

# MASDIFF: LARGE SCALE MULTI-AGENT SYSTEM EMERGENCE CONTROL VIA EVOLUTIONARY DIFFUSION

**Anonymous authors**

Paper under double-blind review

## ABSTRACT

Reward model learning methods are primarily divided into implicit reward modeling (IRM) and explicit reward modeling. Implicit reward modeling aims to learn the intrinsic reward of each agent to facilitate better exploration while explicit reward modeling (ERM) aims to learn the behavioral preferences of agents. The biggest difference between implicit and explicit reward modeling is that ERM is transferable to other scenarios but IRM cannot. Currently, few methods can simultaneously archive the global target and also learn the ERM. However, the problem addressed in this paper requires the integration of the two objectives. This paper use the diffusion model to generate the expert data for learning ERM. Since the traditional diffusion model can only generate data according to the given expert data, we introduce evolutionary diffusion model to generate data in the absence of any expert data. To steer the collaboration of all agents towards the specified macro-level objective, the macro-level objective is adopted as the fitness for the population. This mechanism transforms the multi-agent exploration based on intrinsic reward into the evolutionary exploration based on genetic operators. Moreover, the optimal individual retention strategy in the evolutionary diffusion model can address the non-stationary problem in MARL. In the experiments, we demonstrate that MASDiff can simultaneously archive the two objectives. Furthermore, we demonstrate that the ability to conduct counterfactual reasoning with the transferable ERM in different scenarios. We propose several ‘What if’ questions to indicate the change of scenarios and obtain relatively accurate counterfactual reasoning results.

## 1 INTRODUCTION

Reward model learning methods are primarily divided into implicit reward modeling and explicit reward modeling. Implicit reward modeling aims to learn the intrinsic reward of each agent to facilitate better exploration, ultimately achieving the overall outcomes. Intrinsic reward shaping (Cao et al., 2023), self-supervised RL (Pecháč et al., 2024) and credit assignment (Zhou et al., 2020) can be classified as implicit reward modeling. Explicit reward modeling aims to learn the behavioral preferences of agents by explicitly learning a reward model. The purpose is to ensure that the distribution of behavioral patterns derived from the reward model aligns with the statistical distribution observed in the expert data. Inverse RL (Wang et al., 2024), behavior cloning (Zhou et al., 2024) and Reinforcement Learning with Human Feedback (RLHF) (Casper et al., 2023) can be classified as implicit reward modeling. Currently, there is hardly any method capable of simultaneously achieving the two aforementioned objectives, namely, learning explicit reward models to realize macro objectives in multi-agent systems. This is primarily because existing approaches are mainly designed to address conventional reinforcement learning problems, which do not require the integration of these two objectives. However, the problem addressed in this paper requires the integration of the two objectives, i.e. to learn the explicit reward model for each agent on the one hand and to archive the overall outcome. We first present the specific problem to be solved in this study. Then, we give the contributions of this paper.

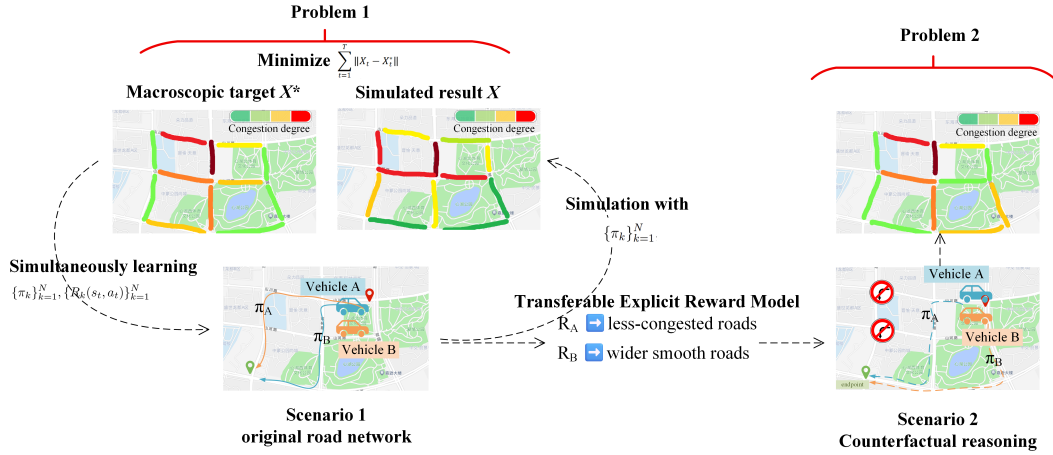


Figure 1: Example of the problems to be addressed by this paper

Consider the following two problems illustrated in Fig 1:

**Problem 1:** In scenario 1, can we learn the navigation policy of each vehicle such that the simulated congestion conditions aligns with the real congestion status? We assume that the number of vehicles is known and the origin-destination pairs for each vehicle are given.

**Problem 2:** In scenario 2, what the congestion conditions will be if we prohibit the left turns at two intersections?

If only problem 1 needs to be solved, we can use the implicit reward modeling approach in MARL. This method can learn the navigation policy instead of the explicit reward model for each vehicle to reconstruct the real congestion condition. But in order to solve problem 2 which is more important, the explicit reward model for each vehicle also need to be learned when solving problem 1. Because although the vehicle’s optimal route may be different in various scenarios, its decision preference remains the same. For example, vehicle A prefers less-congested roads, while vehicle B prefers wider arterial roads. In scenario 1, vehicle A’s route is shown as the blue line and vehicle B’s as the orange line. In scenario 2, vehicle A stays on the same route (blue dashed line), while vehicle B shifts to another road (orange dashed line). So to learn the explicit reward model, we need to obtain the expert data for each vehicle in scenario 1 which is transferable in scenario 2. To answer the what-if questions question 2, vehicles simply need to learn the new navigation policy based on the reward model as illustrated in Figure 1. In real world, many issues are similar to problem 1 and 2 (Jiaoling et al., 2025). Problem 1 can be generalized as Formula 1.

$$\begin{aligned}
 \{\pi_k\}_{k=1}^N, \{R_k(s_t, a_t)\}_{k=1}^N = \text{Argmin} \sum_{t=1}^T \|X_t - X_t^*\| \\
 \text{s.t.} \\
 \pi_k = \arg \max \mathbb{E} \left[ \sum_{t=1}^T \gamma^t r_{k,t} \right] \quad \text{and} \quad r_{k,t} = R_k(s_t, a_t)
 \end{aligned} \tag{1}$$

We need to learn  $\{\pi_k\}_{k=1}^N$  to cooperatively minimize  $\sum_{t=1}^T \|X_t - X_t^*\|$ . And also we need to learn  $R_k(s_t, a_t)$  for so that  $\pi_k$  can maximize  $\sum_{t=1}^T \gamma^t r_{k,t}$ . Traditional cooperative MARL methods such as MAPPO (Yu et al., 2022) or QMIX (Rashid et al., 2020) only require  $\{\pi_k\}_{k=1}^N$  instead of  $R_k(s_t, a_t)$  to be learned. Traditional explicit reward modeling methods such as GAIL (Ho & Ermon, 2016), RLHF (Christiano et al., 2017) only require  $R_k(s_t, a_t)$  instead of  $\{\pi_k\}_{k=1}^N$  to be learned. When the number of agents is 1000+, problem 1 becomes the emergence control problem. Through interactions of agents guided by microscopic policies, different types of macroscopic properties, such as different types of congestion patterns, may emerge. Although from a research perspective,

1000+ agents pose significant challenges for both algorithmic design and hardware implementation. From a practical perspective, it is highly prevalent. For example, regional traffic networks typically involves 1000+ vehicles. Airport terminals usually accommodate 1000+ pedestrians. Problem 2 is actually the counterfactual reasoning based on the explicit reward model obtained by solving problem 1. To distinguish from the problems typically addressed by existing MARL, we collectively refer to the solutions for problem 1 and 2 as emergence control.

The solution to this problem primarily involves two main challenges: (1) To maximize  $\sum_{t=1}^T \|X_t - X_t^*\|$  aligns with the goal of implicit reward modeling (Kontogiannis et al., 2025). To explicitly learn the navigation preference  $R_k(s_t, a_t)$  for each vehicle  $k$  aligns with the goal of explicit reward modeling. However, a simple combination of these two methods cannot effectively solve the integrated problem posed by their superposition. (2) Currently, the primary computational frameworks capable of supporting large-scale multi-agent reinforcement learning (involving 1000+ agents) are population-based MARL and Swarm RL. However, population-based MARL operate on the parameters of policies (Li et al., 2023). It cannot learn an explicit reward model for each individual agent. Swarm RL such as the Mean Field Reinforcement Learning requires the agents in the swarm are interchangeable. Neither of these frameworks can solve the emergence control problem.

We propose the MASDiff training framework that utilizes an evolutionary diffusion model. The evolutionary diffusion model can simultaneously learn the explicit reward model for 1000+ agents and enhance the cooperation of multi-agents to reach the macroscopic target. The contributions of this paper are: (1) Since we don't have the expert dataset for training explicit reward model. We use the diffusion model to generate the expert data. (2) Since the traditional diffusion model can only generate data according to the given data, we introduce evolutionary diffusion model to generate data in an evolutionary manner. This mechanism dose not require any training data and guarantees that we can obtain the explicit reward model  $\{R_k(s_t, a_t)\}_{k=1}^N$ . (3) To steer the collaboration of all agents towards the specified macro-level objective, MASDiff employs the macro-level objective, i.e.  $\text{Argmin} \sum_{t=1}^T \|X_t - X_t^*\|$  in Formula 1, as the fitness for the evolutionary population. It transforms the multi-agent exploration based on intrinsic reward into the evolutionary exploration based on genetic operators. (4) The optimal individual retention strategy in the evolutionary diffusion model can address the non-stationary problem in MARL. (5) MASDiff adopt the DTDE and offline-to-online-to-offline training framework to support MARL with 1000+ agents. Since the evolutionary diffusion model can generate the training data for each agent, the learning process of each agent can be executed in parallel by using offline RL for each agent. Training of instance in population can be executed in parallel in an offline manner. Only sample collection are conducted online.

## 2 RELATED WORK

**Explicit reward modeling.** Current approaches for explicit reward modeling fall into three main categories. (1) **Behavioral Cloning (BC)**. BC method learns intelligent body actions directly (Hussein et al., 2017)(Le Mero et al., 2022)(Zhou et al., 2024). It does not require learning reward functions but needs a large amount of expert trajectory data, split into state-policy pairs to represent the experts actions in the current state. (2) **Inverse Reinforcement Learning (IRL)**. IRL infers a reward function by modeling expert trajectories, and learning a reward function that emulates expert decision patterns (Adams et al., 2022)(Fernando et al., 2021)(Wang et al., 2024). (3) **Reinforcement Learning with Human Feedback (RLHF)**. RLHF aligns agents preferences with real-world behavior by using human feedback to optimize reward functions (Casper et al., 2023)(Kaufmann et al., 2025)(Cao et al., 2023). All three methods require large-scale trajectory datasets, which is problematic in our case due to the lack of granular trajectory data for individual vehicles.

**Implicit Reward modeling.** Implicit reward modeling aims to learn the intrinsic reward of an agent to facilitate better exploration in multi-agent reinforcement learning, ultimately achieving superior overall outcomes. (1) **Single-Agent Exploration**. Reward bonuses based on the inverse state-action count have been shown to be effective in accelerating learning (Strehl & Littman, 2008). In order to scale count-based approaches to large state spaces, state counts are use as reward bonuses (Ostrovski et al., 2017)(Tang et al., 2017). Some work has focused on defining intrinsic rewards for exploration based on inspiration from psychology (Oudeyer & Kaplan, 2007). (2) **Multi-Agent Exploration**. MAVEN (Mahajan et al., 2019) adopts hierarchical control and the policies of the agents are conditioned on the shared latent variable generated by a hierarchical policy. Wang et al. (2020)

propose two exploration methods, EITI and EDTI, to induce cooperative exploration by capturing the influence of one agents on other agents. CMAE (Liu et al., 2021) proposes that reward function only depends on a small subset of the large state space. EMC (Zheng et al., 2021) proposes that local Q function of each agent can capture the novelty of states and the influences between agents.

Both of the explicit and implicit reward model learning methods are unable to solve our problem. The explanation are given in the above paragraphs.

**MARL training architecture for 1000+ agents.** Currently, the primary computational frameworks capable of supporting large-scale multi-agent reinforcement learning are population-based MARL and Swarm MARL. **(1)Population-based MARL.** Population-based MARL is a feasible approach to support distributed learning. Prior works include population-based training (Carroll et al., 2019)(Jaderberg et al., 2019), self-play (Vinyals et al., 2019)(Heinrich et al., 2015) and meta-game (Muller et al., 2020)(Omidshafiei et al., 2019). However, population-based MARL cannot learn an explicit reward model for each individual agent. **(2)Swarm MARL.** MARL for swarm system is described as a swarm MDP environment (Hüttenrauch et al., 2019). Among the MARL methods for swarm systems, mean field theory are wildly adopted. Mean Field Reinforcement Learning is proposed to solve the scalability and the interaction of the population of agents(Guo et al., 2023)(Gu et al., 2023)(Anahtarci et al., 2023)(Perrin et al., 2020). However, those methods require two important properties of swarm systems. (1) the agents in the swarm are interchangeable.(2) the exact number of agents in the swarm is irrelevant. It also cannot learn an explicit reward model for each individual agent.

**Diffusion models for reinforcement learning.** Offline reinforcement learning is a powerful tool for addressing the problem of action sequence combination explosion. However, its performance is ultimately constrained by the offline dataset. Many studies have introduced diffusion models into reinforcement learning(Kang et al., 2023),(Hansen-Estruch et al., 2023),(Ren et al., 2025). Zhu et al. (2023) summarizes the application of diffusion models in reinforcement learning-related fields. Janner et al. (2022) adopts diffusion models as planners to replace traditional autoregressive-based planners. SafeDiffuser (Xiao et al., 2025) introduces safety constraints to generate safe trajectories. Diffusion-QL (Wang et al., 2023) integrates diffusion policies into the Q-learning framework to address the distribution shift issue in offline RL. CEP (Lu et al., 2023) uses contrastive performance prediction to guide energy function sampling. MADiff (Zhu et al., 2024) uses an attention-based diffusion model to model the complex coordination among behaviors of multiple agents. All those methods require training data and thus is not applicable for solvin our problem.

### 3 METHODOLOGY

#### 3.1 IDEAL TRAINING FRAMEWORK VIA ORACLE

We define a Markov Decision Process as  $(S, A, R, \pi)$ , where  $S = \{s_1, s_2, \dots, s_N\}$ ,  $A = \{a_1, a_2, \dots, a_N\}$  and  $R = \{r_1, r_2, \dots, r_N\}$  are the joint set of individual states, actions and rewards, respectively. Once an agent reaches its next state, it receives a reward  $r$ .  $\pi = \{\pi_i\}_{i=1}^N$  are the combination of policies. Let  $\tau_{\pi_i} = \{(s, a, s')_t\}_{t=1}^T$  denote trajectory sampled by  $\pi_i$  where  $s, s' \in S_i$  and  $a \in A_i$ . Let  $\tau = \{\tau_{\pi_i}\}_{i=1}^N$ ,  $\tau \in \mathbb{R}^{|S|*|A|*|S|*|T|*|N|}$  denote the combined trajectory collected by  $\{\pi_i\}_{i=1}^N$ . Let  $R = \{r_i\}_{i=1}^N$ ,  $r_i = \{r_t\}_{t=1}^T$ ,  $r_i$  denotes the corresponding rewards for  $\tau_{\pi_i} \in \tau$ .

Each  $\pi_i \in \pi$  can be obtained by conducting offline reinforcement learning on  $\{\tau_{\pi_i} \cup r_i\} \in \{\tau \cup R\}$  via the DTDE framework. Let  $\{\{X_t\}_{t=1}^T\}_{\{\tau \cup R\}}$  denotes the macroscopic status by conducting simulation on  $\pi$  deriving from  $\{\tau \cup R\}$ . Let  $\{X_t^*\}_{t=1}^T$  denotes the target macroscopic state. If we can obtain the optimal  $\{\tau \cup R\}$ , then the optimal  $\pi$  can be derived. This paper firstly designs an ideal training framework based on an oracle  $\mathcal{O}$  to obtain the optimal  $\{\tau \cup R\}$ .

**Oracle reward model.** We denote by  $\mathcal{O}(\tau)$  an oracle reward model satisfying the following guarantee for  $\tau$ . For input  $\tau$ , the oracle outputs  $R \leftarrow \mathcal{O}(\tau)$ . For any other  $R', R' \neq R$ , we have  $\sum_{t=1}^T \|X_t - X_t^*\| \leq \sum_{t=1}^T \|X'_t - X_t^*\|$ , where  $X_t \in \{\{X_t\}_{t=1}^T\}_{\{\tau \cup R\}}$  and  $X'_t \in \{\{X'_t\}_{t=1}^T\}_{\{\tau \cup R'\}}$ .

**Proposition 1.**  $\sum_{t=1}^T \|X_t - X_t^*\|$  is monotonically decreasing in each iteration of Algorithm 1.

**Algorithm 1** Ideal training framework

- 
- 1: Input: Oracle reward model  $\mathcal{O}$ , target macro emergent phenomenon  $\{X_t\}_{t=1}^T$
  - 2: Initializing policy combination  $\pi = \{\pi_i\}_{i=1}^N$
  - 3: **while**  $\sum_{t=1}^T \|X_t - X_t^*\|$  not converged **do**
  - 4:   Conduct simulation with  $\pi = \{\pi_i\}_{i=1}^N$ , obtaining  $\tau$  and  $\{X_t\}_{t=1}^T$
  - 5:   Query Oracle  $\mathcal{O}$  to obtain  $R$ ,  $R \leftarrow \mathcal{O}(\tau)$
  - 6:   Update  $\pi = \{\pi_i\}_{i=1}^N$  by offline RL via DTDE framework with  $\{\tau \cup R\}$
  - 7: **end while**
- 

**Proof.** Based on oracle reward model  $\mathcal{O}$ , the MARL in Algorithm 1 is actually transformed into a single-agent reinforcement learning.  $\{X_i\}_{i=1}^T$  can be viewed as the ordinary reward value. The entire process is equivalent to classical policy iteration, which has been proven to yield monotonic improvement in returns, corresponding to a monotonic decrease of  $\{X_i\}_{i=1}^T$  in Algorithm 1.

Based on Proposition 1, we introduce the MASDiff framework that use a diffusion model to approximate oracle  $\mathcal{O}$ . MASDiff shift the focus from exploring the optimal  $\pi$  to exploring the optimal  $R$  on condition of  $\tau$ . In the article ‘‘Reward is enough’’ (Silver et al., 2021), the author believes rewards are sufficient to drive behaviors. Therefore, the diffusion model essentially generates the decision preferences for agents.

### 3.2 APPROXIMATING ORACLE $\mathcal{O}$ WITH DIFFUSION MODEL

To approximate oracle  $\mathcal{O}$ , the diffusion model is to generate  $R$  on condition of  $\tau$ . From a practical perspective, it is also necessary to treat  $\tau$  as a condition obtained from the environment. Because the generated  $\tau$  may be inconsistent with reality. This paper adopts the denoising diffusion probability model (Ho et al., 2020) as the basic diffusion model. We firstly give the description of the basic diffusion model. Then, we describe how to train the model such that it can approximate  $\mathcal{O}$ .

#### 3.2.1 CONDITIONAL DENOISING DIFFUSION PROBABILITY MODEL

To generate  $R$  by incorporating  $\tau$  as conditions, the original distribution  $R_0$  of the diffusion process can be transformed into an isotropic Gaussian distribution  $\mathcal{N}(0, 1)$  through  $T$  steps of forward diffusion. The diffusion model denoted as  $p_\theta$  attempts to recover the original distribution of  $R_0$  under the condition of  $\tau$ . Moreover, the dimension of  $\tau$  is  $\mathbb{R}^{|S|*|a|*|S|*|T|*|N|}$  which is quite high when  $N > 1000$  and  $T > 100$ . We reduce the dimensions of  $\tau$  and  $R$  with encoders by mapping  $\tau$  and  $R$  to latent space representation  $Z_\tau$  and  $Z_R$ .

The forward process is parameterized as:

$$q(Z_{R_{1:T}} | Z_{R_0}, Z_\tau) = \prod_{t=1}^T q(Z_{R_t} | Z_{R_{t-1}}, Z_\tau) \quad (2)$$

$$q(Z_{R_t} | Z_{R_{t-1}}, Z_\tau) = \mathcal{N}\left(Z_{R_t}; \sqrt{1 - \beta_t} Z_{R_{t-1}}, \beta_t I\right) \quad (3)$$

where  $q(Z_{R_t} | Z_{R_{t-1}}, Z_\tau)$  represents the conditional Gaussian distribution, and the variance  $\beta_t$  adjusts the noise level.  $Z_{R_t}$  at any time step  $t$  can be directly represented by  $Z_{R_0}$  as:

$$Z_{R_t} = \sqrt{\alpha_t} Z_{R_0} + \sqrt{1 - \alpha_t} \epsilon, \quad \epsilon \sim \mathcal{N}(0, 1) \quad (4)$$

$$\alpha_t = \prod_{s=1}^t (1 - \beta_s) \quad (5)$$

The backward process is parameterized as:

$$p_\theta(Z_{R_0}, \dots, Z_{R_{t-1}} | Z_{R_T}, Z_\tau) = p_\theta(Z_{R_T}) \prod_{t=1}^T p_\theta(Z_{R_{t-1}} | Z_{R_t}, Z_\tau) \quad (6)$$

where  $Z_{R_T} \sim \mathcal{N}(0, 1)$  and  $p_\theta(Z_{R_{t-1}} | Z_{R_t}, Z_\tau)$  is assumed to follow a normal distribution with learnable parameters. The backward process can be trained by optimizing the following objective:

$$L = \min_{\theta} \mathbb{E}_{Z_{R_t}, Z_\tau, \epsilon, t} \left[ \|\epsilon - \epsilon_\theta(Z_{R_t}, Z_\tau, t)\|_2^2 \right] \quad (7)$$

where  $\epsilon_\theta(Z_{R_t}, Z_\tau, t)$  predicts the noise using a neural network.

### 3.2.2 TRAINING OF DIFFUSION MODEL

Training of diffusion model implies two layers of meaning. Firstly, let  $(\{\tau \cup R\}, \rho(\{\tau \cup R\}))$  denote a training sample for  $p_\theta$ .  $\rho(\{\tau \cup R\}) = \sum_{t=1}^T \|X_t - X_t^*\|$  denotes the weight of  $\{\tau \cup R\}$ , where  $X_t \in \{\{X_t\}_{t=1}^T\}_{\{\tau \cup R\}}$ .  $\rho(\{\tau \cup R\})$  is normalized based on Formula 8.  $\min_\rho$  and  $\max_\rho$  denotes the minimum and maximum value of  $\rho$  in the sample set. To obtain high-quality samples,  $\gamma$  is further applied for scaling. If we have enough samples, we can train the model according to Formula 9 (Kumar & Levine, 2020). The weighted training can let the diffusion model generate sample with higher  $\rho$  values.

$$\rho(\{\tau \cup R\}) = \left(1 - \frac{\rho - \min_\rho}{\max_\rho - \min_\rho}\right) * \gamma \quad (8)$$

$$L = \min_\theta \mathbb{E}_{R_t, \tau, \epsilon, t} \left[ \rho(\{\tau \cup R\}) * \|\epsilon - \epsilon_\theta(R_t, \tau, t)\|_2^2 \right] \quad (9)$$

Secondly, since there is not any samples at the beginning, we need to train  $p_\theta$  in a bootstrapping manner. In other words,  $p_\theta$  need to generate training sample for itself. Moreover,  $p_\theta$  need to explore samples that have higher  $\rho$  value. We design the evolution mechanism to let  $p_\theta$  explore better new samples. The evolutionary strategy consists of four steps illustrated as follows.

**Population.** Let  $\mathcal{T} = \{\{\tau \cup R\}_j\}_{j=1}^M$  denote the population of  $\{\tau \cup R\}$ .  $\mathcal{T}$  consists  $M$  instances of  $\{\tau \cup R\}$  where  $R \sim p_\theta(R|\tau)$ . Let  $\Pi = \{\pi_j\}_{j=1}^M$ ,  $\pi = \{\pi_i\}_{i=1}^N$  denote the policy population.  $\Pi$  consists  $M$  instances of  $\pi$ . We randomly initialize the policy population  $\Pi$ . The initial population of  $\mathcal{T}$  is the rollout data by  $\Pi$ .

**Sample scoring.** At each iteration  $k$ , we score each  $\{\tau \cup R\} \in \mathcal{T}^k$  with  $\rho(\{\tau \cup R\})$ .

**Selection.** We use  $\rho(\{\tau \cup R\})$  to decide which samples in  $\mathcal{T}^k$  should be selected to propagate to the next iteration. Let  $E^k$  denotes the elite set selected from  $\mathcal{T}^k$  according to the roulettwheel strategy (Lipowska, 2012).

**Mutation using truncated diffusion-denoising.** We apply randomized mutations to  $E^k$  for exploration. A truncated diffusion-denoising process is adopted to mutate trajectories. We run the first  $t$  steps of the forward diffusion process to add noise to elite trajectories in  $E^k$  based on Formula 10.

$$E_{noisy}^{k+1} = \{\tau \cup R_{noisy}\}, R_{noisy} = \{\sqrt{\alpha_t}R + \sqrt{1 - \alpha_t}\epsilon, \epsilon \sim N(0, 1) | \{\tau \cup R\} \in E^k\} \quad (10)$$

Then we run the last  $t$  steps of the reverse diffusion process to denoise  $E_{noisy}^{k+1}$  to obtain clean  $E_{clean}^{k+1}$  based on Formula 11. Finally,  $E_{clean}^{k+1}$  is the next generation of  $\mathcal{T}^k$ , i.e.  $\mathcal{T}^{k+1}$ .

$$E_{clean}^{k+1} = \{\tau \cup R_{clean}\}, R_{clean} \sim \{p_\theta(R_{noisy}|\tau) | \{\tau \cup R_{noisy}\} \in E_{noisy}^{k+1}\} \quad (11)$$

**Notation.** Both the diffusion model and  $\{X_t\}_{t=1}^T$  incorporate  $T$ .  $T$  denotes the diffusion steps when describing the diffusion model, whereas  $\bar{T}$  represents the time-series length when referring to  $\{X_t\}_{t=1}^T$ . In Formula 10 and 11, we use  $R$  and  $\tau$  to denote  $Z_R$  and  $Z_\tau$  for brevity.

## 3.3 TRAINING FRAMEWORK OF MASDIFF TO SUPPORT MARL WITH 1000+ AGENTS

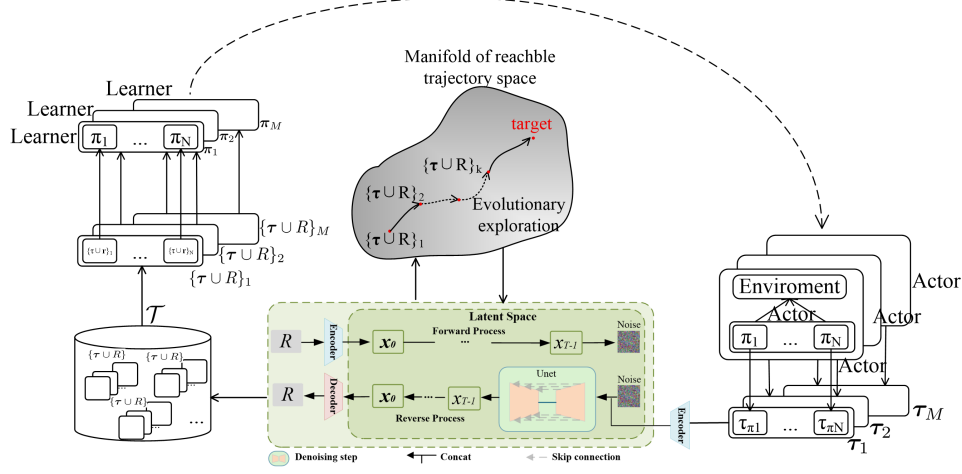


Figure 2: Training framework of MASDiff

**Training framework of MASDiff.** Algorithm 2 illustrates how to transform the ideal training framework into the practically feasible training framework MASDiff. MASDiff utilizes policy population and sample population to leverage evolutionary search over the trajectory manifold.

**Algorithm 2** Training framework of MASDiff

- 1: Input: Target macro emergent phenomenon  $\{X_t\}_{t=1}^T$ , number of agents  $N$
- 2: Initializing policy population  $\Pi, \Pi = \{\pi_j\}_{j=1}^M, \pi = \{\pi_i\}_{i=1}^N$
- 3: Initializing diffusion model  $p_\theta$
- 4: **while**  $\sum_{t=1}^T \|X_t - X_t^*\|$  not converged or max iteration not reached **do**
- 5: Simulating with  $\Pi = \{\pi_j\}_{j=1}^M$ , obtaining  $\{\tau_j\}_{j=1}^M$  for each  $\pi_j$   $\triangleright$  collect trajectories online
- 6: Query  $p_\theta$  to obtain  $R_j$  for each  $\tau_j \in \{\tau_j\}_{j=1}^M$ ,  $R_j \sim p_\theta(R_j|\tau_j)$
- 7: Update each  $\pi_j \in \Pi$  by independent offline RL with  $\{\tau_j \cup R_j\}$   $\triangleright$  update  $\pi_j \in \Pi$  offline
- 8: Simulating with  $\Pi = \{\pi_j\}_{j=1}^M$ , obtaining  $\{\{X_t\}_{t=1}^T\}_{j=1}^M$  for each  $\pi_j$   $\triangleright$  collect simulation result for trajectories in line 5 online
- 9: Update  $p_\theta$  with weighted sample  $\{\{\tau_j \cup R_j\}, \rho(\{\tau_j \cup R_j\})\}_{j=1}^M$
- 10: Evolving the sample population  $\{\tau_j \cup R_j\}_{j=1}^M$  to the next generation  $\{\tau_j \cup R'_j\}_{j=1}^M$
- 11: Update each  $\pi_j \in \Pi$  by independent offline RL with  $\{\tau_j \cup R'_j\}$   $\triangleright$  update  $\pi_j \in \Pi$  offline
- 12: Simulating with  $\Pi = \{\pi_j\}_{j=1}^M$ , obtaining  $\{\{X_t\}_{t=1}^T\}_{j=1}^M$  for each  $\pi_j$   $\triangleright$  collect simulation result for trajectories in line 10 online
- 13: Update  $p_\theta$  with weighted sample  $\{\{\tau_j \cup R'_j\}, \rho(\{\tau_j \cup R'_j\})\}_{j=1}^M$
- 14: **end while**

In line 5, we use the policy population to collect trajectories. In line 6, rewards are generated by query the diffusion model with the trajectories. In line 7~9, we update the diffusion model by samples collected in line 5~6. In line 10~13, the diffusion model are further updated by samples collected with evolutionary strategy. MASDiff also adopt the offline-to-online-to-offline training framework. Training are all conducted offline. Only sample collection and simulation result collection are conducted online.

**Recovering explicit reward model  $R_k(s_t, a_t)$  from weighted sample  $\{\{\tau_j \cup R'_j\}, \rho(\{\tau_j \cup R'_j\})\}_{j=1}^M$ .** Assume  $\{\{\tau_j \cup R'_j\}, \rho(\{\tau_j \cup R'_j\})\}$  has the highest  $\rho$  value. Assume  $\{\tau_{\pi_i} \cup r_k\} \in \{\{\tau_j \cup R'_j\}, \tau_{\pi_k} = \{(s, a, s')_t\}_{t=1}^T$  and  $r_k = \{r_t\}_{t=1}^T$ . We can parameterize  $R_k(s_t, a_t)$  using a single-time-step neural network  $r_t = g_\theta(s_t, a_t)$  to predict a reward for a state-action pair  $(s_t, a_t)$ .

**Parallel training paradigm for MASDiff.** Although the dimension of the trajectory space is  $\mathbb{R}^{|S|*|a|*|S|*|T|*|N|}$  which is quite large. The reachable trajectory space is a low dimensional manifold embedded in  $\mathbb{R}^{|S|*|a|*|S|*|T|*|N|}$ . MASDiff is essentially to efficiently explore the reachable trajectory manifold that can emerge the target macroscopic phenomenon. Figure 2 illustrates the training framework of MASDiff. It decouples the training and rollout tasks using the Actor-Diffusion-Learner model. The Learner operates policy training and the Actor operates data collecting. The diffusion model interleaves the Learners and Actors by providing better samples to the Learner instead of those collected directly by Actors.

To conduct fully decentralized multi-agent learning, MASDiff uses  $\tau_{\pi_i} \in \tau$  to train  $\pi_i \in \pi$  independently. One learner for one policy, policy learning is independent to other agents. Actor is to collect the rollout data  $\{\tau_{\pi_i}\}_{i=1}^N$  using  $\{\pi_i\}_{i=1}^N$ . Multiple Actors can be dispatched in parallel. The Actor-Diffusion-Learner architecture are built on top of Ray which allows tasks to be distributed over a large cluster.

## 4 EXPERIMENT

The following three research questions illustrate the purpose of the experiment:

**RQ1:** Can MASDiff minimize  $\sum_{t=1}^T \|X_t - X_t^*\|$ ?

**RQ2:** Can the explicit reward model  $R_k(s_t, a_t)$  learned by MASDiff answer the ‘What if’ counterfactual questions in new scenarios?

**RQ3:** Since the search space of multi-agent reinforcement learning grows exponentially as the number of agents increases, how is the scalability of the MASDiff framework?

### 4.1 EXPERIMENTAL SETTINGS

**Comparison methods.** Since existing CTDE or DTDE MARL framework can not work on such large number of vehicles with local reward completely unknown. We design methods for comparison by change or remove part of the algorithms in MASDiff framework.

CGAN replaces the conditional denoising diffusion probabilistic model in MASDiff with Conditional Generative Adversarial Network.

LDM.IMIT adopts imitation learning instead of offline DQN in learning vehicle navigation policy.

LDM.CFG use classifier free guidance in the conditional diffusion model. It incorporates  $\rho(\{\tau \cup R\})$  as an additional condition for the diffusion model, i.e.  $R \sim p_\theta(R|\tau_j, \rho(\{\tau \cup R\}))$ . Instead of utilizing evolutionary strategy to explore new sample, LDM.CFG uses the additional condition  $\rho(\{\tau \cup R\})$  to explore new sample.

LDM.PCFG add a proxy in LDM.CFG. The proxy uses  $\sum_{t=1}^T \|X_t - X_t^*\|$  to predict  $\rho(\{\tau \cup R\})$ . The condition  $\rho(\{\tau \cup R\})$  in LDM.PCFG is given by the proxy in stead of the real value of  $\rho(\{\tau \cup R\})$ .

**Experimental environment.** Experiments were conducted on a high-performance server equipped with an AMD EPYC 9554 CPU, four NVIDIA RTX 4090 GPUs (24GB each), 256GB of DDR5 RAM, and a 2TB NVMe SSD. The operation system is Ubuntu 22.04 LTS. The software stack include Python 3.10, the SUMO(Simulation of Urban MObility) simulation environment, and the Ray framework for parallel execution. The population size in Algorithm 2 is 500. The traffic network in 3 which is part of a certain district of Chongqing. This region includes 114 roads and 45 intersections.

**4.2 RQ1: CAN MASDIFF MINIMIZE  $\sum_{t=1}^T \|X_t - X_t^*\|$ ?**

**Experimental design.** This experiment is conducted in scenario 1 illustrated in Figure 3(b). Figure 3(a) is the original Baidu map. Figure 3(b) is the SUMO traffic road network derived from Figure 3(a). The surveillance camera data from 8:30 to 9:00 am on August 21, 2023 in this region are used for this experiment. There are totally 1557 vehicles in the camera data. We adopt the largest functional cluster of the traffic network (*LCC*) as the emergent macro traffic state denoted as  $\{LCC_t^*\}_{t=1}^T$ . Functional cluster is the connected roads with relatively high speed (Zeng et al.,

2020). We calculate  $LCC$  using the velocity of the 1557 vehicles recorded by surveillance camera.  $LCC$  is calculated every minute and thus  $T = 30$  in  $\{LCC_t^*\}_{t=1}^T$ . MASDiff learns the navigation policy for each vehicle according to  $\{LCC_t^*\}_{t=1}^{30}$ .  $\{LCC_t\}_{t=1}^{30}$  is obtained by conduct simulation with the learned policies.

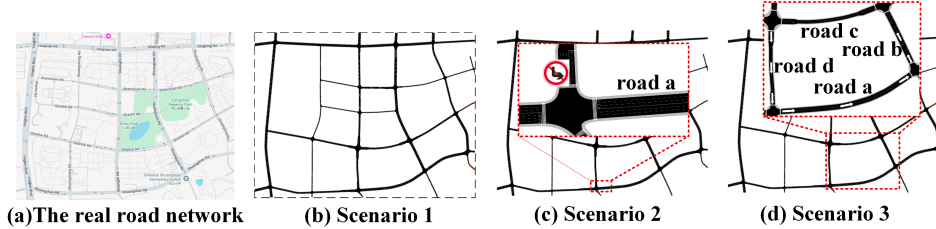


Figure 3: Different experiment scenarios

**Experimental results.** The results are shown in Figure 4. The horizontal axis represents the iteration times. The vertical axis represents the minimal value of  $\|\{LCC_t^*\}_{t=1}^{30} - \{LCC_t\}_{t=1}^{30}\|$  during the past iterations. All methods are iterated 100 rounds. MASDiff obtain the best results. LDM\_CFG and LDM\_PCFG produce the worst results. This indicates that directly using  $\sum_{t=1}^T \|X_t - X_t^*\|$  as a condition provides limited guidance for the diffusion model’s exploration.

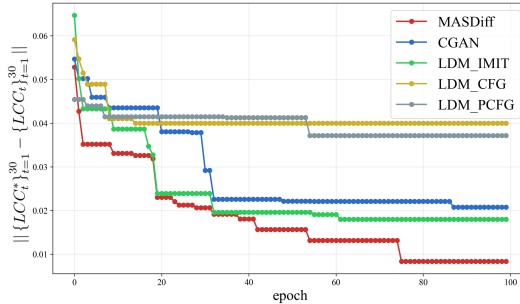


Figure 4: Minimal value of  $\|\{LCC_t^*\}_{t=1}^{30} - \{LCC_t\}_{t=1}^{30}\|$  per iteration.

Figure 5 shows the detail of  $\{LCC_t^*\}_{t=1}^{30}$  in each minute. The orange curve is  $\{LCC_t^*\}_{t=1}^{30}$ . The blue curve is  $\{LCC_t\}_{t=1}^{30}$ . The curve by MASDiff best approximates the ground truth.

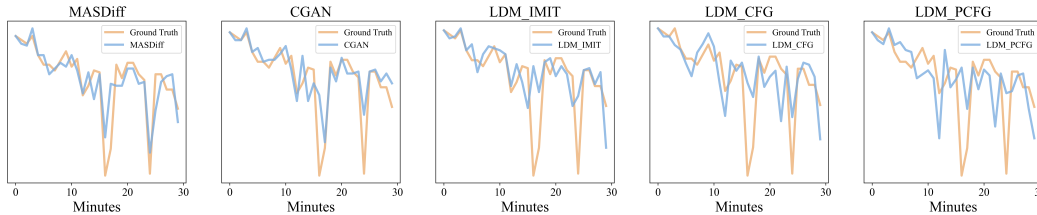


Figure 5:  $\{LCC_t^*\}_{t=1}^{30}$  and  $\{LCC_t\}_{t=1}^{30}$  in scenario 1.

### 4.3 RQ2: CAN WE CONDUCT COUNTERFACTUAL REASONING BY THE LEARNED EXPLICIT REWARD MODEL?

**Experimental design.** This experiment is conducted in scenario 2~3 illustrated in Figure 3(c)~(d).

Scenario 2: What is the  $\{MQL_t\}_{t=1}^T$  if we prohibiting the left turns onto road a at the intersections?

Scenario 3: What is the  $\{MQL_t\}_{t=1}^T$  if we establish a one-way traffic loop in the specific area?

$MQ_L$  is adopted as the macro emergent traffic state which is the maximum queue length of a road segment in a certain period.  $\{MQ_L^*\}_{t=1}^{30}$  represents the ground truth data.  $\{MQ_L^*\}_{t=1,s1}^{30} \sim \{MQ_L\}_{t=1,s3}^{30}$  denote the ground truth data in scenario 1  $\sim$  3.  $\{MQ_L\}_{t=1,s1}^{30} \sim \{MQ_L^*\}_{t=1,s3}^{30}$  denote the results obtained through different methods in scenario 1  $\sim$  3.

Due to the change of scenario,  $\{MQ_L^*\}_{t=1,s2}^{30}$  and  $\{MQ_L^*\}_{t=1,s3}^{30}$  cannot be obtained from the original surveillance camera data or be predicted based on historical data. Therefore, we synthesize the ground truth data through randomly simulated reward models. First, we randomly generate the ground truth reward model for those 1557 vehicles. Based on the generated reward models, we conduct distributed RL for each vehicle in scenario 1  $\sim$  3 to obtain the navigation policies. By simulating with the policies we can obtain  $\{MQ_L^*\}_{t=1,s2}^{30}, \{MQ_L^*\}_{t=1,s3}^{30}$ . To obtain  $\{MQ_L\}_{t=1,s2}^{30}, \{MQ_L\}_{t=1,s3}^{30}$ , We first learn the explicit reward model and the navigation policy according to  $\{MQ_L^*\}_{t=1,s1}^{30}$ . Then,  $\{MQ_L\}_{t=1,s2}^{30}$  and  $\{MQ_L\}_{t=1,s3}^{30}$  can be obtained by conduct simulation with the learned policies in scenario 2 and 3.

**Experimental results.** In scenario 1, roads a and b are congested. Roads c and d experience lower traffic volumes. In scenario 2, congestion on road a is partially alleviated. In scenario 3, congestion on both roads a and b is reduced. Figures 7, 8 and 9 show heatmaps of  $\{MQ_L\}_{t=1,s1}^{30} \sim \{MQ_L\}_{t=1,s2}^{30}$  for each method in different scenarios. The horizontal axis represents the minutes. The vertical axis represents the  $MQ_L$  value of a specific road in each minute. The top 50 roads with the highest  $MQ_L$  in the original scenario were selected. Due to space constraints, we place Figures 7, 8 and 9 in the appendix. In different scenarios,  $\{MQ_L\}_{t=1,s1}^{30} \sim \{MQ_L\}_{t=1,s4}^{30}$  by MASDiff are more consistent with  $\{MQ_L^*\}_{t=1,s1}^{30} \sim \{MQ_L^*\}_{t=1,s3}^{30}$  than other methods. This means that the learned policies can be used to conduct counterfactual reasoning in hypothetical scenarios.

#### 4.4 RQ3: HOW IS THE SCALABILITY OF THE MASDIFF FRAMEWORK

Scenario 1 is chose as the simulation environment. We increase the number of vehicles from 1,557 to 2,578 and 4,423 at time intervals of 0.5, 1, and 1.5 hours from 8:30 am. We adopt  $MQ_L$  as the macro metric. The population size is also 500. The results are shown in Tables 1 and 2. In Table 2, running time denotes the time elapsed when the algorithm reaches convergence. The experimental results demonstrates that MASDiff significantly outperforms other methods both in performance and running time. The results obtained by LDM\_CFG and LDM\_PCFG are nearly identical to that from random sampling.

Table 1:  $\|\{MQ_L^*\}_{t=1}^T - \{MQ_L\}_{t=1}^T\|$

	1557	2758	4423
MASDiff	1.69	2.19	2.46
CGAN	1.76	5.32	7.24
LDM.IMIT	2.02	2.40	2.91
LDM_CFG	-	-	-
LDM_PCFG	-	-	-

Table 2: Running time(minute)

	1557	2758	4423
MASDiff	101.25	139.17	218.33
CGAN	172.08	245.83	259.17
LDM.IMIT	113.75	123.33	140.83
LDM_CFG	-	-	-
LDM_PCFG	-	-	-

## 5 CONCLUSION

This paper tries to control the emergence of large scale multi-agent system. To solve this problem, the MASDiff framework are proposed. The experimental results verifies that MASDiff can achieve emergence control fairly well. To our best knowledge, MASDiff is the first framework targeting at controlling emergence in MAS. In future work, each part of the framework will be optimized.

## REFERENCES

- 540  
541  
542 Stephen C. Adams, Tyler Cody, and Peter A. Beling. A survey of inverse reinforcement learning.  
543 *Artif. Intell. Rev.*, 55(6):4307–4346, 2022. doi: 10.1007/S10462-021-10108-X. URL <https://doi.org/10.1007/s10462-021-10108-x>.  
544
- 545 Berkay Anahtarci, Can Deha Kariksiz, and Naci Saldi. Q-learning in regularized mean-field games.  
546 *Dyn. Games Appl.*, 13(1):89–117, 2023. doi: 10.1007/S13235-022-00450-2. URL <https://doi.org/10.1007/s13235-022-00450-2>.  
547  
548
- 549 Yulong Cao, Boris Ivanovic, Chaowei Xiao, and Marco Pavone. Reinforcement learning with human  
550 feedback for realistic traffic simulation. *CoRR*, abs/2309.00709, 2023. doi: 10.48550/ARXIV.  
551 2309.00709. URL <https://doi.org/10.48550/arXiv.2309.00709>.
- 552 Micah Carroll, Rohin Shah, Mark K. Ho, Thomas L. Griffiths, Sanjit A. Seshia, P. Abbeel, and  
553 Anca D. Dragan. On the utility of learning about humans for human-ai coordination. *ArXiv*,  
554 abs/1910.05789, 2019. URL [https://api.semanticscholar.org/CorpusID:  
555 202770731](https://api.semanticscholar.org/CorpusID:202770731).  
556
- 557 Stephen Casper, Xander Davies, Claudia Shi, Thomas Krendl Gilbert, Jérémy Scheurer, Javier Rando,  
558 Rachel Freedman, Tomasz Korbak, David Lindner, Pedro Freire, Tony Tong Wang, Samuel  
559 Marks, Charbel-Raphaël Ségerie, Micah Carroll, Andi Peng, Phillip J. K. Christoffersen, Mehul  
560 Damani, Stewart Slocum, Usman Anwar, Anand Siththaranjan, Max Nadeau, Eric J. Michaud, Ja-  
561 cob Pfau, Dmitrii Krashennikov, Xin Chen, Lauro Langosco, Peter Hase, Erdem Biyik, Anca D.  
562 Dragan, David Krueger, Dorsa Sadigh, and Dylan Hadfield-Menell. Open problems and funda-  
563 mental limitations of reinforcement learning from human feedback. *CoRR*, abs/2307.15217, 2023.  
564 doi: 10.48550/ARXIV.2307.15217. URL [https://doi.org/10.48550/arXiv.2307.  
565 15217](https://doi.org/10.48550/arXiv.2307.15217).
- 566 Paul F Christiano, Jan Leike, Tom Brown, Miljan Martic, Shane Legg, and Dario Amodei. Deep  
567 reinforcement learning from human preferences. *Advances in neural information processing sys-  
568 tems*, 30, 2017.
- 569 Tharindu Fernando, Simon Denman, Sridha Sridharan, and Clinton Fookes. Deep inverse rein-  
570 forcement learning for behavior prediction in autonomous driving: Accurate forecasts of vehicle  
571 motion. *IEEE Signal Process. Mag.*, 38(1):87–96, 2021. doi: 10.1109/MSP.2020.2988287. URL  
572 <https://doi.org/10.1109/MSP.2020.2988287>.  
573
- 574 Haotian Gu, Xin Guo, Xiaoli Wei, and Renyuan Xu. Dynamic programming principles for mean-  
575 field controls with learning. *Oper. Res.*, 71(4):1040–1054, 2023. doi: 10.1287/OPRE.2022.2395.  
576 URL <https://doi.org/10.1287/opre.2022.2395>.  
577
- 578 Xin Guo, Anran Hu, Renyuan Xu, and Junzi Zhang. A general framework for learning mean-  
579 field games. *Math. Oper. Res.*, 48(2):656–686, 2023. doi: 10.1287/MOOR.2022.1274. URL  
580 <https://doi.org/10.1287/moor.2022.1274>.
- 581 Philippe Hansen-Estruch, Ilya Kostrikov, Michael Janner, Jakub Grudzien Kuba, and Sergey  
582 Levine. IDQL: implicit q-learning as an actor-critic method with diffusion policies. *CoRR*,  
583 abs/2304.10573, 2023. doi: 10.48550/ARXIV.2304.10573. URL [https://doi.org/10.  
584 48550/arXiv.2304.10573](https://doi.org/10.48550/arXiv.2304.10573).
- 585 Johannes Heinrich, Marc Lanctot, and David Silver. Fictitious self-play in extensive-form games.  
586 In *Proceedings of the 32nd International Conference on International Conference on Machine  
587 Learning - Volume 37, ICML’15*, pp. 805813. JMLR.org, 2015.  
588
- 589 Jonathan Ho and Stefano Ermon. Generative adversarial imitation learning. *Advances in neural  
590 information processing systems*, 29, 2016.  
591
- 592 Jonathan Ho, Ajay Jain, and Pieter Abbeel. Denoising diffusion probabilistic models. In *Proceed-  
593 ings of the 34th International Conference on Neural Information Processing Systems, NIPS ’20*,  
Red Hook, NY, USA, 2020. Curran Associates Inc. ISBN 9781713829546.

- 594 Ahmed Hussein, Mohamed Medhat Gaber, Eyad Elyan, and Chrisina Jayne. Imitation learning: A  
595 survey of learning methods. *ACM Comput. Surv.*, 50(2):21:1–21:35, 2017. doi: 10.1145/3054912.  
596 URL <https://doi.org/10.1145/3054912>.  
597
- 598 Maximilian Hüttenrauch, Adrian Sosic, and Gerhard Neumann. Deep reinforcement learning for  
599 swarm systems. *J. Mach. Learn. Res.*, 20:54:1–54:31, 2019. URL [https://jmlr.org/  
600 papers/v20/18-476.html](https://jmlr.org/papers/v20/18-476.html).
- 601 Max Jaderberg, Wojciech M. Czarnecki, Iain Dunning, Luke Marris, Guy Lever, Antonio Garci-  
602 a Castaeda, Charles Beattie, Neil C. Rabinowitz, Ari S. Morcos, Avraham Ruderman, Nico-  
603 las Sonnerat, Tim Green, Louise Deason, Joel Z. Leibo, David Silver, Demis Hassabis, Ko-  
604 ray Kavukcuoglu, and Thore Graepel. Human-level performance in 3d multiplayer games with  
605 population-based reinforcement learning. *Science*, 364(6443):859–865, 2019. doi: 10.1126/  
606 science.aau6249. URL [https://www.science.org/doi/abs/10.1126/science.  
607 aau6249](https://www.science.org/doi/abs/10.1126/science.aau6249).
- 608 Michael Janner, Yilun Du, Joshua B. Tenenbaum, and Sergey Levine. Planning with diffusion for  
609 flexible behavior synthesis. In Kamalika Chaudhuri, Stefanie Jegelka, Le Song, Csaba Szepesvári,  
610 Gang Niu, and Sivan Sabato (eds.), *International Conference on Machine Learning, ICML 2022,  
611 17-23 July 2022, Baltimore, Maryland, USA*, volume 162 of *Proceedings of Machine Learn-  
612 ing Research*, pp. 9902–9915. PMLR, 2022. URL [https://proceedings.mlr.press/  
613 v162/janner22a.html](https://proceedings.mlr.press/v162/janner22a.html).
- 614 Zheng Jiaoling, Xu Weifeng, Luo Qian, Dang Wanli, Geng Long, Gao Guokang, Ren Yulin,  
615 and Fan Xingyu. Decision preference alignment for large-scale agents based on reward mod-  
616 el generation. In *ICLR 2025 Workshop on Bidirectional Human-AI Alignment*, 2025. URL  
617 <https://openreview.net/forum?id=mQ1pLtdjbg>.  
618
- 619 Bingyi Kang, Xiao Ma, Chao Du, Tianyu Pang, and Shuicheng Yan. Efficient diffu-  
620 sion policies for offline reinforcement learning. In Alice Oh, Tristan Naumann, Amir  
621 Globerson, Kate Saenko, Moritz Hardt, and Sergey Levine (eds.), *Advances in Neu-  
622 ral Information Processing Systems 36: Annual Conference on Neural Information Pro-  
623 cessing Systems 2023, NeurIPS 2023, New Orleans, LA, USA, December 10 - 16,  
624 2023*, 2023. URL [http://papers.nips.cc/paper\\_files/paper/2023/hash/  
625 d45e0bfb5a39477d56b55c0824200008-Abstract-Conference.html](http://papers.nips.cc/paper_files/paper/2023/hash/d45e0bfb5a39477d56b55c0824200008-Abstract-Conference.html).
- 626 Timo Kaufmann, Paul Weng, Viktor Bengs, and Eyke Hüllermeier. A survey of reinforcemen-  
627 t learning from human feedback. *Trans. Mach. Learn. Res.*, 2025, 2025. URL [https:  
628 //openreview.net/forum?id=f70kIurx4b](https://openreview.net/forum?id=f70kIurx4b).
- 629 Andreas Kontogiannis, Konstantinos Papathanasiou, Yi Shen, Giorgos Stamou, Michael M. Za-  
630 vlanos, and George A. Vouros. Enhancing cooperative multi-agent reinforcement learning with s-  
631 tate modelling and adversarial exploration. *CoRR*, abs/2505.05262, 2025. doi: 10.48550/ARXIV.  
632 2505.05262. URL <https://doi.org/10.48550/arXiv.2505.05262>.  
633
- 634 Ilya Kostrikov, Ashvin Nair, and Sergey Levine. Offline reinforcement learning with implicit q-  
635 learning. In *The Tenth International Conference on Learning Representations, ICLR 2022, Vir-  
636 tual Event, April 25-29, 2022*. OpenReview.net, 2022. URL [https://openreview.net/  
637 forum?id=68n2s9ZJWF8](https://openreview.net/forum?id=68n2s9ZJWF8).
- 638 Aviral Kumar and Sergey Levine. Model inversion networks for model-based optimization.  
639 In H. Larochelle, M. Ranzato, R. Hadsell, M.F. Balcan, and H. Lin (eds.), *Advances in  
640 Neural Information Processing Systems*, volume 33, pp. 5126–5137. Curran Associates, Inc., 2020. URL [https://proceedings.neurips.cc/paper\\_files/paper/2020/  
641 file/373e4c5d8edfa8b74fd4b6791d0cf6dc-Paper.pdf](https://proceedings.neurips.cc/paper_files/paper/2020/file/373e4c5d8edfa8b74fd4b6791d0cf6dc-Paper.pdf).  
642
- 643 Luc Le Mero, Dewei Yi, Mehrdad Dianati, and Alexandros Mouzakitis. A survey on imitation  
644 learning techniques for end-to-end autonomous vehicles. *IEEE Transactions on Intelligent Trans-  
645 portation Systems*, 23(9):14128–14147, 2022. doi: 10.1109/TITS.2022.3144867.  
646
- 647 Pengyi Li, Jianye Hao, Hongyao Tang, Yan Zheng, and Xian Fu. RACE: improve multi-agent  
reinforcement learning with representation asymmetry and collaborative evolution. In Andreas

- 648 Krause, Emma Brunskill, Kyunghyun Cho, Barbara Engelhardt, Sivan Sabato, and Jonathan Scarlett (eds.), *International Conference on Machine Learning, ICML 2023, 23-29 July 2023, Honolulu, Hawaii, USA*, volume 202 of *Proceedings of Machine Learning Research*, pp. 19490–19503. PMLR, 2023. URL <https://proceedings.mlr.press/v202/li23i.html>.
- 652
- 653 Lipowski Dorota Lipowska. Roulette-wheel selection via stochastic acceptance. *Physica A: Statistical Mechanics and its Applications*, 2012.
- 654
- 655 Iou-Jen Liu, Unnat Jain, Raymond A. Yeh, and Alexander G. Schwing. Cooperative exploration for multi-agent deep reinforcement learning. In Marina Meila and Tong Zhang (eds.), *Proceedings of the 38th International Conference on Machine Learning, ICML 2021, 18-24 July 2021, Virtual Event*, volume 139 of *Proceedings of Machine Learning Research*, pp. 6826–6836. PMLR, 2021. URL <http://proceedings.mlr.press/v139/liu21j.html>.
- 660
- 661 Cheng Lu, Huayu Chen, Jianfei Chen, Hang Su, Chongxuan Li, and Jun Zhu. Contrastive energy prediction for exact energy-guided diffusion sampling in offline reinforcement learning. In Andreas Krause, Emma Brunskill, Kyunghyun Cho, Barbara Engelhardt, Sivan Sabato, and Jonathan Scarlett (eds.), *International Conference on Machine Learning, ICML 2023, 23-29 July 2023, Honolulu, Hawaii, USA*, volume 202 of *Proceedings of Machine Learning Research*, pp. 22825–22855. PMLR, 2023. URL <https://proceedings.mlr.press/v202/lu23d.html>.
- 666
- 667 Anuj Mahajan, Tabish Rashid, Mikayel Samvelyan, and Shimon Whiteson. MAVEN: multi-agent variational exploration. In Hanna M. Wallach, Hugo Larochelle, Alina Beygelzimer, Florence d’Alché-Buc, Emily B. Fox, and Roman Garnett (eds.), *Advances in Neural Information Processing Systems 32: Annual Conference on Neural Information Processing Systems 2019, NeurIPS 2019, December 8-14, 2019, Vancouver, BC, Canada*, pp. 7611–7622, 2019. URL <https://proceedings.neurips.cc/paper/2019/hash/f816dc0acface7498e10496222e9db10-Abstract.html>.
- 673
- 674 Paul Muller, Shayegan Omidshafiei, Mark Rowland, Karl Tuyls, Julien Pérolat, Siqi Liu, Daniel Hennes, Luke Marris, Marc Lanctot, Edward Hughes, Zhe Wang, Guy Lever, Nicolas Heess, Thore Graepel, and Rémi Munos. A generalized training approach for multiagent learning. In *8th International Conference on Learning Representations, ICLR 2020, Addis Ababa, Ethiopia, April 26-30, 2020*. OpenReview.net, 2020. URL <https://openreview.net/forum?id=Bk15kxrKDr>.
- 680
- 681 Shayegan Omidshafiei, Christos Papadimitriou, Georgios Piliouras, Karl Tuyls, Mark Rowland, Jean-Baptiste Lespiau, Wojciech M Czarnecki, Marc Lanctot, Julien Perolat, and Remi Munos.  $\alpha$ -rank: Multi-agent evaluation by evolution. *Scientific reports*, 9(1):9937, July 2019. ISSN 2045-2322. doi: 10.1038/s41598-019-45619-9. URL <https://europepmc.org/articles/PMC6617105>.
- 685
- 686 Georg Ostrovski, Marc G. Bellemare, Aäron van den Oord, and Rémi Munos. Count-based exploration with neural density models. *ArXiv*, abs/1703.01310, 2017. URL <https://api.semanticscholar.org/CorpusID:2924063>.
- 689
- 690 Pierre-Yves Oudeyer and Frédéric Kaplan. What is intrinsic motivation? A typology of computational approaches. *Frontiers Neurobotics*, 1:6, 2007. doi: 10.3389/NEURO.12.006.2007. URL <https://doi.org/10.3389/neuro.12.006.2007>.
- 692
- 693 Matej Pecháč, Michal Chovanec, and Igor Farkaš. Self-supervised network distillation: An effective approach to exploration in sparse reward environments. *Neurocomputing*, 599:128033, 2024.
- 695
- 696 Sarah Perrin, Julien Pérolat, Mathieu Laurière, Matthieu Geist, Romuald Elie, and Olivier Pietquin. Fictitious play for mean field games: Continuous time analysis and applications. In Hugo Larochelle, Marc Aurelio Ranzato, Raia Hadsell, Maria-Florina Balcan, and Hsuan-Tien Lin (eds.), *Advances in Neural Information Processing Systems 33: Annual Conference on Neural Information Processing Systems 2020, NeurIPS 2020, December 6-12, 2020, virtual*, 2020. URL <https://proceedings.neurips.cc/paper/2020/hash/995ca733e3657ff9f5f3c823d73371e1-Abstract.html>.
- 701

- 702 Tabish Rashid, Mikayel Samvelyan, Christian Schroeder De Witt, Gregory Farquhar, Jakob Foerster,  
703 and Shimon Whiteson. Monotonic value function factorisation for deep multi-agent reinforcement  
704 learning. *Journal of Machine Learning Research*, 21(178):1–51, 2020.  
705
- 706 Allen Z. Ren, Justin Lidard, Lars Lien Ankile, Anthony Simeonov, Pulkit Agrawal, Anirudha Ma-  
707 jumdar, Benjamin Burchfiel, Hongkai Dai, and Max Simchowitz. Diffusion policy policy opti-  
708 mization. In *The Thirteenth International Conference on Learning Representations, ICLR 2025,*  
709 *Singapore, April 24-28, 2025*. OpenReview.net, 2025. URL [https://openreview.net/](https://openreview.net/forum?id=mEpqHvbD2h)  
710 [forum?id=mEpqHvbD2h](https://openreview.net/forum?id=mEpqHvbD2h).
- 711 David Silver, Satinder Singh, Doina Precup, and Richard S. Sutton. Reward is enough. *Artif.*  
712 *Intell.*, 299:103535, 2021. doi: 10.1016/J.ARTINT.2021.103535. URL [https://doi.org/](https://doi.org/10.1016/j.artint.2021.103535)  
713 [10.1016/j.artint.2021.103535](https://doi.org/10.1016/j.artint.2021.103535).  
714
- 715 Alexander L. Strehl and Michael L. Littman. An analysis of model-based interval estimation for  
716 markov decision processes. *J. Comput. Syst. Sci.*, 74(8):1309–1331, 2008. doi: 10.1016/J.JCSS.  
717 2007.08.009. URL <https://doi.org/10.1016/j.jcss.2007.08.009>.
- 718 Haoran Tang, Rein Houthoofd, Davis Foote, Adam Stooke, Xi Chen, Yan Duan, John Schul-  
719 man, Filip De Turck, and Pieter Abbeel. #exploration: A study of count-based explo-  
720 ration for deep reinforcement learning. In Isabelle Guyon, Ulrike von Luxburg, Samy  
721 Bengio, Hanna M. Wallach, Rob Fergus, S. V. N. Vishwanathan, and Roman Garnett (ed-  
722 s.), *Advances in Neural Information Processing Systems 30: Annual Conference on Neu-  
723 ral Information Processing Systems 2017, December 4-9, 2017, Long Beach, CA, USA*, pp.  
724 2753–2762, 2017. URL [https://proceedings.neurips.cc/paper/2017/hash/](https://proceedings.neurips.cc/paper/2017/hash/3a20f62a0af1aa152670bab3c602feed-Abstract.html)  
725 [3a20f62a0af1aa152670bab3c602feed-Abstract.html](https://proceedings.neurips.cc/paper/2017/hash/3a20f62a0af1aa152670bab3c602feed-Abstract.html).
- 726
- 727 Oriol Vinyals, Igor Babuschkin, Wojciech M. Czarnecki, Michaël Mathieu, Andrew Joseph Dudzik,  
728 Junyoung Chung, David Choi, Richard Powell, Timo Ewalds, Petko Georgiev, Junhyuk Oh, Dan  
729 Horgan, Manuel Kroiss, Ivo Danihelka, Aja Huang, L. Sifre, Trevor Cai, John P. Agapiou, Max  
730 Jaderberg, Alexander Sasha Vezhnevets, Rémi Leblond, Tobias Pohlen, Valentin Dalibard, David  
731 Budden, Yury Sulsky, James Molloy, Tom Le Paine, Caglar Gulcehre, Ziyun Wang, Tobias Pfaff,  
732 Yuhuai Wu, Roman Ring, Dani Yogatama, Dario Wünsch, Katrina McKinney, Oliver Smith, Tom  
733 Schaul, Timothy P. Lillicrap, Koray Kavukcuoglu, Demis Hassabis, Chris Apps, and David Silver.  
734 Grandmaster level in starcraft ii using multi-agent reinforcement learning. *Nature*, 575:350 – 354,  
735 2019. URL <https://api.semanticscholar.org/CorpusID:204972004>.
- 736 Tonghan Wang, Jianhao Wang, Yi Wu, and Chongjie Zhang. Influence-based multi-agent explo-  
737 ration. In *8th International Conference on Learning Representations, ICLR 2020, Addis Ababa,*  
738 *Ethiopia, April 26-30, 2020*. OpenReview.net, 2020. URL [https://openreview.net/](https://openreview.net/forum?id=BJgy96EYvr)  
739 [forum?id=BJgy96EYvr](https://openreview.net/forum?id=BJgy96EYvr).
- 740 Yongjie Wang, Yuchen Niu, Wenying Zhu, Wenqiang Chen, Qiong Li, and Tao Wang. Predicting  
741 pedestrian crossing behavior at unsignalized mid-block crosswalks using maximum entropy deep  
742 inverse reinforcement learning. *IEEE Trans. Intell. Transp. Syst.*, 25(5):3685–3698, 2024. doi: 10.  
743 1109/TITS.2023.3326276. URL <https://doi.org/10.1109/TITS.2023.3326276>.  
744
- 745 Zhendong Wang, Jonathan J. Hunt, and Mingyuan Zhou. Diffusion policies as an expressive policy  
746 class for offline reinforcement learning. In *The Eleventh International Conference on Learn-  
747 ing Representations, ICLR 2023, Kigali, Rwanda, May 1-5, 2023*. OpenReview.net, 2023. URL  
748 <https://openreview.net/forum?id=AHvFDPi-FA>.
- 749 Wei Xiao, Tsun-Hsuan Wang, Chuang Gan, Ramin M. Hasani, Mathias Lechner, and Daniela Rus.  
750 Safediffuser: Safe planning with diffusion probabilistic models. In *The Thirteenth International  
751 Conference on Learning Representations, ICLR 2025, Singapore, April 24-28, 2025*. OpenRe-  
752 view.net, 2025. URL <https://openreview.net/forum?id=ig2wk7kK9J>.  
753
- 754 Chao Yu, Akash Velu, Eugene Vinitzky, Jiaxuan Gao, Yu Wang, Alexandre Bayen, and Yi Wu. The  
755 surprising effectiveness of ppo in cooperative multi-agent games. *Advances in neural information  
processing systems*, 35:24611–24624, 2022.

- Guanwen Zeng, Jianxi Gao, Louis Shekhtman, Shengmin Guo, Weifeng Lv, Jianjun Wu, Hao Liu, Orr Levy, Daqing Li, Ziyou Gao, H. Eugene Stanley, and Shlomo Havlin. Multiple metastable network states in urban traffic. *Proceedings of the National Academy of Sciences*, 117(30):17528–17534, 2020. doi: 10.1073/pnas.1907493117. URL <https://www.pnas.org/doi/abs/10.1073/pnas.1907493117>.
- Lulu Zheng, Jiarui Chen, Jianhao Wang, Jiamin He, Yujing Hu, Yingfeng Chen, Changjie Fan, Yang Gao, and Chongjie Zhang. Episodic multi-agent reinforcement learning with curiosity-driven exploration. In Marc’Aurelio Ranzato, Alina Beygelzimer, Yann N. Dauphin, Percy Liang, and Jennifer Wortman Vaughan (eds.), *Advances in Neural Information Processing Systems 34: Annual Conference on Neural Information Processing Systems 2021, NeurIPS 2021, December 6-14, 2021, virtual*, pp. 3757–3769, 2021. URL <https://proceedings.neurips.cc/paper/2021/hash/1e8ca836c962598551882e689265c1c5-Abstract.html>.
- Meng Zhou, Ziyu Liu, Pengwei Sui, Yixuan Li, and Yuk Ying Chung. Learning implicit credit assignment for cooperative multi-agent reinforcement learning. In *Neural Information Processing Systems*, 2020.
- Mingyan Zhou, Biao Wang, and Xiatao Sun. Developing trajectory planning with behavioral cloning and proximal policy optimization for path-tracking and static obstacle nudging. *CoRR*, abs/2409.05289, 2024. doi: 10.48550/ARXIV.2409.05289. URL <https://doi.org/10.48550/arXiv.2409.05289>.
- Zhengbang Zhu, Hanye Zhao, Haoran He, Yichao Zhong, Shenyu Zhang, Yong Yu, and Weinan Zhang. Diffusion models for reinforcement learning: A survey. *CoRR*, abs/2311.01223, 2023. doi: 10.48550/ARXIV.2311.01223. URL <https://doi.org/10.48550/arXiv.2311.01223>.
- Zhengbang Zhu, Minghuan Liu, Liyuan Mao, Bingyi Kang, Minkai Xu, Yong Yu, Stefano Ermon, and Weinan Zhang. Madiff: Offline multi-agent learning with diffusion models. In Amir Globersons, Lester Mackey, Danielle Belgrave, Angela Fan, Ulrich Paquet, Jakub M. Tomczak, and Cheng Zhang (eds.), *Advances in Neural Information Processing Systems 38: Annual Conference on Neural Information Processing Systems 2024, NeurIPS 2024, Vancouver, BC, Canada, December 10 - 15, 2024*, 2024. URL [http://papers.nips.cc/paper\\_files/paper/2024/hash/07e278a120830b10aae20cc600a8c07b-Abstract-Conference.html](http://papers.nips.cc/paper_files/paper/2024/hash/07e278a120830b10aae20cc600a8c07b-Abstract-Conference.html).

## A APPENDIX

### A.1 IMPLEMENTATION DETAIL

This subsection describes the specific Markov process used in the experiment. Before giving the Markov Decision Process, we first introduce the topology of road network.

**Directed topology of a road network.** The directed topology of a road network is defined as  $G = (V, E, X)$ .  $V$  represents the roads in the network, assume there are  $m$  roads in  $G$ .  $E$  denotes the connections between roads, assume there are  $|E|$  connections.  $X$  indicates traffic status for each road.

Figure 6(a) illustrates the physical road network. Figure 6(b) shows its corresponding topological graph  $G = (V, E, X)$ .  $V = \{road_1, \dots, road_{12}\}$ ,  $E = \{road_7 \rightarrow road_4, road_7 \rightarrow road_2, \dots, road_4 \rightarrow road_{11}\}$ .  $road_7 \rightarrow road_4$  indicates vehicles can travel from  $road_7$  to  $road_4$ .  $X = \{queue_{road_1}, \dots, queue_{road_{12}}\}$  represents the traffic status of roads, for example queue length of each road. Actually,  $G$  can be viewed as the complete MDP for learning  $Q_\theta(x, a)$ . We detail each component below.

**State  $s$ .**  $s$  represents the current state, indicating the road and queue length of the vehicle,  $s = (road, queue_{road})$ .

**Action  $a$ .**  $a$  is the navigation actions such as left turn, straight ahead, or right turn on the current road to the next road. For example, if the vehicle is in  $road_4$ ,  $a = (road_{10}, road_{11})$ .

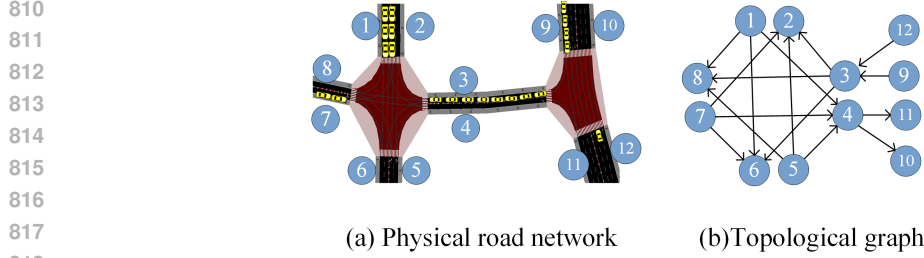


Figure 6: Topological graph of road network.

**Reward**  $r$ .  $r$  is the reward obtained by steering to the next road at the current road.

**Policy**  $Q_\theta(s, a)$ .  $Q_\theta(s, a)$  determines action  $a$  for each state  $s$ .  $Q_\theta(s, a)$  maximizes cumulative discounted rewards, such that  $Q_\theta(s, a) = \text{Argmax} \left[ \sum_{t=1}^T \gamma^t r(s_t, a_t) \right]$ .  $\gamma \in [0, 1)$  represents the discount factor.

For example,  $((road_9, queue_{road_9}), road_3, (road_3, queue_{road_3}), r)$  is a piece of experience. It means that if vehicle in state  $(road_9, queue_{road_9})$  takes action  $road_3$ , it will reach state  $(road_3, queue_{road_3})$  and will get reward  $r$ . Let  $\tau = \{(s_i, a_i, s'_i)_{t=1}^T\}_{i=1}^{|E|}$  denote the offline trajectory for the vehicle.  $|E|$  is the number of edges in  $G$ . Since the connection from road  $s$  to  $s'$  is fixed in a given traffic map, the number of experiences is exactly  $|E|$ .

However, if the local reward is not based on distance (such as “avoiding traffic lights” or “taking the favourite path”), reinforcement learning alone does not necessarily guarantee that the vehicle will reach the destination. This paper adopts a hybrid navigation policy of A\* + Q-learning. For example, in Figure 6, when vehicle A travels from  $road_9$  to  $road_2$  using the A\* alone. It selects the shortest path  $road_9 \rightarrow road_3 \rightarrow road_2$ . However, this path is the most congested. By combining A\* with  $Q_\theta(x, a)$ , roads with longer queue lengths are assigned with lower rewards. The optimal path changes to a longer but less congested route:  $road_9 \rightarrow road_{11} \rightarrow \dots \rightarrow road_5 \rightarrow road_2$ .

**A\* + Q-learning navigation policy.** The traditional A\* navigation policy defines the total cost function as  $r(s) = g(s) + h(s)$ .  $g(s)$  represents the actual cost from the starting point to the current node.  $h(s)$  is the heuristic estimation function that estimates the potential cost to the endpoint. We designed  $h(s)$  in Formula 12 by using an A\* + Q-learning hybrid mechanism that not only meets routing requirements but also reflecting the vehicle’s unique preferences.

$$h(s) = \|loc(s_{\text{current}}) - loc(s_{\text{end}})\|_2 - \lambda \cdot Q_\theta(s, a) \quad (12)$$

where  $\|loc(s_{\text{current}}) - loc(s_{\text{end}})\|_2$  is the Manhattan distance between the current road of the vehicle and the destination of the vehicle.  $\lambda$  is an adjustable hyper-parameter.  $Q_\theta(s, a)$  is the preference policy learned by reinforcement learning.  $s$  is the current state of the vehicle.  $a$  is the navigation actions such as turning left to the next road. A\* tries to minimize the cost of the path, while DQN tries to maximize the reward. So, we use the negative values of  $Q_\theta(s, a)$  in Formula 12.

We employ offline Deep Q-learning (Kostrikov et al., 2022) to train  $Q_\theta(x, a)$ . The training process is not performed in real-time but completed once the vehicle first enters the road network by using the offline dataset.

## A.2 EXPERIMENTAL RESULTS IN SCENARIO 2 AND SCENARIO 3

Figure 7~9 shows the results. The horizontal axis represents the minutes. The vertical axis represents the  $MQL$  value of a specific road.

864  
865  
866  
867  
868  
869  
870  
871  
872  
873  
874  
875  
876  
877  
878  
879  
880  
881  
882  
883  
884  
885  
886  
887  
888  
889  
890  
891  
892  
893  
894  
895  
896  
897  
898  
899  
900  
901  
902  
903  
904  
905  
906  
907  
908  
909  
910  
911  
912  
913  
914  
915  
916  
917

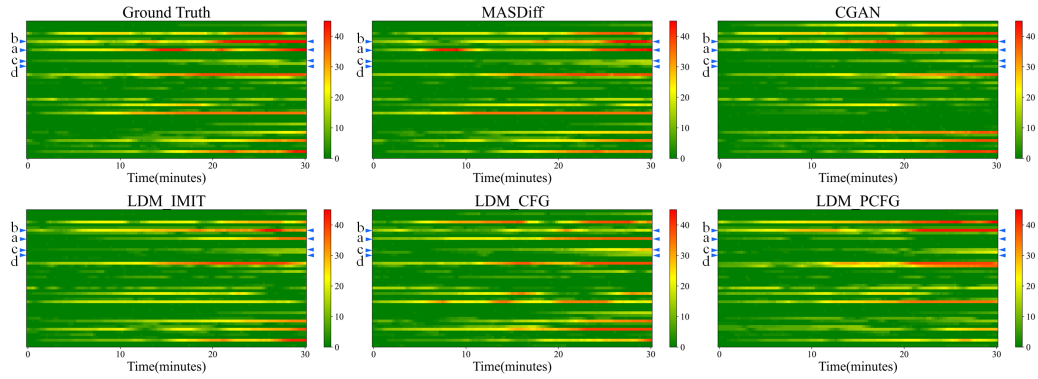


Figure 7:  $\{MQL_t^*\}_{t=1}^{30}$  and  $\{MQL_t\}_{t=1}^{30}$  of selected roads in scenario 1.

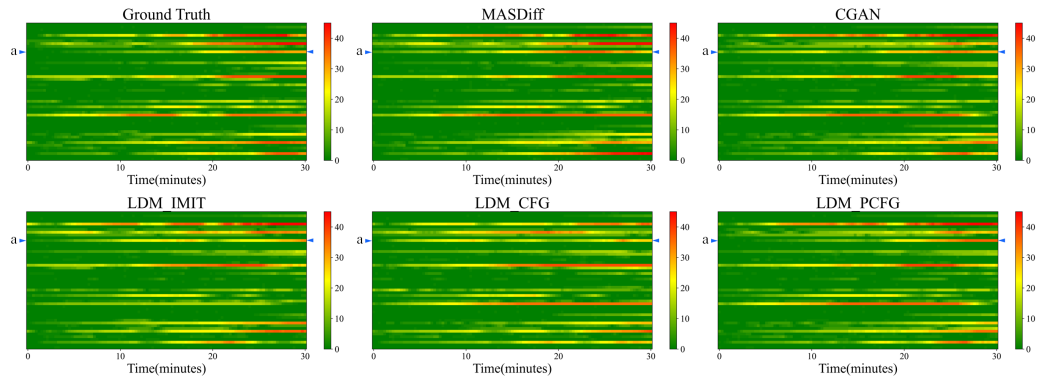


Figure 8:  $\{MQL_t^*\}_{t=1}^{30}$  and  $\{MQL_t\}_{t=1}^{30}$  of selected roads in scenario 2.

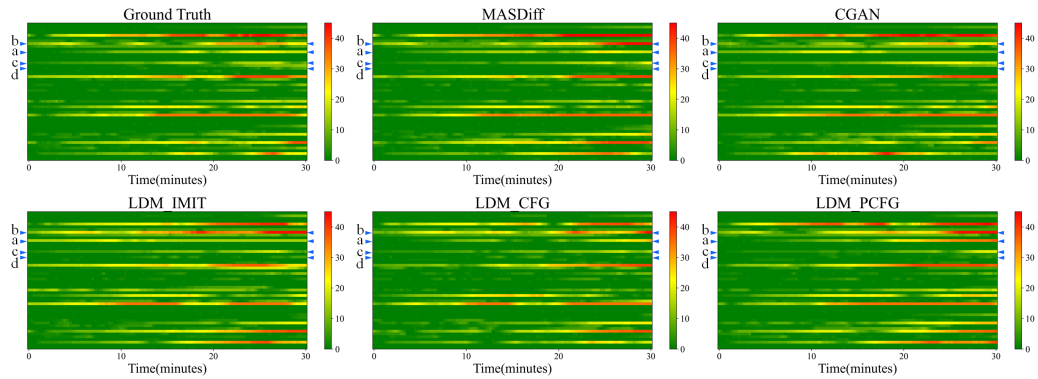


Figure 9:  $\{MQL_t^*\}_{t=1}^{30}$  and  $\{MQL_t\}_{t=1}^{30}$  of selected roads in scenario 3.