log(k^2)$ with $k = 25$. In (a), the policy profiles trained for satisfying $h_{0.5 \log(k^2)}(\cdot) \geq 0$ never violate $h_{0.85 \log(k^2)}(\cdot) \geq 0$, which shows the conservative behavior of our model. In (b), we show the data efficiency of the learning protocol (Algorithm 1) by comparing learning curves observed during training models that satisfy $h_{0.85 \log(k^2)}(\cdot) \geq 0$ when using one, five and ten representative agents (RA) for data collection. We see that the model trained with 1-RA converges to the performance of the model under known transitions in around 80 episodes, while it takes 25 and 15 episodes for 5-RA and 10-RA models, respectively. Note that we use log-reward to emphasize learning speeds on a comparable scale.

Table 1: Learning protocol hyperparameters for vehicle repositioning

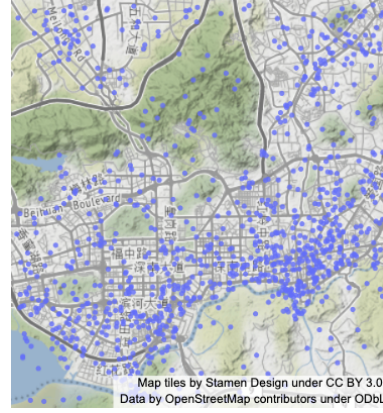| Hyperparameter | Value | Description |
|---|---|---|
| $N$ | 200 | Number of episodes |
| $T$ | 12 | Number of steps |
| $k$ | 25 | Number of discretization segments per axis |
| $\sigma$ | 0.0175 | Standard deviation of the system noise |
| $L_h$ | 0.1 | Lipschitz constant in Equation (9) |
| $\lambda$ | 1 | Log-barrier hyperparameter in Equation (9) |
| # of representative agents | 1 to 10 | By default 1, but in some experiments, we use more |

Table 2: Policy hyperparameters for vehicle repositioning

| Hyperparameter | Value | Description |
|---|---|---|
| # of hidden layers | 2 | |
| # of neurons | 256 | Number of neurons per hidden layer |
| hidden activations | Leaky-ReLU | |
| output activation | Tanh | |
| $\alpha$ | $10^{-4}$ | Learning rate |
| $w$ | $5 \cdot 10^{-4}$ | Weight decay |
| initialization | Xavier uniform | |
| bias initialization | 0 | |
| $n$ | 20,000 | Number of epochs |
| early stopping | 0.5% | If not improved after 500 epochs |

(a)

Objective from Equation (9) achieved in the finite regime by
approximating mean-field distribution with empirical distribution



(b)

Vehicles' locations after repositioning action
in the finite regime

Figure 9: To showcase Safe-M$^3$-UCRL performance in the finite regime, we instantiate a finite number of vehicles, each following a policy profile $\pi_n^*$ learned in the infinite regime and satisfying the entropic safety constraint $h_C(\cdot) \geq 0$ for $C = 0.85 \log(k^2)$ with $k = 25$. We perform 100 randomly initialized runs for each of the various fleet sizes. In (a), we see that increasing the number of vehicles leads to better performance due to the increased precision of mean-field distribution approximation. Further, we see that Safe-M$^3$-UCRL learned under unknown transitions and applied in the finite regime converges to the solution achieved in the infinite regime under known transitions. We also see that the curves for Safe-M$^3$-UCRL trained under unknown and known transitions almost overlap, i.e., the value of knowing transitions has little to no insignificance in the finite regime. In (b), we showcase the performance for the realistic number of vehicles operating in Shenzhen (10,000 to 20,000). We display the positions of a randomly chosen subset of 1,000 vehicles (out of 10,000) at step $T = 12$ after the final repositioning. We observe that the majority of vehicles are repositioned to areas of high demand, while some of them are sent to residential zones in the northwest and northeast to enforce accessibility.

Table 3: Probabilistic neural network ensemble hyperparameters for vehicle repositioning

| Hyperparameter | Value | Description |
|---|---:|---|
| # of ensemble members | 10 | |
| # of hidden layers | 2 | |
| # of neurons | 16 | Number of neurons per hidden layer |
| hidden activations | Leaky-ReLU | |
| mean output activation | Linear | |
| variance output activation | Softplus | |
| $\alpha$ | $10^{-4}$ | Learning rate |
| $w$ | $5 \cdot 10^{-4}$ | Weight decay |
| $\beta$ | 1 | Assumption 4 hyperparameter |
| initialization | Xavier uniform | |
| bias initialization | 0 | |
| $n$ | 10,000 | Number of epochs |
| early stopping | 0.5% | If not improved for 100 consecutive epochs |
| train-validation split | 90%-10% | We use a validation set for early stopping |
| replay buffer size | 10-100 | By default 100 |
| batch size | 8-128 | Increasing with the replay buffer size |

|                                                                    | (a)                                                     |                                                                   | (b)                                                   |
|:---:|:---:|

**Safe-M³-UCRL under unknown transitions in the infinite regime**    **Safe-M³-UCRL under unknown transitions in the finite regime**
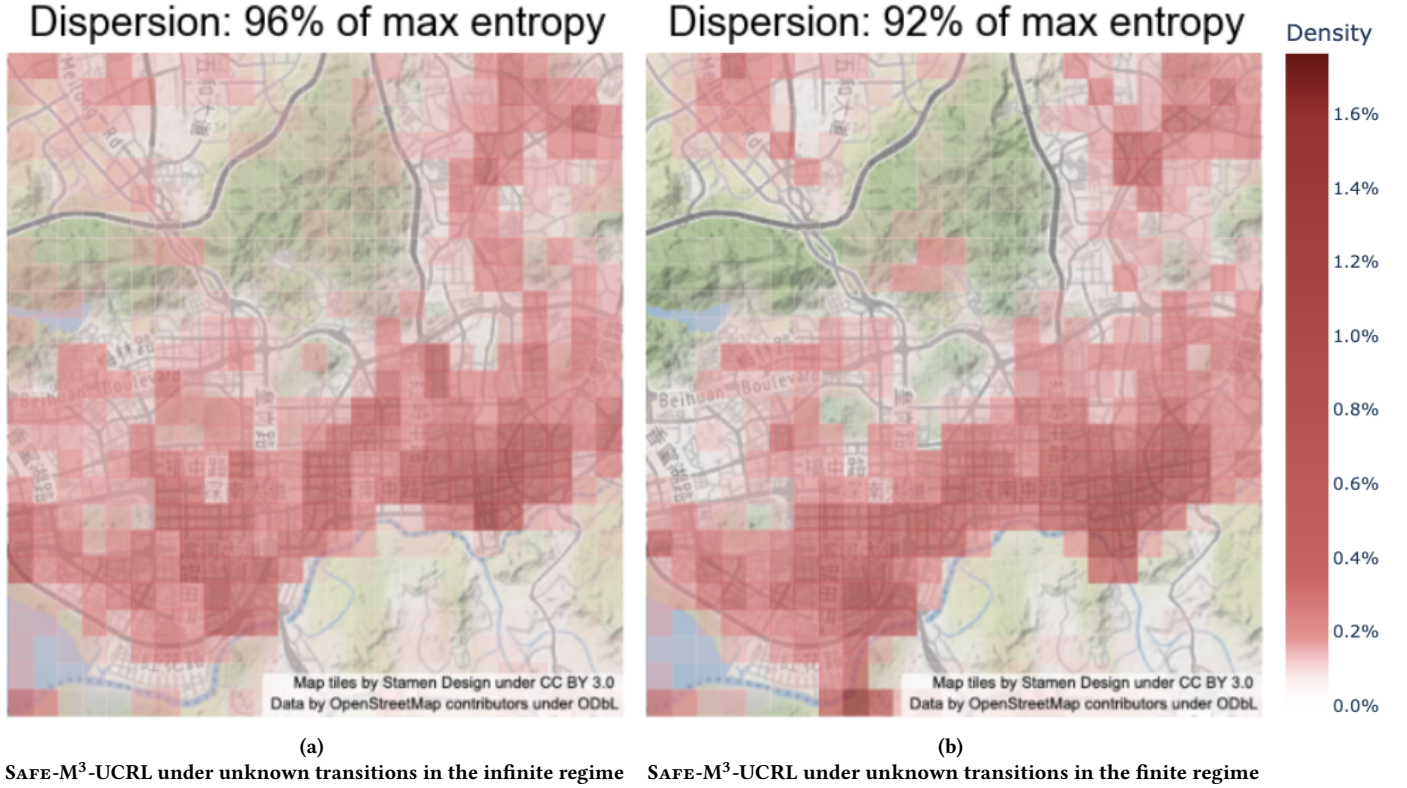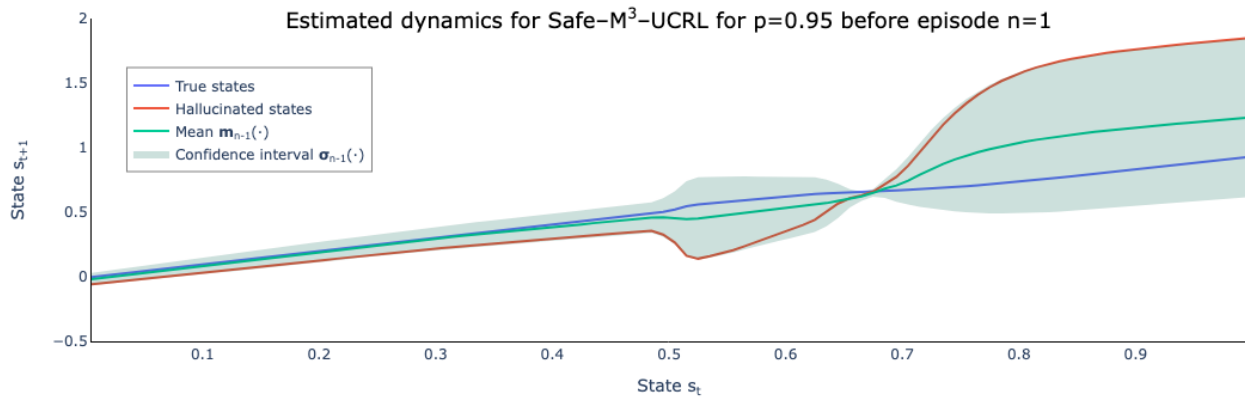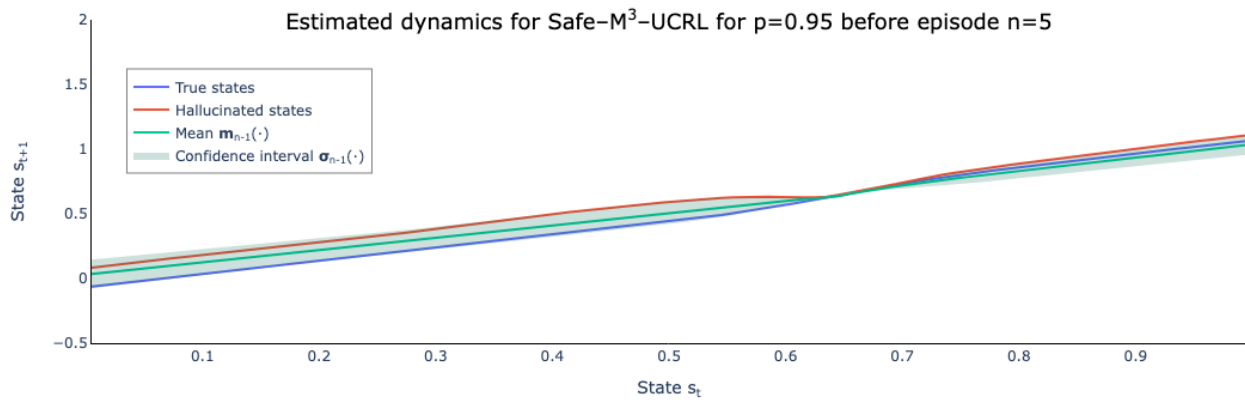
**Figure 10: We show the difference in the performance introduced by approximating mean-field distribution with the finite number of vehicles as described in Appendix D.2. In (a), we see that the dispersion of Safe-M³-UCRL in the infinite regime is $p = 0.96$ as elaborated in Figure 3. In (b), the policy profile $\pi_n^*$ trained in the infinite regime is used to control a fleet of 10,000 vehicles. We observe a slight decrease in the dispersion from $p = 0.96$ to $p = 0.92$ due to the mean-field distribution approximation errors.**

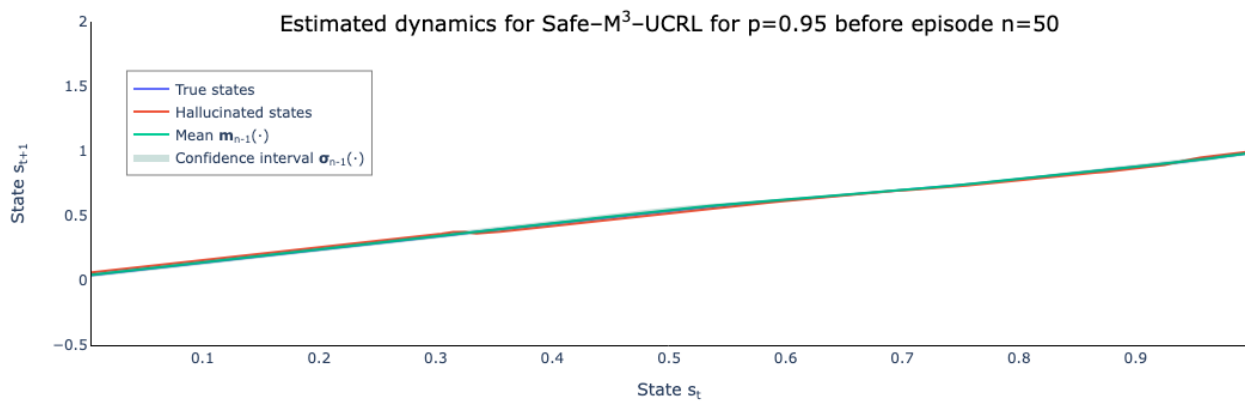**Table 4: Learning protocol hyperparameters for swarm motion**

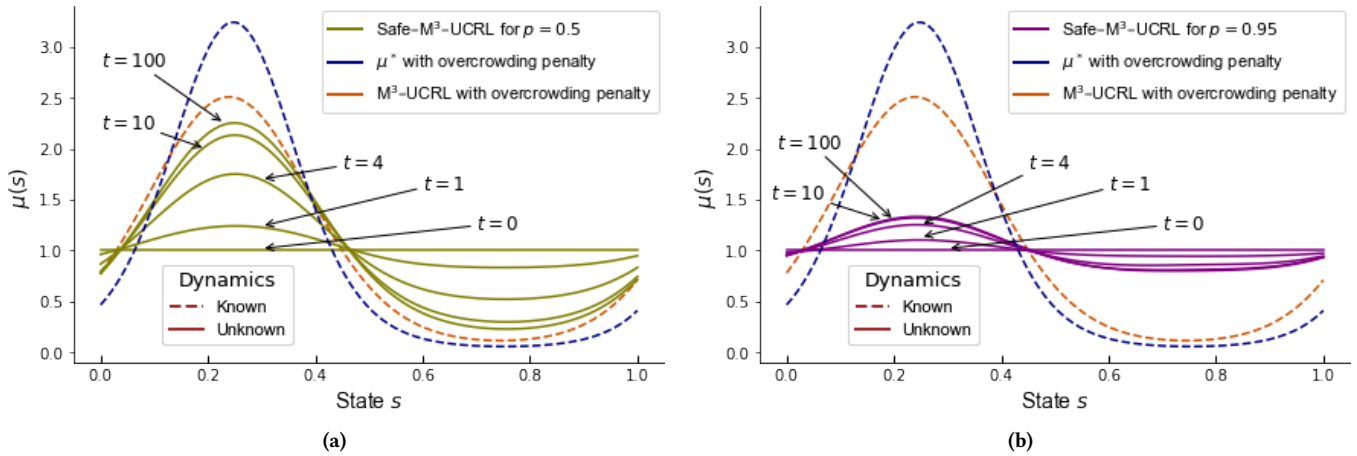| Hyperparameter            | Value     | Description                                   |
|---------------------------|-----------|-----------------------------------------------|
| $N$                       | 200       | Number of episodes                            |
| $T$                       | 100       | Number of steps                               |
| $k$                       | 100       | Number of discretization segments per axis    |
| $\sigma$                  | 1         | Standard deviation of the system noise        |
| $L_h$                     | $10^{-4}$ | Lipschitz constant in Equation (9)            |
| $\lambda$                 | 15        | Log-barrier hyperparameter in Equation (9)    |
| # of representative agents | 1        |                                               |

**(a)**



**(b)**



**(c)**

**Figure 12: Mean-field distributions progression over time.**

**Table 5: Policy hyperparameters for swarm motion**

| Hyperparameter | Value | Description |
| --- | ---: | --- |
| # of hidden layers | 2 | |
| # of neurons | 16 | Number of neurons per hidden layer |
| hidden activations | Leaky-ReLU | |
| output activation | Tanh | |
| $\alpha$ | $5 \cdot 10^{-3}$ | Learning rate |
| $w$ | $5 \cdot 10^{-4}$ | Weight decay |
| initialization | Xavier uniform | |
| bias initialization | 0 | |
| $n$ | 50,000 | Number of epochs |
| early stopping | 0.5% | If not improved after 100 epochs |

**Table 6: Probabilistic neural network ensemble hyperparameters for swarm motion**

| Hyperparameter | Value | Description |
| --- | ---: | --- |
| # of ensemble members | 10 | |
| # of hidden layers | 2 | |
| # of neurons | 16 | Number of neurons per hidden layer |
| hidden activations | Leaky-ReLU | |
| mean output activation | Linear | |
| variance output activation | Softplus | |
| $\alpha$ | $5 \cdot 10^{-3}$ | Learning rate |
| $w$ | $5 \cdot 10^{-4}$ | Weight decay |
| $\beta$ | 1 | Assumption 4 hyperparameter |
| initialization | Xavier uniform | |
| bias initialization | 0 | |
| $n$ | 10,000 | Number of epochs |
| early stopping | 0.5% | If not improved for 30 consecutive epochs |
| train-validation split | 90%-10% | We use a validation set for early stopping |
| replay buffer size | 10,000 | |
| batch size | 8-512 | Increasing with the replay buffer size |