

LLM-Driven Stationarity-Aware Expert Demonstrations for Multi-Agent Reinforcement Learning in Mobile Systems

Tianyang Duan, Zongyuan Zhang, Zheng Lin, Songxiao Guo, Xiuxian Guan, Guangyu Wu, Zihan Fang, Haotian Meng, Xia Du, Ji-Zhe Zhou, Heming Cui, *Member, IEEE*, Jun Luo, *Fellow, IEEE*, and Yue Gao, *Fellow, IEEE*

Abstract—Multi-agent reinforcement learning (MARL) has been increasingly adopted in many real-world applications. While MARL enables decentralized deployment on resource-constrained edge devices, it suffers from severe non-stationarity due to the synchronous updates of agent policies. This non-stationarity results in unstable training and poor policy convergence, especially as the number of agents increases. In this paper, we propose RELED, a scalable MARL framework that integrates large language model (LLM)-driven expert demonstrations with autonomous agent exploration. RELED incorporates a Stationarity-Aware Expert Demonstration module, which leverages theoretical non-stationarity bounds to enhance the quality of LLM-generated expert trajectories, thus providing high-reward and training-stable samples for each agent. Moreover, a Hybrid Expert-Agent Policy Optimization module adaptively balances each agent’s learning from both expert-generated and agent-generated trajectories, accelerating policy convergence and improving generalization. Extensive experiments with real city networks based on OpenStreetMap demonstrate that RELED achieves superior performance compared to state-of-the-art MARL methods.

Index Terms—Multi-agent reinforcement learning, large language model, expert demonstration generation, non-stationarity.

I. INTRODUCTION

With the significant enhancement of computational capabilities in edge devices such as vehicles and robots, their ability

to execute complex tasks and achieve real-time interactions has been substantially improved [1]–[11]. Multi-agent systems (MAS) have attracted increasing attention due to their unique advantages in distributed collaboration and scalability. Among various MAS paradigms, multi-agent reinforcement learning (MARL) has emerged as a prominent direction in recent years [12], [13]. MARL leverages lightweight neural networks as policy functions for each agent, markedly reducing inference overhead and latency, and thus facilitating deployment on resource-constrained edge devices [14]. Furthermore, MARL typically employs centralized value estimation for joint policy optimization during training, while enabling decentralized decision-making based on local observations during inference [15]. This paradigm effectively achieves low-latency responses and minimizes inter-agent communication overhead, leading to its successful application in domains such as computation offloading [16]–[19], drone scheduling [20], [21], and urban traffic management [22], [23].

Despite these advantages, MARL still faces significant non-stationarity challenges [24]. Existing MARL training frameworks can be broadly categorized into fully decentralized frameworks [25], [26] and centralized training with decentralized execution (CTDE) [27], [28]. In fully decentralized settings, each agent optimizes its policy from local observations and individual rewards. This independence ignores cross-agent behavioral dependencies and results in policy conflicts that degrade overall system performance. CTDE mitigates these issues by using centralized value estimation during training to promote cooperative behavior. Nevertheless, as the scale of the multi-agent system increases, the joint state-action space grows exponentially, leading to a dramatic increase in the cost of policy evaluation and exacerbating estimation errors of each agent’s action value [29]. The resulting bias often drives agents toward homogenized policies, reducing behavioral diversity and hindering exploration of globally optimal solutions. Fundamentally, these issues stem from learning-induced non-stationarity. Each agent’s environment includes other agents that are simultaneously updating their policies. Any update by one agent alters the environment dynamics observed by the others, violating the stationarity assumptions underlying standard reinforcement learning algorithms [30].

Fortunately, large language models (LLMs) trained on extensive corpora have shown strong reasoning and decision-making in complex settings [31]–[36], making them promising

T. Duan, Z. Zhang, S. Guo, X. Guan, and H. Cui are with the Division of Computer Science, The University of Hong Kong, Hong Kong SAR, China (e-mail: tyduan@cs.hku.hk; zyzhang2@cs.hku.hk; sxguo@connect.hku.hk; xxguan@cs.hku.hk; heming@cs.hku.hk).

Z. Lin is with the Department of Electrical and Electronic Engineering, The University of Hong Kong, Hong Kong SAR, China (e-mail: linzheng@eee.hku.hk).

G. Wu is with the Department of Computer Science and Technology, Peking University, Beijing, China (e-mail: gywu9908@163.com).

Z. Fang is with the Department of Computer Science, City University of Hong Kong, Hong Kong SAR, China (e-mail: zihanfang3-c@my.cityu.edu.hk).

H. Meng is an employee in the China Unicom Digital Technology, China Unicom co.,Ltd, China (e-mail: 790498573@qq.com).

X. Du is with the School of Computer and Information Engineering, Xiamen University of Technology, Xiamen, China (email: duxia@xmut.edu.cn).

J. Zhou is with the School of Computer Science, Engineering Research Center of Machine Learning and Industry Intelligence, Sichuan University, Chengdu, China (email: yb87409@um.edu.mo).

J. Luo is with the College of Computing and Data Science, Nanyang Technological University, Singapore (e-mail: junluo@ntu.edu.sg).

Y. Gao is with the Institute of Space Internet, Fudan University, Shanghai, China, and the School of Computer Science, Fudan University, Shanghai, China (e-mail: gao.yue@fudan.edu.cn).

tools for generating expert knowledge and demonstrations to improve policy optimization in MARL. However, LLMs operate from prompts and prior knowledge rather than direct interaction with the environment, they offer limited exploration and thus tend to recommend actions only for familiar or well-specified states [37]. This reduces demonstration diversity and state-space coverage, which in turn hampers policy generalization to unvisited or poorly observed states in MARL. Moreover, the quality of LLM-generated demonstrations is highly dependent on the input context [38]. In partially observable environments, where agents can only obtain incomplete observations of the true state, LLM-generated demonstrations tend to guide agents toward suboptimal behaviors, impeding convergence to optimal policies.

To address these research gaps, we propose RELED, a scalable multi-agent policy optimization framework that integrates LLM-driven expert demonstrations with agent autonomy. Firstly, we develop a Stationarity-Aware Expert Demonstration (SED) module in RELED based on LLMs, which quantifies and theoretically bounds the environmental non-stationarity induced by simultaneous multi-agent policy updates. Refined by feedback metrics derived from our theoretical analysis, the LLM in SED module generates targeted demonstration results for each agent, aiming to maximize high cumulative rewards while constraining non-stationarity. Secondly, we develop a Hybrid Expert-Agent Policy Optimization (HPO) module in RELED. In HPO module, each agent combines LLM-generated demonstrations from the SED module with its own autonomous exploration experiences by independently optimizing a hybrid policy loss function. This hybrid policy optimization approach not only accelerates policy convergence but also mitigates the sample diversity and performance limitations in purely LLM-generated demonstrations, thus improving both learning efficiency and generalization for each agent. The key contributions of this work are summarized as follows:

- We propose RELED, a scalable MARL framework that couples LLM-based expert demonstrations with decentralized exploration, achieving efficient training.
- We derive a performance bound with two computable metrics, namely the reward volatility index and the policy divergence index, to quantify non-stationarity.
- We propose a stationarity-aware LLM refinement mechanism, which uses the above metrics as feedback to iteratively update instructions to generate high-quality demonstrations.
- We propose a hybrid policy optimization mechanism that transitions from imitation to exploration by adaptively balancing the weights of expert and agent-generated samples.
- We empirically evaluate RELED on real city networks based on OpenStreetMap. The results demonstrate that RELED outperforms state-of-the-art MARL baselines.

The rest of the paper is organized as follows. Section II discusses related work. Section III presents the RELED framework. Section IV describes the implementation of RELED. Section V reports the evaluation results. Finally, Section VI concludes the paper.

II. RELATED WORK

MARL has attracted significant attention for its ability to address complex cooperative and competitive tasks involving multiple agents [24]. A prominent category of MARL algorithms is value-based methods, which estimate the expected returns of joint actions through action-value functions (Q-functions) and derive optimal joint policies accordingly [39]–[41]. QMIX enhances decentralized policy learning in MARL by employing mixing networks to combine individual agent Q-values into a joint action-value function while maintaining monotonicity [27]. Building on QMIX, HMDQN integrates hierarchical reinforcement learning to mitigate sparse reward challenges in multi-robot tasks [42]. Another key category is policy-based methods, which directly optimize stochastic policies represented as probability distributions over actions conditioned on states [25], [43]. HATRPO extends trust region policy optimization to multi-agent scenarios, ensuring stable and efficient collaborative learning through advantage decomposition and sequential policy updates [44]. MACPO incorporates safety constraints into MARL, ensuring agents consistently satisfy constraints during policy updates via trust region restrictions [45]. Other methods employing deterministic policies output actions directly rather than probability distributions, improving policy robustness by leveraging global information for centralized training [28], [46].

Recent advances in LLMs have broadened their knowledge and reasoning abilities, leading to research on using LLMs to improve the performance of single-agent reinforcement learning [47]. Some studies focus on achieving effective representation mapping between LLMs’ text-based input-output and the RL environment’s state and action spaces [48], [49]. Qiu et al. [50] propose a multimodal vision-text framework for aligning non-text observations with actions via joint fine-tuning. LANCAR addresses ambiguous robotic locomotion instructions by leveraging LLMs to generate context-aware embeddings [51]. LLaRP uses frozen pre-trained LLMs with learnable adapters to enable strong generalization in vision-language policies [52]. Other studies have explored LLMs for designing reward functions, addressing the limitations of traditional hand-crafted rewards in terms of sparsity, bias, and scalability [53], [54]. Kwon et al. proposed using LLMs as proxy reward functions, allowing users to define RL objectives with natural language prompts and few-shot examples, thus reducing the need for expert demonstrations [55]. Song et al. developed a self-optimizing reward framework where LLMs iteratively generate and refine reward functions from natural language instructions, replacing manual reward design in continuous control tasks [56]. Several studies [57], [58] have leveraged LLMs to enhance cooperation or improve sample efficiency in MARL. However, leveraging LLMs to address the non-stationarity bottleneck in MARL systems has received little attention and remains an unexplored direction.

III. RELED FRAMEWORK

A. System Model

Decentralized Partially Observable Markov Process. In most MAS application scenarios, agents collaboratively ac-

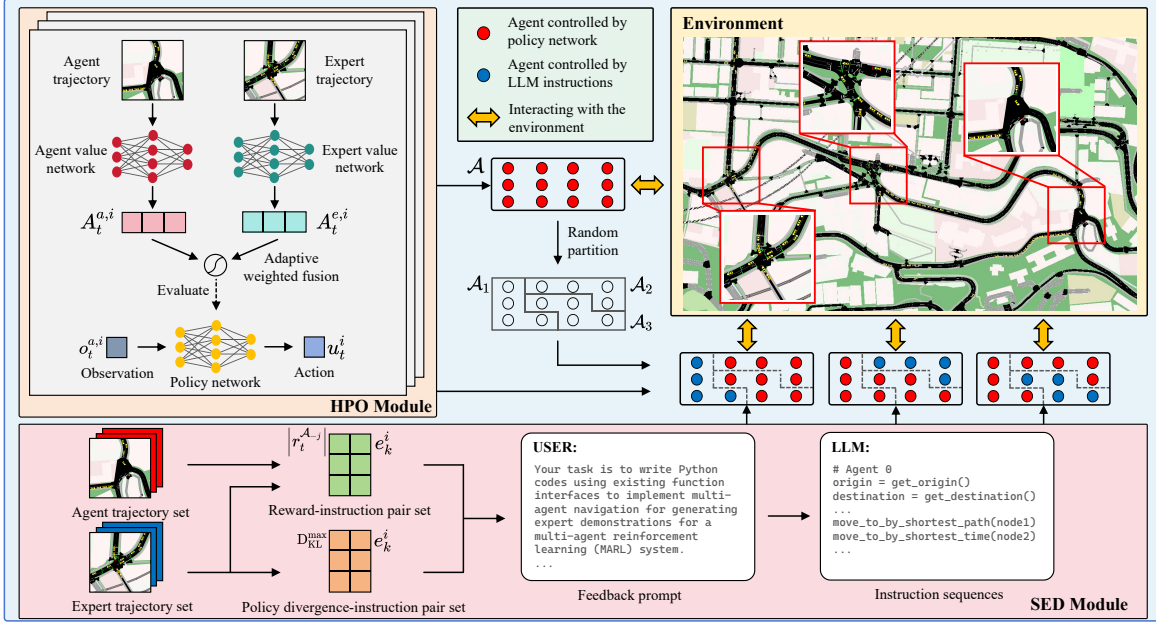


Fig. 1. An overview of RELED framework.

comply tasks in a shared environment while only receiving partial local observations. As a result, MARL is formulated as a decentralized partially observable Markov process (Dec-POMDP) [59], which is formally defined by the tuple $\langle \mathcal{A}, \mathcal{S}, \mathcal{U}, P, \mathcal{R}, \Omega, \mathcal{O}, \rho_0, \gamma \rangle$. Specifically, $\mathcal{A} = \{1, \dots, n\}$ represents the set of n agents, and \mathcal{S} denotes the state space. The joint action space is given by $\mathcal{U} = \prod_{i=1}^n \mathcal{U}^i$, where \mathcal{U}^i denotes the action space of agent i . The state transition function $P: \mathcal{S} \times \mathcal{U} \times \mathcal{S} \rightarrow [0, 1]$ defines the probability of transitioning to the next state given the current state and joint action. The reward function $\mathcal{R} = \{R^i\}_{i \in \mathcal{A}}$, where $R^i: \mathcal{S} \times \mathcal{U} \times \mathcal{S} \rightarrow \mathbb{R}$ assigns a reward to each agent. Ω denotes the observation space, and $\mathcal{O} = \{O^i\}_{i \in \mathcal{A}}$, with $O^i: \mathcal{S} \times \mathcal{U}^i \rightarrow \Omega$ as the observation function for agent i . The initial state distribution $\rho_0: \mathcal{S} \rightarrow [0, 1]$ defines the probability of each initial state. The discount factor $\gamma \in [0, 1]$ determines the relative importance of future rewards.

Interaction Process and Agent Trajectory. In Dec-POMDP, the policy $\pi^i: \Omega \times \mathcal{U}^i \rightarrow [0, 1]$ denotes the probability distribution over actions for agent i . An episode begins with the environment sampling an initial state $s_0 \sim \rho_0$. At each time step t , agent i receives a local observation $o_t^i \in \Omega$ according to its observation function $O^i(s_t)$. Each agent then selects an action $u_t^i \sim \pi^i(\cdot | o_t^i)$, forming the joint action $\mathbf{u}_t = (u_t^1, \dots, u_t^n)$. The environment transitions to $s_{t+1} \sim P(\cdot | s_t, \mathbf{u}_t)$ and provides each agent with a reward $r_t^i = R^i(s_t, \mathbf{u}_t, s_{t+1})$. The process continues, with each agent receiving an observation-action-reward tuple at each time step, which is formalized as follows:

$$\tau^a \sim \text{Env}(\pi | s_0, P), \quad (1)$$

where $\tau^a = \{\tau^i | i \in \mathcal{A}\}$ denotes the agent trajectory set, $\tau^i = (o_0^i, u_0^i, r_0^i, o_1^i, u_1^i, r_1^i, \dots)$ denotes the observation tra-

jectory obtained by agent i , and $\text{Env}(\cdot)$ denotes the environment.

Objective Function. The objective for each agent is to maximize the expected return, defined as the expected cumulative discounted reward over its trajectory:

$$J^i(\pi) = \mathbb{E}_{\rho_0, \pi, P} \left[\sum_{t=0}^{\infty} \gamma^t r_t^i \right], \quad (2)$$

where $\pi(\cdot | \mathbf{o}_t) = \prod_{i \in \mathcal{A}} \pi^i(\cdot | o_t^i)$ denotes the joint policy of all agents. The overall objective is to maximize the sum of individual agent objectives:

$$\begin{aligned} J(\pi) &= \sum_{i \in \mathcal{A}} J^i(\pi) \\ &= \mathbb{E}_{\rho_0, \pi, P} \left[\sum_{t=0}^{\infty} \sum_{i \in \mathcal{A}} \gamma^t r_t^i \right]. \end{aligned} \quad (3)$$

B. System Overview

In this section, we propose RELED, a scalable policy optimization framework that integrates LLM-driven expert demonstrations with agent-driven autonomous exploration. The framework consists of two key modules: a Stationarity-Aware Expert Demonstration (SED) module and a Hybrid Expert-Agent Policy Optimization (HPO) module, as outlined below:

- To mitigate the intrinsic non-stationarity arising from simultaneous policy updates in MARL systems, we propose the Stationarity-Aware Expert Demonstration (SED) module (see Section III-D). Based on a theoretically derived upper bound on environmental non-stationarity, the SED module quantitatively estimates the impact of policy changes and leverages LLMs to generate targeted expert demonstration trajectories for each agent. By incorporating both the reward volatility index and

the policy divergence index, which are derived from our theoretical analysis, as feedback, the generated expert demonstrations can effectively stabilize the environment and accelerate policy convergence for all agents.

- To address the limited coverage of the state-action space and the poor generalization associated with pure imitation learning from expert demonstrations, we propose the Hybrid Expert-Agent Policy Optimization (HPO) module (see Section III-E). The HPO module separately evaluates policy losses on both expert-generated and agent-generated trajectories, and integrates them through a dynamically adjusted weighting scheme. This adaptive mechanism enables agents to gradually shift from expert-driven imitation to autonomous exploration as training progresses, accelerating policy convergence and mitigating suboptimal solutions.

As illustrated in Figure 1, during each iteration of expert demonstration generation, the agent set \mathcal{A} is randomly partitioned into m disjoint subsets. Agents within each subset interact with the environment by following instruction sequences generated by the SED module, while the remaining agents execute their current policies. Within the SED module, reward-instruction and policy divergence-instruction pair sets are constructed by analyzing the resulting trajectory sets. These pair sets then provide feedback to the LLM, enabling iterative refinement of the instruction sequences. The HPO module adopts a fully decentralized training framework, where each agent independently optimizes its policy by leveraging both its own trajectories and those generated via instruction sequences. The agent and expert advantages, $A_t^{a,i}$ and $A_t^{e,i}$, are computed based on their respective value functions and are used to define separate policy losses. An adaptive weighted fusion mechanism dynamically adjusts the contributions of agent and expert losses according to the similarity between agent and expert trajectories, achieving hybrid policy optimization.

C. Theoretical Analysis of Non-Stationarity in MARL System

In this section, we present a theoretical analysis of environmental non-stationarity in MARL systems, serving as the foundation for the SED module.

A fundamental challenge in MARL arises from the inherent non-stationarity of the environment, induced by the simultaneous policy updates of multiple agents. As each agent updates its policy, the environment perceived by other agents changes dynamically, leading to biased estimates of its objective functions. This bias stems from the difficulty in distinguishing whether reward fluctuations originate from an agent's own actions or from the evolving policies of other agents. Therefore, it is essential to quantify the impact of each agent's policy changes on the value estimations of other agents. Such quantification enables the targeted refinement of LLM-generated expert demonstrations, thereby enhancing the stability and performance of MARL systems.

However, directly assessing the impact of each agent's policy changes on the objectives of other agents results in sample complexity and computational overhead growing multiplicatively with the number of agents. To address this problem, we partition the agent set \mathcal{A} into m disjoint subsets

$\{\mathcal{A}_1, \dots, \mathcal{A}_m\}$, and focus on the joint objective function of the other agents when the policies of agents within a specific subset \mathcal{A}_j are changed. Specifically, for a given subset of agents \mathcal{A}_j , we first define the joint objective of the external agents as follows:

$$J^{\mathcal{A}-j}(\pi) = \sum_{i \notin \mathcal{A}_j} J^i(\pi) \\ = \mathbb{E}_{\rho_0, \pi, P} \left[\sum_{t=0}^{\infty} \gamma^t r_t^{\mathcal{A}-j} \right], \quad (4)$$

where $\mathcal{A}-j = \mathcal{A} \setminus \mathcal{A}_j$ denotes the external agents of subset \mathcal{A}_j , i.e., all agents not in \mathcal{A}_j , and $r_t^{\mathcal{A}-j} = \sum_{i \notin \mathcal{A}_j} r_t^i$ denotes the sum of rewards obtained by external agents of subset \mathcal{A}_j at time step t . We then derive the variation bounds of the external agents' objective function under different policy configurations of an agent subset, thereby quantifying the resulting environmental non-stationarity, as formalized in the following theorem:

Theorem 1. (Performance Bound on Group-Level Non-Stationarity) Consider two arbitrary sets of joint policies for an agent subset \mathcal{A}_j :

$$\hat{\pi}^{\mathcal{A}_j}(\cdot | \mathbf{o}_t^j) = \prod_{i \in \mathcal{A}_j} \hat{\pi}^i(\cdot | o_t^i), \\ \tilde{\pi}^{\mathcal{A}_j}(\cdot | \mathbf{o}_t^j) = \prod_{i \in \mathcal{A}_j} \tilde{\pi}^i(\cdot | o_t^i),$$

where $\mathbf{o}_t^j = (o_t^k)_{k \in \mathcal{A}_j}$ denotes the joint observation of the agent subset \mathcal{A}_j . Let $\pi^{\mathcal{A}-j} = \{\pi^i | i \in \mathcal{A} \setminus \mathcal{A}_j\}$ denote the set of policies for agents external to the subset \mathcal{A}_j , and define the joint policies $\hat{\pi}$ and $\tilde{\pi}$ as:

$$\hat{\pi}(\cdot | \mathbf{o}_t) = \hat{\pi}^{\mathcal{A}_j}(\cdot | \mathbf{o}_t^j) \prod_{k \notin \mathcal{A}_j} \pi^k(\cdot | o_t^k), \\ \tilde{\pi}(\cdot | \mathbf{o}_t) = \tilde{\pi}^{\mathcal{A}_j}(\cdot | \mathbf{o}_t^j) \prod_{k \notin \mathcal{A}_j} \pi^k(\cdot | o_t^k).$$

Then, the following bound holds:

$$J^{\mathcal{A}-j}(\hat{\pi}) - J^{\mathcal{A}-j}(\tilde{\pi}) \\ \leq \frac{\sqrt{2}}{(1-\gamma)^2} \max_{\hat{\tau}, \tilde{\tau}} |r_t^{\mathcal{A}-j}| \sum_{k \in \mathcal{A}_j} D_{\text{KL}}^{\max}(\hat{\pi}^k \| \tilde{\pi}^k), \quad (5)$$

where $\hat{\tau}$ and $\tilde{\tau}$ denote trajectory rollouts of the joint policies $\hat{\pi}$ and $\tilde{\pi}$, respectively, and $D_{\text{KL}}^{\max}(\hat{\pi}^k \| \tilde{\pi}^k) = \max_{o_t^k} D_{\text{KL}}(\hat{\pi}^k \| \tilde{\pi}^k)[o_t^k]$ denotes the maximal Kullback-Leibler (KL) divergence between policies $\hat{\pi}^k$ and $\tilde{\pi}^k$.

Proof. We build upon intermediate results from prior work on performance-difference bounds, extend them to partially observable multi-agent settings, and further derive a broader bound to facilitate subsequent estimation in practice. Specifically, we first assume $\gamma f(s') \equiv f(s)$ in Theorem 1 of Achiam et al. [60], and extend it to joint policies $\hat{\pi}$ and $\tilde{\pi}$, which yields:

$$J^{\mathcal{A}_j}(\hat{\pi}) - J^{\mathcal{A}_j}(\tilde{\pi}) \\ \leq \frac{1}{1-\gamma} \mathbb{E}_{\substack{s_t \sim d^{\hat{\pi}} \\ \mathbf{u}_t \sim \hat{\pi} \\ s_{t+1} \sim P}} \left[\left(\frac{\hat{\pi}(\mathbf{u}_t | \mathbf{o}_t)}{\tilde{\pi}(\mathbf{u}_t | \mathbf{o}_t)} - 1 \right) r_t^{\mathcal{A}-j} \right] \\ + \frac{2\gamma}{(1-\gamma)^2} \max_{s_t} \left| \mathbb{E}_{\substack{\mathbf{u}_t \sim \hat{\pi} \\ s_{t+1} \sim P}} [r_t^{\mathcal{A}-j}] \right| \mathbb{E}_{s_t \sim d^{\tilde{\pi}}} [D_{\text{TV}}(\hat{\pi} \| \tilde{\pi})[\mathbf{o}_t]], \quad (6)$$

where $d^{\tilde{\pi}}$ denotes the discounted future state distribution under the joint policy $\tilde{\pi}$:

$$d^{\tilde{\pi}}(s) = (1 - \gamma) \sum_{t=0}^{\infty} \Pr(s_t = s \mid \rho_0, P, \tilde{\pi}), \quad (7)$$

and $D_{TV}(\hat{\pi} \parallel \tilde{\pi})[\mathbf{o}_t]$ denotes the total variation distance between $\hat{\pi}$ and $\tilde{\pi}$ at \mathbf{o}_t :

$$D_{TV}(\hat{\pi} \parallel \tilde{\pi})[\mathbf{o}_t] = \frac{1}{2} \sum_{\mathbf{u}_t} |\hat{\pi}(\mathbf{u}_t \mid \mathbf{o}_t) - \tilde{\pi}(\mathbf{u}_t \mid \mathbf{o}_t)|. \quad (8)$$

We bound the first term on the right side of Eqn. 6:

$$\begin{aligned} & \frac{1}{1 - \gamma} \mathbb{E}_{\substack{s_t \sim d^{\tilde{\pi}} \\ \mathbf{u}_t \sim \tilde{\pi} \\ s_{t+1} \sim P}} \left[\left(\frac{\hat{\pi}(\mathbf{u}_t \mid \mathbf{o}_t)}{\tilde{\pi}(\mathbf{u}_t \mid \mathbf{o}_t)} - 1 \right) r_t^{\mathcal{A}-j} \right] \\ &= \frac{1}{1 - \gamma} \mathbb{E}_{\substack{s_t \sim d^{\tilde{\pi}} \\ s_{t+1} \sim P}} \left[\sum_{\mathbf{u}_t} [\hat{\pi}(\mathbf{u}_t \mid \mathbf{o}_t) - \tilde{\pi}(\mathbf{u}_t \mid \mathbf{o}_t)] r_t^{\mathcal{A}-j} \right] \\ &\leq \frac{1}{1 - \gamma} \max_{\mathbf{o}_t} \sum_{\mathbf{u}_t} |\hat{\pi}(\mathbf{u}_t \mid \mathbf{o}_t) - \tilde{\pi}(\mathbf{u}_t \mid \mathbf{o}_t)| \max_{\tilde{\pi}, \tilde{\tau}} |r_t^{\mathcal{A}-j}| \\ &= \frac{2}{1 - \gamma} \max_{\mathbf{o}_t} D_{TV}(\hat{\pi} \parallel \tilde{\pi})[\mathbf{o}_t] \max_{\tilde{\pi}, \tilde{\tau}} |r_t^{\mathcal{A}-j}| \\ &\leq \frac{\sqrt{2}}{1 - \gamma} \max_{\mathbf{o}_t} D_{KL}(\hat{\pi} \parallel \tilde{\pi})[\mathbf{o}_t] \max_{\tilde{\pi}, \tilde{\tau}} |r_t^{\mathcal{A}-j}|, \end{aligned} \quad (9)$$

where $D_{KL}(\hat{\pi} \parallel \tilde{\pi})[\mathbf{o}_t]$ denotes the KL divergence between $\hat{\pi}$ and $\tilde{\pi}$ at \mathbf{o}_t :

$$D_{KL}(\hat{\pi} \parallel \tilde{\pi})[\mathbf{o}_t] = \mathbb{E}_{\mathbf{u}_t \sim \hat{\pi}} \left[\log \frac{\hat{\pi}(\mathbf{u}_t \mid \mathbf{o}_t)}{\tilde{\pi}(\mathbf{u}_t \mid \mathbf{o}_t)} \right]. \quad (10)$$

Then, we bound the second term in Eqn. 6:

$$\begin{aligned} & \frac{2\gamma}{(1 - \gamma)^2} \max_{s_t} \left| \mathbb{E}_{\substack{\mathbf{u}_t \sim \hat{\pi} \\ s_{t+1} \sim P}} [r_t^{\mathcal{A}-j}] \right| \mathbb{E}_{s_t \sim d^{\tilde{\pi}}} [D_{TV}(\hat{\pi} \parallel \tilde{\pi})[\mathbf{o}_t]] \\ &\leq \frac{2\gamma}{(1 - \gamma)^2} \max_{\tilde{\pi}} |r_t^{\mathcal{A}-j}| \max_{\mathbf{o}_t} D_{TV}(\hat{\pi} \parallel \tilde{\pi})[\mathbf{o}_t] \\ &\leq \frac{\sqrt{2}\gamma}{(1 - \gamma)^2} \max_{\tilde{\pi}} |r_t^{\mathcal{A}-j}| \max_{\mathbf{o}_t} D_{KL}(\hat{\pi} \parallel \tilde{\pi})[\mathbf{o}_t]. \end{aligned} \quad (11)$$

By combining Eqn. 9 and Eqn. 11 and comparing with Eqn. 6, we obtain:

$$\begin{aligned} & J^{\mathcal{A}_j}(\hat{\pi}) - J^{\mathcal{A}_j}(\tilde{\pi}) \\ &\leq \frac{\sqrt{2}}{(1 - \gamma)^2} \max_{\tilde{\pi}, \tilde{\tau}} |r_t^{\mathcal{A}-j}| \max_{\mathbf{o}_t} D_{KL}(\hat{\pi} \parallel \tilde{\pi})[\mathbf{o}_t] \\ &\leq \frac{\sqrt{2}}{(1 - \gamma)^2} \max_{\tilde{\pi}, \tilde{\tau}} |r_t^{\mathcal{A}-j}| \sum_{k \in \mathcal{A}_j} D_{KL}^{\max}(\hat{\pi}^k \parallel \tilde{\pi}^k) \end{aligned} \quad (12)$$

This concludes the proof. \square

Remark. Theorem 1 decomposes the impact of non-stationarity induced by policy changes within an agent subset \mathcal{A}_j into two components: (i) the maximal fluctuation in the cumulative rewards of external agents, quantified by the *reward volatility index* $\max_{\tilde{\pi}, \tilde{\tau}} |r_t^{\mathcal{A}-j}|$, and (ii) the aggregate maximal shift in the policy distributions within the subset, quantified by the *policy divergence index* $\sum_{k \in \mathcal{A}_j} D_{KL}^{\max}(\hat{\pi}^k \parallel \tilde{\pi}^k)$. Theorem 1 demonstrates that even minor deviations in individual agent policies can collectively exacerbate instability in the objective functions of other agents. This highlights the importance of controlling individual policy changes to ensure stability in MARL systems.

Building on Theorem 1, we derive an immediate corollary that bounds the agent-level non-stationarity.

Corollary 1. (Performance Bound on Agent-Level Non-Stationarity) Consider two arbitrary policies $\hat{\pi}^i$ and $\tilde{\pi}^i$ for agent i , and a policy sets for $\mathcal{A} \setminus \{i\}$:

$$\pi^{-i} = \{\pi^1, \dots, \pi^{i-1}, \pi^{i+1}, \dots, \pi^n\}. \quad (13)$$

Define the joint policies $\hat{\pi}$ and $\tilde{\pi}$ as:

$$\begin{aligned} \hat{\pi}(\cdot \mid \mathbf{o}_t) &= \hat{\pi}^i(\cdot \mid \mathbf{o}_t^i) \prod_{j \neq i} \pi^j(\cdot \mid \mathbf{o}_t^j), \\ \tilde{\pi}(\cdot \mid \mathbf{o}_t) &= \tilde{\pi}^i(\cdot \mid \mathbf{o}_t^i) \prod_{j \neq i} \pi^j(\cdot \mid \mathbf{o}_t^j). \end{aligned}$$

Then, the following inequality holds:

$$\begin{aligned} & J^i(\hat{\pi}) - J^i(\tilde{\pi}) \\ &\leq \frac{\sqrt{2}}{(1 - \gamma)^2} \max_{\tilde{\pi}, \tilde{\tau}} \left| \sum_{j \neq i} r_t^j \right| \max_{\mathbf{o}_t^i} D_{KL}(\hat{\pi}^i \parallel \tilde{\pi}^i)[\mathbf{o}_t^i]. \end{aligned} \quad (14)$$

Proof. The result follows from Theorem 1 by setting $\mathcal{A}_j = \{i\}$. \square

Remark. Corollary 1 establishes an upper bound on the impact of an agent's policy change on the objective functions of other agents. However, applying this upper bound necessitates individual evaluation for each agent, resulting in a total of n estimations and a computational complexity that scales linearly with the number of agents. In contrast, Theorem 1 introduces a subset-based bound, reducing the required evaluations to m . Therefore, the bound estimation based on Theorem 1 reduces computational cost and thus is more scalable and practical for MARL systems.

D. Stationarity-Aware Expert Demonstration Module

In this section, we propose the SED module, which leverages LLMs to generate expert demonstrations for each agent. The SED module incorporates both the reward volatility index and the policy divergence index, derived from the theoretical results in Section III-C, as feedback to iteratively query LLMs for demonstrations with high expected returns and enhanced environmental stability.

Conventional expert demonstrations for MARL typically relies on rule-crafted or human-curated demonstrations [61], [62], which are costly to obtain and tailored to specific task scenarios. In contrast, LLMs, with their strong contextual reasoning capabilities and extensive prior knowledge, have shown great potential in generating customized expert demonstrations for MARL agents [63], [64]. Furthermore, by continuously providing environmental feedback to the LLM, the generated demonstrations can be dynamically adapted for each agent, thereby producing high-quality demonstration samples and further accelerating the convergence of agent policies.

Therefore, we employ an LLM-driven expert demonstration generation mechanism in the SED module. For clarity, we denote the LLM as $\Pi(d)$, where d is the prompt input. The objective of the LLM is to generate an instruction sequence for each agent, formalized as $\mathcal{E} = \{e_k^i \mid i \in \mathcal{A}, k = 0, 1, \dots\}$. Each instruction sequence e^i comprises ordered task steps that guide agent i to interact with the environment. To ensure

executability, each instruction e_j^i is selected from a predefined set \mathcal{C} . The environment then maps each instruction to an executable action sequence via an execution function $f: \mathcal{C} \rightarrow \mathcal{U}^l$, defined as:

$$f(e_k^i) = (u_{t_k}^i, u_{t_k+1}^i, \dots, u_{t_k+l_k^i-1}^i) \quad (15)$$

where $0 < l_k^i \leq l$ denotes the number of timesteps required for agent i to complete instruction e_k^i . In implementation, f corresponds to a set of programmatic functions that realize the semantics of instructions in \mathcal{C} . For example, in vehicle routing, an instruction such as `move_to_by_shortest_path(node_id)` is operationalized by invoking a function that directs the agent to the specified node via the shortest path. Further implementation details are provided in Section IV.

In Theorem 1, we quantitatively characterize the impact of policy changes within a subset of agents on the objective functions of other agents. This provides a theoretical bound for analyzing the non-stationarity induced in the environment when that subset updates its policy by imitating expert demonstrations. Accordingly, the SED module employs the reward volatility index and policy divergence index, both derived from this bound, as two feedback metrics for the LLM. Specifically, let $\tilde{\pi}$ and $\hat{\pi}$ denote the joint policies of the agent subset \mathcal{A}_j before and after the update, respectively. By minimizing the upper bound established in Theorem 1, the SED module constrains the fluctuations in the objective functions of external agents caused by policy shifts within \mathcal{A}_j . To this end, during each iteration, the SED module randomly partitions the agent set \mathcal{A} into m disjoint subsets $\{\mathcal{A}_1, \dots, \mathcal{A}_m\}$ to comprehensively evaluate the non-stationarity induced by different agent combinations. For each subset \mathcal{A}_j , agents in the subset select actions according to the instruction sequence $\mathcal{E}^{\mathcal{A}_j} = \{e^i \mid i \in \mathcal{A}_j, k = 0, 1, \dots\}$ generated by the LLM, executed via function f , while the remaining agents follow their respective policies. All agents interact with the environment to generate the trajectory set $\hat{\tau}$, formalized as:

$$\hat{\tau} = \tau^{e,j} \sim \text{Env}(\mathcal{E}^{\mathcal{A}_j}, \pi^{\mathcal{A}-j} \mid s_0, P, f), \quad (16)$$

where $\pi^{\mathcal{A}-j} = \prod_{i \notin \mathcal{A}_j} \pi^i$ denotes the joint policy of agents outside \mathcal{A}_j . For comparison, the agent trajectory set $\tilde{\tau}$ is obtained via rollout under the current joint policy (i.e., $\tilde{\tau} = \tau^a$).

To minimize the reward volatility index, expert demonstrations should ensure that updating the policy of the agent subset \mathcal{A}_j results in a decrease in the reward fluctuations of external agents, which can be formalized as the following constraint:

$$\max \left\{ \left| \min_{\tilde{\tau}} r_t^{\mathcal{A}-j} \right|, \left| \max_{\tilde{\tau}} r_t^{\mathcal{A}-j} \right| \right\} \leq \max_{\tilde{\tau}} \left| r_t^{\mathcal{A}-j} \right|. \quad (17)$$

However, in cooperative settings with coupled or global rewards, multiple agents receive identical rewards. In these scenarios, policy improvement within the agent subset \mathcal{A}_j can inevitably cause an increase in $\max_{\tilde{\tau}} r_t^{\mathcal{A}-j}$. Therefore, the SED module focuses on constraining the term $\left| \min_{\tilde{\tau}} r_t^{\mathcal{A}-j} \right|$ within the range $\max_{\tilde{\tau}} \left| r_t^{\mathcal{A}-j} \right|$, thereby optimizing worst-case reward fluctuations. Specifically, the SED module constructs

a timestep set \mathcal{T}^j , consisting of all timesteps in the expert demonstration trajectory set $\hat{\tau}$ where the external agent's reward is lower than the maximum absolute reward in the agent trajectory set $\tilde{\tau}$:

$$\mathcal{T}^j = \left\{ t \mid r_t^{\mathcal{A}-j} < -\max_{\tilde{\tau}} \left| r_t^{\mathcal{A}-j} \right|, r_t^{\mathcal{A}-j} \in \hat{\tau} \right\}, \quad (18)$$

The SED module further establishes the correspondence between the reward volatility index and LLM-generated instructions by constructing a reward-instruction pair set \mathcal{R}^j , which enables iterative feedback on the instruction sequence, as follows:

$$\mathcal{R}^j = \left\{ \left| r_t^{\mathcal{A}-j} \right|, e_k^i \mid t \in \mathcal{T}^j, i \in \mathcal{A}_j, l_{k-1}^i \leq t < l_k^i \right\}. \quad (19)$$

To minimize the policy divergence index, we estimate the KL divergence between the pre-update policy $\tilde{\pi}^i$ and the post-update policy $\hat{\pi}^i$ for each agent $i \in \mathcal{A}_j$. Given the trajectory set $\hat{\tau}$ and the current policy π^i , we approximate the updated policy $\hat{\pi}^i$ as follows:

$$\hat{\pi}^i(u^i \mid o_t^i) = \begin{cases} \delta(u^i - u_t^i), & \exists o_t^i, u_t^i \in \hat{\tau} \\ \pi^i(u^i \mid o_t^i), & \text{otherwise} \end{cases}, \quad (20)$$

where δ denotes the Dirac delta function. This approximation treats the observation-action pairs (o_t^i, u_t^i) in the trajectory $\hat{\tau}$ as deterministic action outputs of the policy, while retaining the original policy distribution for all other observation inputs. In this way, the influence of expert demonstrations on the updated policy of agent i is explicitly captured. Consequently, the KL divergence between the pre- and post-update policies for agent i at an observation o_t^i is given by:

$$\begin{aligned} D_{\text{KL}}(\hat{\pi}^i \parallel \tilde{\pi}^i)[o_t^i] &= \mathbb{E}_{u_t^i \sim \tilde{\pi}^i} \left[\ln \frac{\hat{\pi}^i(u_t^i \mid o_t^i)}{\tilde{\pi}^i(u_t^i \mid o_t^i)} \right] \\ &= \mathbb{E}_{o_t^i, u_t^i \in \hat{\tau}} \left[\ln \frac{1}{\tilde{\pi}^i(u_t^i \mid o_t^i)} \right] + \mathbb{E}_{o_t^i, u_t^i \notin \hat{\tau}} \left[\ln \frac{\pi^i(u_t^i \mid o_t^i)}{\tilde{\pi}^i(u_t^i \mid o_t^i)} \right] \\ &= \mathbb{E}_{o_t^i, u_t^i \in \hat{\tau}} \left[\ln \frac{1}{\tilde{\pi}^i(u_t^i \mid o_t^i)} \right]. \end{aligned} \quad (21)$$

Therefore, the maximum policy divergence for agent i can be expressed as:

$$D_{\text{KL}}^{\max}(\hat{\pi}^i \parallel \tilde{\pi}^i) = \max_{o_t^i, u_t^i \in \hat{\tau}} \ln \frac{1}{\tilde{\pi}^i(u_t^i \mid o_t^i)}, \quad (22)$$

and the timestep t_i^* corresponding to the maximum KL divergence in the trajectory set $\hat{\tau}$ is given by:

$$t_i^* = \underset{t}{\operatorname{argmax}} \ln \frac{1}{\tilde{\pi}^i(u_t^i \mid o_t^i)}, \quad (23)$$

where $(o_t^i, u_t^i) \in \hat{\tau}$. Based on the above, the SED module constructs the policy divergence-instruction pair set \mathcal{K}^j for each agent subset \mathcal{A}_j as:

$$\mathcal{K}^j = \{D_{\text{KL}}^{\max}(\hat{\pi}^i \parallel \tilde{\pi}^i), e_k^i \mid i \in \mathcal{A}_j, l_{k-1} \leq t_i^* < l_k\} \quad (24)$$

Through the above mechanism, the SED module constructs a reward-instruction pair set $\{\mathcal{R}^j\}_{j=1}^m$ and a policy divergence-instruction pair set $\{\mathcal{K}^j\}_{j=1}^m$ for each agent subset. The expert trajectory set are aggregated as $\tau^e =$

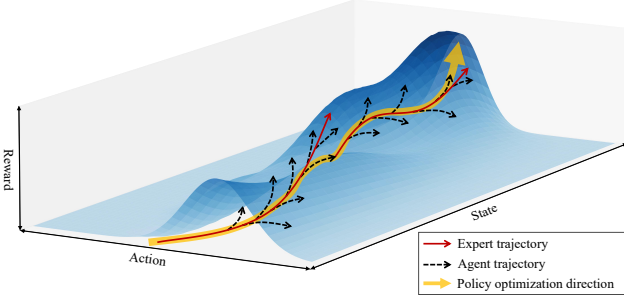


Fig. 2. The hybrid policy optimization mechanism adaptively combines expert demonstration guidance (red solid line) and autonomous exploration (black dashed line), facilitating the gradual convergence of agent’s policy towards the optimal solution.

$\{\tau^i \mid \tau^i \in \tau^{e,j}, i \in \mathcal{A}_j, j = 1, \dots, m\}$. Accordingly, we design two LLM inputs d : (i) initial prompt d_{init} , which provides a detailed description of the environment and reward structure to guide the LLM in generating task-targeted instruction sequences for each agent; and (ii) feedback prompt d_{fb} , which incorporates observation information of the expert trajectory set τ^e and instructs the LLM to revise instructions associated with significant reward fluctuations or excessive policy divergence, as identified via $\{\mathcal{R}^j\}_{j=1}^m$ and $\{\mathcal{K}^j\}_{j=1}^m$, respectively. A complete description of the prompts is provided in the Appendix A.

E. Hybrid Expert-Agent Policy Optimization Module

In this section, we propose the HPO module, a fully decentralized MARL training framework that adaptively combines expert demonstrations from Section III-D with autonomous exploration samples, thus accelerating policy convergence and enhancing agent performance.

The HPO module adopts a fully decentralized framework, where the training and execution of each agent are independent to ensure scalability. Specifically, as described in Section III-D, the SED module generates the expert trajectory set τ^e and the agent trajectory set τ^a . For each agent i , we denote its own trajectories by $\tau^{e,i} \in \tau^e$ and $\tau^{a,i} \in \tau^a$, where the superscripts a and e indicate trajectories derived from agent-environment interactions and expert instruction executions. To perform policy evaluation and subsequent policy fusion, we define the agent and expert value functions as follows:

$$V_i^x(o_t^{x,i}) = \mathbb{E}_{\tau^{x,i}} \left[\sum_{k=t}^{\infty} \gamma^{k-t} r_k^{x,i} \right], \quad (25)$$

where $x \in \{a, e\}$, $o_t^{x,i}$ and $r_k^{x,i}$ are the observation and reward at time step k in trajectory $\tau^{x,i}$. Subsequently, we compute the finite-horizon bootstrapped return for each trajectory using the agent and expert value functions, V_i^a and V_i^e , respectively:

$$R_t^{x,i} = \sum_{k=0}^{T^x-t-1} \gamma^k r_{t+k}^{x,i} + \gamma^{T^x-t} V_i^x(o_{T^x}^{x,i}) \quad (26)$$

where $r_{t+k}^{x,i}$ and $o_{T^x}^{x,i}$ are the reward and observation in trajectory $\tau^{x,i}$, and T^x represents the truncation length of the trajectory $\tau^{x,i}$. The agent value function V_i^a and the expert

Algorithm 1 RELED Training Procedure

Input: $K, \Pi, d_{\text{init}}, d_{\text{fb}}, q, n, m$.

Output: π .

```

1: Random initialization:  $\pi^i, V^{a,i}, V^{e,i}, i = 1, \dots, n$ .
2: Initialization:  $\Pi, d = d_{\text{init}}$ .
3: for  $k = 1$  to  $K$  do
4:    $\tau^a \sim \text{Env}(\pi \mid s_0, P)$  // Sample agent trajectories
5:   /* Stationarity-Aware Expert Demonstration Module: */
6:   if  $k \bmod q = 0$  then
7:     Randomly partition  $\mathcal{A}$  into  $\{\mathcal{A}_1, \dots, \mathcal{A}_m\}$ .
8:      $\mathcal{E} \sim \Pi(d)$  // Generate instruction sequences
9:      $d \leftarrow d_{\text{fb}}, \tau^e \leftarrow \emptyset$ 
10:    for  $j = 1$  to  $m$  do
11:       $\tau^{e,j} \sim \text{Env}(\mathcal{E}^{\mathcal{A}_j}, \pi^{\mathcal{A}-j} \mid s_0, P, f)$  // Sample expert trajectories
12:      Calculate  $\mathcal{R}^j$  and  $\mathcal{K}^j$  via Eqn. 19 and 24.
13:       $d \leftarrow d \cup \{\mathcal{R}^j, \mathcal{K}^j\}$ 
14:       $\tau^e \leftarrow \tau^e \cup \{\tau^i \mid \tau^i \in \tau^{e,j}, i \in \mathcal{A}_j\}$ 
15:    end for
16:     $d \leftarrow d \cup \tau^e$ 
17:  end if
18:  /* Hybrid Expert-Agent Policy Optimization Module: */
19:  for  $i = 1$  to  $n$  do
20:     $\tau^{a,i} \leftarrow \tau^a(i), \tau^{e,i} \leftarrow \tau^e(i)$  // Extract the trajectory of the  $i$ -th agent
21:    Calculate  $R_t^{a,i}, R_t^{e,i}$  and  $A_t^{a,i}, A_t^{e,i}$  via Eqn. 26 and 28.
22:    Update  $V_i^a$  and  $V_i^e$  via Eqn. 27.
23:    Update  $\pi^i$  via Eqn. 31.
24:  end for
25: end for

```

value function V_i^e are optimized by minimizing the mean squared error:

$$\mathcal{L}_{\text{value}}(V_i^x) = \mathbb{E}_{\tau^{x,i}} \left[\left(V_i^x(o_t^{x,i}) - R_t^{x,i} \right)^2 \right]. \quad (27)$$

Relying solely on expert demonstration samples often results in limited sample diversity, which constrains the agent’s policy generalization to unknown states. Furthermore, in partially observable environments, LLMs are unable to access complete state information, which may lead to suboptimal expert demonstrations. To address these limitations, we propose a hybrid policy optimization mechanism that integrates both expert demonstration samples and agent autonomous exploration samples. This approach enables the agent to further optimize its policy beyond what can be achieved by learning from expert knowledge alone. Specifically, we first compute the agent and expert advantages for the respective trajectories as follows:

$$A_t^{x,i} = R_t^{x,i} - V_i^x(o_t^{x,i}). \quad (28)$$

We then adopt the clipped surrogate objective [65] to construct the agent and expert policy losses using the agent and expert advantages, respectively:

$$\mathcal{L}^x(\pi^i) = \mathbb{E}_{\tau^{x,i}} \left[\min(\omega_t^{x,i}, \text{clip}(\omega_t^{x,i}, 1 - \epsilon, 1 + \epsilon)) A_t^{x,i} \right], \quad (29)$$

where $\omega_t^{x,i} = \pi^i(u_t^{x,i}|o_t^{x,i})/\pi^{\text{old},i}(u_t^{x,i}|o_t^{x,i})$ denotes the importance sampling ratio, $\pi^{\text{old},i}$ denotes the policy of agent i from the previous iteration, and $\epsilon \in (0, 1)$ represents the clipping coefficient. Consequently, we define the hybrid policy loss function \mathcal{L}_{mix} as follows:

$$\mathcal{L}_{\text{mix}}(\pi^i) = \alpha \mathcal{L}^a(\pi^i) + (1 - \alpha) \mathcal{L}^e(\pi^i), \quad (30)$$

where $\alpha = \exp(-k/K \cdot \text{D}_{\text{DTW}}(\tau^{a,i}, \tau^{e,i}))$ dynamically adjusts the weighting between the agent and expert policy losses, with k and K denoting the current and total training epochs, respectively, and $\text{D}_{\text{DTW}}(\tau^{a,i}, \tau^{e,i})$ denotes the dynamic time warping (DTW) distance, which measures the similarity between agent and expert trajectories via nonlinear temporal alignment. To promote autonomous exploration, we incorporate a maximum entropy regularization term into the hybrid policy loss function, which is defined as follows:

$$\mathcal{L}(\pi^i) = \mathcal{L}_{\text{mix}}(\pi^i) + \beta \mathbb{E}_{\tau^{a,i}}[\mathcal{H}(\pi^i(\cdot|o_t^{a,i}))], \quad (31)$$

where $\beta \in (0, 1]$ denotes the entropy regularization coefficient, and \mathcal{H} denotes the entropy of the policy distribution.

As shown in Figure 2, during the initial training phase, the agent lacks sufficient experience, which leads to a substantial discrepancy between its policy and the expert demonstrations. The resulting large DTW distance induces a low hybrid weight α that prioritizes expert-guided policy optimization. As training progresses and the agent's policy converges towards the expert's behavior, the DTW distance decreases, leading to a gradual increase in α and the optimization increasingly emphasizes experiences gained through agent's autonomous exploration. In parallel, the SED module will periodically leverage environmental feedback to refine and provide updated demonstrations, mitigating the non-stationarity resulting from current policy updates. With each refinement, the difference between the previous and new demonstrations enlarges the DTW distance, which reduces α and guides the agent to optimize its policy to follow the new demonstrations. This iterative mechanism, integrating both expert guidance and self-exploration, enables the agent to continuously improve its policy, thereby achieving efficient learning and improved generalization across complex tasks.

Algorithm 1 outlines the training procedure of RELED, where q denotes the sampling interval for expert demonstration. The training process begins with the initialization of the agent policy π^i , agent value function $V^{a,i}$, and expert value function $V^{e,i}$ for each agent i (line 1). Subsequently, the LLM Π is loaded, and the initial prompt d_{init} is set (line 2). At each training epoch, the agent trajectory set τ^a is collected by executing the current joint policy in the environment (line 4). To improve the stability of policy updates, the SED module adopts a sampling interval q , allowing each agent to utilize the same expert demonstration over multiple epochs (lines 6-17). During each demonstration generation, the agent set \mathcal{A} is randomly partitioned into m disjoint subsets (line 7). For each subset \mathcal{A}_j , expert trajectories $\tau^{e,j}$ are obtained by applying LLM-generated instructions to agents in \mathcal{A}_j , while the remaining agents follow their own policy decisions (line 11). The reward-instruction pair set \mathcal{R}^j and the policy divergence-instruction pair set \mathcal{K}^j are calculated to quantitatively assess the impact

of the demonstration on environment stationarity (line 12), and these results are integrated into the feedback prompt for LLM refinement (line 13). Subsequently, the observation information from the expert trajectory set τ^e is incorporated into the feedback prompts to provide the LLM with detailed context on expert instruction executions (line 16). In the HPO module, each agent extracts its own agent trajectory and expert trajectory for policy evaluation and improvement (line 20). Bootstrapped returns $R_t^{a,i}$, $R_t^{e,i}$ and advantages $A_t^{a,i}$, $A_t^{e,i}$ are calculated (line 21), serving as inputs to update agent and expert value functions via mean squared error loss (line 22), and to optimize policy π^i using the hybrid loss with entropy regularization (line 23).

IV. IMPLEMENTATION

This section describes the experimental setups and implementation details of RELED.

A. Environment Setups

1) *Simulation*: We evaluate RELED in multi-agent navigation tasks in realistic urban traffic environments constructed from OpenStreetMap (OSM) data [66] and simulated in SUMO [67]. OSM data are converted to SUMO via a standard netconvert pipeline, including geometry simplification, ramp inference, junction merging, standardizing intersections and signals while preserving lane counts, one-way attributes, and speed limits.

Each agent controls a single vehicle navigating from an assigned origin to a designated destination. For fair comparison, origin-destination pairs are sampled once at random from a predefined set and then fixed identically for all methods and runs. A single environment interaction step corresponds to one agent perform an action. The simulator advances in discrete steps last up to 2400 simulator steps in total, terminating early if all agents reach their destinations.

For each map, we consider two background traffic regimes to emulate different demand levels: (i) moderate flow, with 90 background vehicles distributed according to edge lengths; and (ii) congested flow, with 190 background vehicles and increased injection frequency on major arterials. Background vehicles has randomness in route assignment, departure times, and lane selection, and their origin-destination pairs are sampled from the same predefined set as the agents. The experiments span two representative urban networks:

- **Orlando**: A canonical grid with uniform intersections (Figure 3(a,b)), featuring 115 regulated intersections and 269 edges across a 1.99 km² area. This environment primarily consists of four-way junctions with consistent lane configurations.
- **Hong Kong**: An urban road network (Figure 3(c,d)) spanning 3.87 square kilometers, comprising 144 regulated intersections and 239 road segments, characterized by roads with extended edge lengths and a high prevalence of one-way traffic routes.

2) *Hyperparameters*: Each independent run consisted of 5×10^2 training epochs and each epoch comprising 10^3 environment interaction steps, accumulating a total of 5×10^5 interaction steps. Unless otherwise specified, all experiments

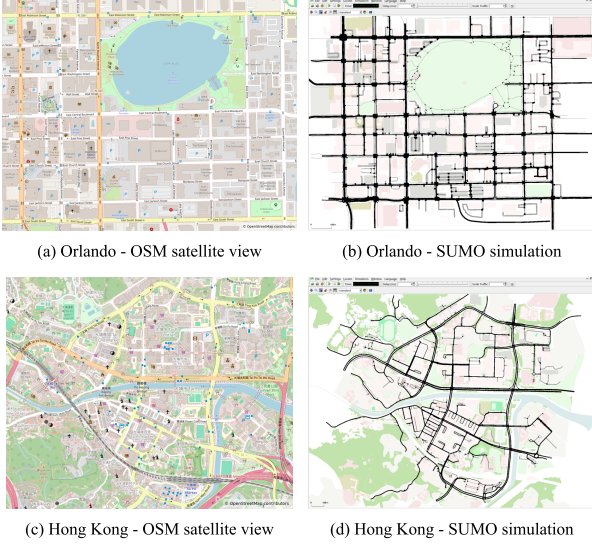


Fig. 3. Experimental scenarios.

TABLE I
ACTION SPACE AND OBSERVATION SPACE.

Action Space (Discrete, shape: 1)	
Index	Note
0	Select the i -th available edge ($i < m_{out}^*$)
Observation Space (Shape: $2m_{out} + 2$)	
Index	Note
[0	Current junction id
1]	Destination junction id
[2 + 2 <i>i</i>	Score of the i -th available edge ($i = 0..m_{out} - 1$)
[2 + 2 <i>i</i> + 1]	End junction id of the i -th edge ($i = 0..m_{out} - 1$)

m_{out} is the maximum number of outgoing edges at any junction.

were conducted with 10 agents operating in the environment. All algorithms are configured with identical optimizer parameters, with the learning rate $= 3 \times 10^{-4}$ and The discount factor $\gamma = 0.99$. Both policy networks and value estimation networks across all algorithms employed the same two-layer neural network structure with 128 neurons per layer. For LLM demonstration sampling, we employed a periodic update strategy with interval $q = 10$ epochs. For RELED, the agent subset partitioning parameter $m = 2$.

3) *Reward Function*: The reward at interaction step k is designed to encourage minimal travel time and progress toward the destination:

$$R_k = -(t_k - t_{k-1}) + \omega_d(d_{k-1} - d_k), \quad (32)$$

where R_k is the reward at step k ; t_k and t_{k-1} are the current and previous simulation times; ω_d is set to 1 in our experiments; and d_k and d_{k-1} are the current and previous Euclidean distances to the destination. The first term penalizes elapsed time, and the second term rewards reductions in distance to the destination—when the distance decreases, the agent receives a positive reward proportional to that decrease. Upon reaching the destination, the agent receives an additional terminal reward $R_{\text{term}} = 0.1 \cdot T_{\text{max}}$, where T_{max} is the maximum number of simulation steps.

TABLE II
PREDEFINED NAVIGATION INTERFACES PROVIDED FOR LLM

Function	Description
<code>move_to_by_shortest_path(node_id) → bool</code>	Move to target node via shortest distance path
<code>move_to_by_shortest_time(node_id) → bool</code>	Move to target node via shortest time path
<code>get_origin() → int</code>	Return the current agent's origin node ID
<code>get_destination() → int</code>	Return the current agent's destination node ID
<code>get_shortest_dist(target_node_id) → float</code>	Calculate shortest path distance to target node
<code>get_shortest_time(target_node_id) → float</code>	Calculate shortest time to target node
<code>get_nearest_node(x, y) → int</code>	Find node ID closest to coordinates (x, y)
<code>get_node_coord(node_id) → tuple[float, float]</code>	Get coordinates of specified node

4) *Observation Space*: Agents operate with partially observation space, with visibility limited to all outgoing edges of the junction they are approaching. The observation space of a single agent contains the current junction ID, the destination junction ID, and for each potential outgoing edge, its score and the end junction ID. The edge score ranges from 0 to 1, calculated as the ratio between the free flow travel time (length/max speed) and the estimated travel time. Detailed environment specifications are summarized in Table 1.

5) *Action Space*: The action space is a discrete selection of available road segments at each junction. When an agent reaches an intersection, it must choose one outgoing edge from m_{out} options. For each map, m_{out} is fixed and set to the maximum number of outgoing roads at any junction in that map. For junctions with fewer than m_{out} edges, the excess dimensions are masked with a -1 . For the Orlando and Hong Kong maps used in our experiments, $m_{out} = 4$.

B. Baselines

To evaluate the performance of RELED, we compare it with the following MARL methods:

- **IPPO** [26]: An extension of PPO where each agent is trained independently with its own policy and value networks, using on-policy data.
- **MAPPO** [25]: A CTDE method that trains decentralized policies together with a centralized value function that has access to global information during training.
- **QMIX** [27]: A value factorization approach that combines individual agent value functions through a monotonic mixing network to form a joint action-value function for decentralized control.

C. RELED Prompt Design

We implement a prompting architecture to generate high-quality LLM-based expert demonstrations for the RELED framework. We employ GPT-3.5 Turbo as our foundational model and structure our prompting approach in a two-phase methodology. Detailed prompt templates and example responses can be found in Appendix A.

1) *Initial Prompt*: The initial prompt defines the task and outlines the context of the environment. The model receives the network topology (node coordinates and a directed adjacency list), each agent's origin–destination pair, and its departure time. We also provide a small set of navigation utility interfaces (Table II), which are thin wrappers around SUMO's built-in path planning and information routines.

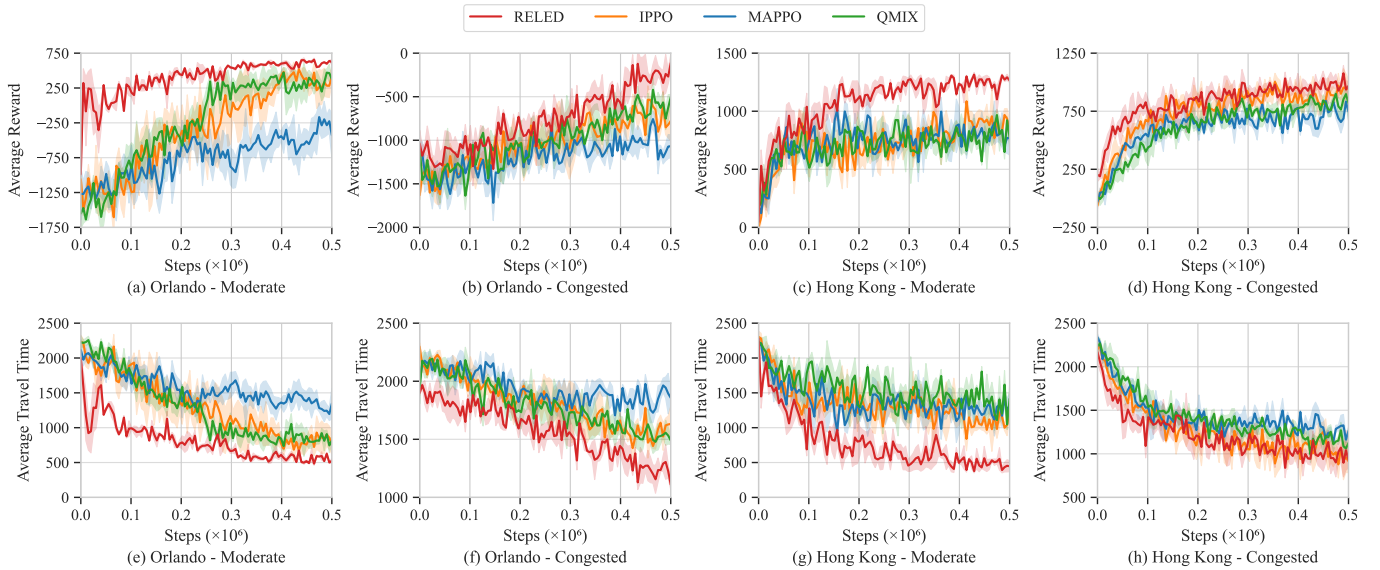


Fig. 4. Sample efficiency across cities and traffic regimes. (a–d) Average episode reward vs. interaction steps under Moderate and Congested settings. (e–h) Average travel time (measured in SUMO simulator time step) vs. steps. Shaded regions denote SD over 3 independent runs.

Finally, we specify output constraints with instructions to generate executable Python code format. In practice, the above prompt serves as the system prompt for the model. To maintain computational efficiency, the context window is limited to 10,000 tokens; earlier historical information is discarded as needed, while the system prompt is preserved throughout the interaction process. In addition, the initial prompt contains a simple command to trigger the instruction.

2) *Feedback Prompt*: Every $q = 10$ epochs, we compile a summary of recent agent exploration and expert-generated plans and deliver it to the LLM as a structured feedback prompt. The summary includes two diagnostics: (i) *reward volatility index (RVI)*: As defined in Eqn. 19, RVI identifies groups whose substituted expert policies produce the largest single-step reward decreases for agents running the current DRL policy, together with the associated states and triggering instructions; (ii) *the policy divergence index (PDI)*: As defined in Eqn. 24, PDI computed as the KL divergence between expert and learned policies at matched decision states, together with the associated states and triggering instructions. Using the combined information, the LLM revises the expert plans to mitigate cross-agent interference and reduce policy mismatch.

V. EVALUATION

In this section, we conduct extensive experiments to evaluate the effectiveness of our proposed method RELED.

A. Comparative Results

This section compare RELED with state-of-the-art MARL methods on sample efficiency and time efficiency.

1) *Sample Efficiency*: Figure 4 demonstrates that RELED consistently outperforms all baseline approaches in both sample efficiency and final performance across all four test scenarios. RELED achieves higher average episode rewards with fewer training steps. In the moderate traffic conditions of (a) Orlando and (c) Hong Kong, RELED rapidly reaches strong

performance levels and continues to improve steadily. In the more challenging congested environments of (b) Orlando and (d) Hong Kong, RELED maintains consistently higher average rewards per interaction step compared to all baseline methods. The complementary metrics average travel time presented in Figure 4 (e)–(h) provide further evidence of RELED’s superior performance. Average travel time directly reflects agent effectiveness in traffic management scenarios as a pure optimization objective, unaffected by reward shaping components. Time unit here is SUMO simulator time step, which is the standard temporal measurement in the SUMO traffic simulation environment. RELED reduces travel times more rapidly and converges to lower values than the baselines. This highlights RELED’s enhanced ability to efficiently translate sample experiences into effective policy improvements across diverse urban networks and traffic conditions.

2) *Time Efficiency*: Figure 5 illustrates the progression of average episode reward over training time. Despite the additional computational overhead introduced by LLM inference, RELED remains competitive in time efficiency. For Orlando – moderate, even with the maximum historical context, the average time for a single refinement inference is only 7.35 ± 4.86 seconds, which is small relative to the overall training duration. More detailed analysis of model inference time is presented in Section V-C. Across both urban environments and traffic conditions, RELED achieves higher average rewards earlier than all baselines, establishing an early performance advantage that is maintained throughout the training process. These results shows that the incremental inference cost is more than offset by the resulting more informative updates, yielding greater performance gains per unit time.

B. Sample Analysis of HPO Module

This section presents empirical evidence that our method effectively aligns agents with expert behavior and transitions

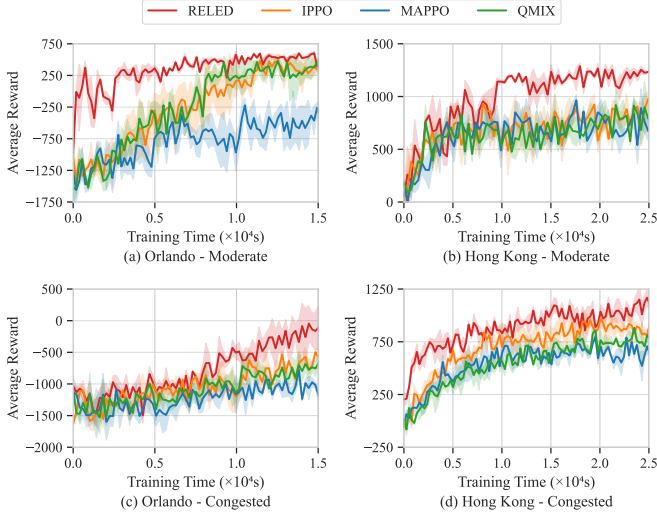


Fig. 5. Time efficiency in wall-clock training. Average episode reward over elapsed time for Orlando and Hong Kong (Moderate and Congested). Shaded regions indicate SD across 3 independent runs.

learning from imitation to self-driven optimization in HPO module.

1) *Trajectories Visualization*: Figure 6 visualizes expert demonstrations and agent-generated trajectories in Orlando - Moderate and Hong Kong - Moderate. The visualization presents samples from the final training phase, showing the last 10 episodes for both expert and agent trajectories. The data is log-transformed and normalized to $[0,1]$ to better represent the right-skewed distribution of traffic flow. In both environments, expert demonstrations reveal concentrated traffic flow along key corridor chains connecting origins and destinations. DRL agents successfully replicate these primary corridors while maintaining comparable flow intensities on critical segments, demonstrating effective policy transfer. Differences primarily emerge in secondary connections near intersections, where the agents redistribute some traffic flow. This reflects the inherent exploration-exploitation trade-off in reinforcement learning. Agents leverage expert demonstrations as a structural foundation while conducting exploration on peripheral edges to discover potentially improved routing decisions.

2) *Policy Alignment*: To quantify this alignment, Figure 7 reports the DTW distance between expert and agent trajectories across training, together with the average episode reward. DTW measures the similarity between temporal sequences by finding the optimal nonlinear alignment between them. Lower DTW denotes closer adherence to the expert policy. In both environment, DTW drops quickly in the early epochs and stays relatively low thereafter, while reward rises. In RELED, a lower DTW leads to a reduced weighting of expert samples in policy updates through the coefficient α , as expressed in Eqn. 30. The result shows that the agents gradually shifts learning toward self-generated experience, continuously improving performance even after achieving close alignment with expert behavior.

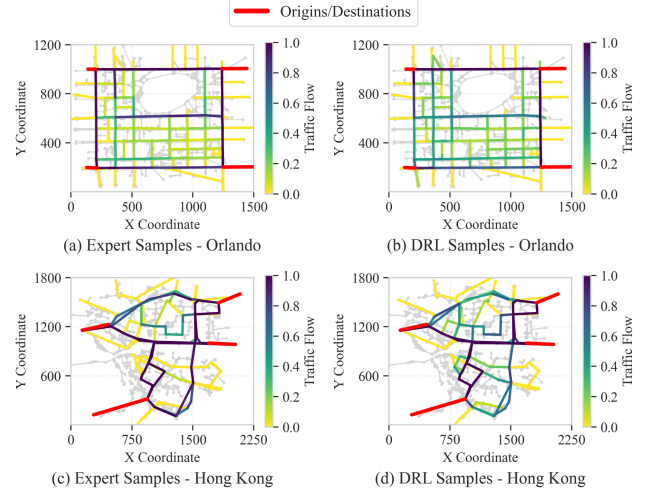


Fig. 6. Qualitative trajectory visualization comparing expert demonstrations (left) and agent-generated trajectories (right) in (a, b) Orlando - Moderate and (c, d) Hong Kong - Moderate.

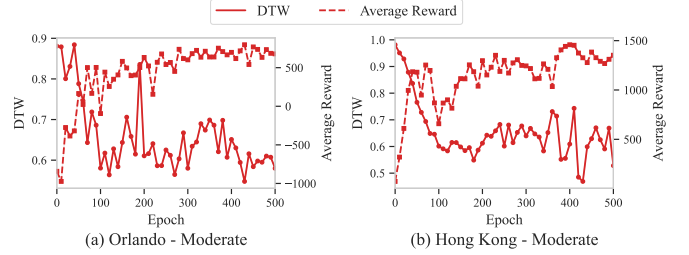


Fig. 7. DTW dynamics with training. DTW distance between expert and agent trajectories (left axis) alongside average episode reward (right axis) during training on (a) Orlando - Moderate and (b) Hong Kong - Moderate.

C. Demonstration Analysis of SED Module

In this section, we evaluate LLM-generated demonstrations across various refinement phases, model choices, and prompt designs. We also identify how the SED module contributes to high-quality demonstrations through its structured interfaces and iterative feedback mechanisms.

1) *LLM Model Choice*: Table III presents our evaluation of LLM-generated demonstrations across three refinement phases in Orlando - Moderate. For each phase, 100 expert trajectories (10 prompts \times 10 agents) are collected. We evaluated four widely-used LLMs of varying parameter sizes for their impact on the quality of expert demonstrations. GPT-4 Turbo achieves the highest rewards (752.84 ± 65.32 with full setting at 50 iterative refinements) but requires substantially longer inference times (19.55 ± 1.70 seconds). GPT-3.5 Turbo and Llama-3.1-70b show competitive performance (684.76 ± 57.34 and 708.46 ± 47.28 rewards) with faster inference (7.35 ± 4.86 and 6.86 ± 1.35 seconds), offering excellent balance between quality and computational efficiency. In our experiments, we selected GPT-3.5 Turbo for this balanced performance profile. Overall, larger models tend to deliver higher rewards and execution rates on average, whereas smaller models like Llama-3.1-8B are attractive with strict computational constraints (e.g.,

TABLE III
LLM-GENERATED DEMONSTRATIONS UNDER DIFFERENT SETTINGS AND REFINEMENT ITERATIONS

Model	IFs	Refn.	Initial			25 Refinements			DTW diff. \uparrow	50 Refinements			DTW diff. \uparrow
			Exec. (%) \uparrow	Reward \uparrow	Time \downarrow	Exec. (%) \uparrow	Reward \uparrow	Time \downarrow		Exec. (%) \uparrow	Reward \uparrow	Time \downarrow	
GPT-4 Turbo	\checkmark	\checkmark	92.00 \pm 16.00	687.92 \pm 78.01	14.36 \pm 2.49	77.00 \pm 10.05	724.63 \pm 72.15	19.83 \pm 2.43	1266.75 \pm 324.20	80.00 \pm 8.94	752.84 \pm 65.32	19.55 \pm 1.70	980.39 \pm 220.59
		\checkmark	74.00 \pm 36.66	643.10 \pm 42.68	13.28 \pm 2.01	79.00 \pm 8.31	685.27 \pm 38.42	19.55 \pm 1.28	928.82 \pm 156.94	74.00 \pm 20.10	710.38 \pm 35.76	16.58 \pm 0.98	949.72 \pm 197.49
		\checkmark	50.00 \pm 36.06	447.96 \pm 80.57	35.55 \pm 4.87	34.00 \pm 39.04	495.74 \pm 73.26	53.61 \pm 9.93	1066.75 \pm 222.53	52.00 \pm 24.41	543.82 \pm 67.38	52.54 \pm 7.30	943.87 \pm 197.46
GPT-3.5 Turbo	\checkmark	\checkmark	92.00 \pm 4.00	608.42 \pm 68.53	5.45 \pm 1.13	92.00 \pm 16.00	645.32 \pm 62.46	7.67 \pm 4.57	1012.38 \pm 187.65	85.00 \pm 6.71	684.76 \pm 57.34	7.35 \pm 4.86	1040.50 \pm 174.15
		\checkmark	86.00 \pm 22.00	556.74 \pm 58.32	4.30 \pm 1.89	87.00 \pm 15.52	584.38 \pm 52.76	7.01 \pm 4.56	817.46 \pm 142.38	80.00 \pm 7.75	615.82 \pm 47.28	6.88 \pm 0.90	836.66 \pm 206.82
		\checkmark	54.00 \pm 18.00	584.86 \pm 72.48	6.51 \pm 1.63	45.00 \pm 15.00	386.35 \pm 68.24	12.60 \pm 5.01	812.54 \pm 174.28	67.00 \pm 11.87	428.92 \pm 62.16	12.46 \pm 1.40	1222.51 \pm 209.50
Llama-3.1-70b	\checkmark	\checkmark	83.00 \pm 14.18	625.38 \pm 57.46	5.33 \pm 0.41	86.00 \pm 8.00	668.75 \pm 52.38	6.46 \pm 0.58	1112.64 \pm 212.38	74.00 \pm 6.63	708.46 \pm 47.28	6.86 \pm 1.35	1017.66 \pm 117.68
		\checkmark	86.00 \pm 8.00	586.42 \pm 48.26	6.02 \pm 0.58	91.00 \pm 9.43	617.38 \pm 42.76	4.47 \pm 0.38	863.54 \pm 152.63	70.00 \pm 14.83	658.74 \pm 38.42	4.78 \pm 0.93	941.86 \pm 104.88
		\checkmark	59.00 \pm 17.58	684.62 \pm 53.84	6.60 \pm 1.66	64.00 \pm 36.66	527.35 \pm 47.62	9.76 \pm 1.43	916.78 \pm 163.42	58.00 \pm 20.88	568.42 \pm 42.16	14.32 \pm 1.23	1098.65 \pm 220.97
Llama-3.1-8b	\checkmark	\checkmark	68.00 \pm 8.93	613.98 \pm 48.55	4.79 \pm 0.31	68.00 \pm 16.61	591.46 \pm 56.19	4.89 \pm 0.40	762.59 \pm 294.72	69.00 \pm 25.08	595.67 \pm 43.05	4.64 \pm 0.60	846.96 \pm 167.29
		\checkmark	79.00 \pm 8.31	557.56 \pm 89.99	4.39 \pm 0.22	63.00 \pm 24.52	512.08 \pm 35.09	3.59 \pm 0.21	781.50 \pm 100.24	69.00 \pm 23.00	526.66 \pm 87.94	3.70 \pm 0.81	710.23 \pm 198.95
		\checkmark	47.00 \pm 8.81	519.18 \pm 38.02	2.79 \pm 0.60	60.00 \pm 48.99	607.84 \pm 64.84	11.49 \pm 4.74	749.58 \pm 163.29	42.00 \pm 20.40	510.51 \pm 94.75	10.06 \pm 5.59	714.42 \pm 153.39

We report execution rate (Exec.), average reward of the demonstration trajectory (Reward), LLM inference time (Time), DTW differences relative to initial trajectories (DTW diff.), existence of navigation interface (IFs) and interactive refinement mechanism (Refn.) as mean \pm SD over runs. The \uparrow indicates higher value is better, while the \downarrow indicates lower value is better.

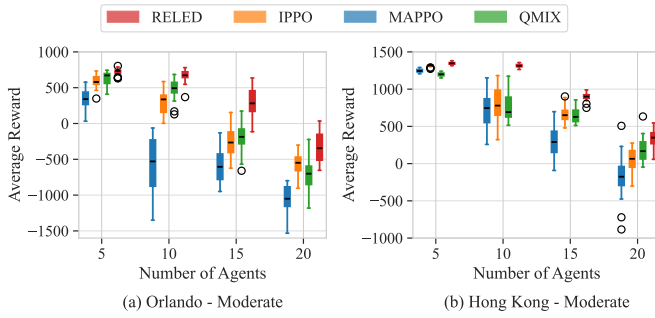


Fig. 8. Performance under different numbers of agents, reported by the distribution of average episode rewards over 20 test runs on (a) Orlando - Moderate and (b) Hong Kong - Moderate.

4.64 \pm 0.60 s at 50 refinements).

2) *Impact of Prompt Design*: The ablation study of prompt design in Table III highlights how different components impact the quality of expert demonstrations. In the absence of refinement feedback (Refn.), DTW difference decreases substantially. For example, with GPT-3.5 Turbo (50 refinements), the DTW difference drops to 836.66 \pm 206.82 (vs. 1040.50 \pm 174.15 with both components). This pronounced reduction in DTW indicates that the feedback mechanism, by iteratively refining model outputs, enhances the consistency between expert demonstrations and reinforcement learning samples. Providing LLMs with predefined navigation interfaces (IFs) generally improves performance, compared to letting LLM generate complete node-connected routes. This advantage is more pronounced in higher refinement iterations where there is more context. At 50 refinements, without IFs, execution rates decline dramatically while rewards decrease substantially — GPT-3.5 Turbo’s execution rate falls from 85% to 67.00% and rewards drop from 684.76 \pm 57.34 to 428.92 \pm 62.16, with inference time also increasing significantly from 7.35 \pm 4.86 to 12.46 \pm 1.40.

D. Scalability

This section evaluates scalability of RELED. Figure 8 evaluates performance as the number of agents increases on the Orlando - Moderate and Hong Kong - Moderate scenarios. Each box shows the distribution of the average episode reward

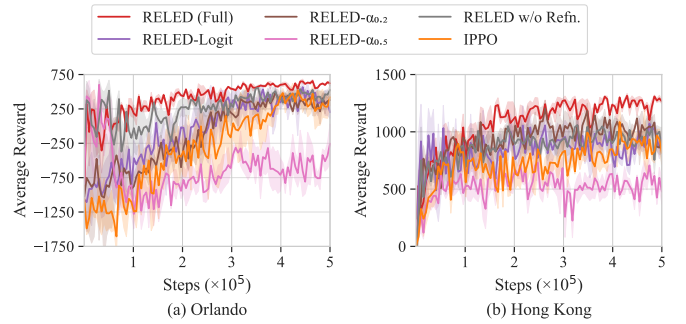


Fig. 9. Ablation study on training dynamics comparing average reward over steps for RELED (Full) and alternative variants on (a) Orlando - Moderate and (b) Hong Kong - Moderate; shaded regions indicate SD across 3 runs.

after 20 test episodes of these models. Across all scales, RELED maintains the highest median reward compared to IPPO, MAPPO, and QMIX. As the population grows from 5 to 20, all methods experience some degradation, but RELED’s decline is more gradual. These results indicate that the proposed approach scales more robustly with the number of agents.

E. Ablation Study

This section evaluates the impact of each module of RELED performance. We compare the following variants: (i) *RELED (Full)*; (ii) *RELED w/o Refn.*: Implemented without iterative prompt refinement; (iii) *RELED- $\alpha_{0.2}$* ; (iv) *RELED- $\alpha_{0.5}$* : Using fixed expert weights $\alpha = 0.2$ and $\alpha = 0.5$ instead of dynamic α calculated by DTW distance; (v) *RELED-Logit*: Replacing LLM-generated demonstrations with a logit-based softmax over shortest-path costs to prioritize lower-cost routes; and (vi) *IPPO*.

1) *Training Performance*: Figure 9 presents our ablation study comparing these RELED variants in training performance. RELED (Full) achieves the fastest performance gains and the highest episode reward in both environments with reduced variance throughout training. This suggests that iteratively updating expert demonstrations can help improve training efficiency. Removing iterative refinement (RELED w/o Refn.) slows late-stage reward improvement. Using fixed expert weights RELED- $\alpha_{0.2}$ /RELED- $\alpha_{0.5}$ further reduces sample

TABLE IV
CONVERGENCE PERFORMANCE OF DIFFERENT VARIANTS.

Method	Orlando - Moderate		Hong Kong- Moderate	
	Reward \uparrow Raw (Drop %)	Time \downarrow Raw (Rise %)	Reward \uparrow Raw (Drop %)	Time \downarrow Raw (Rise %)
RELED (Full)	782.34 \pm 2.81 (-)	505.21 \pm 21.13 (-)	1492.91 \pm 31.46 (-)	519.97 \pm 18.75 (-)
RELED w/o Refn.	627.51 \pm 90.20 (-19.79%)	553.68 \pm 24.57 (+9.59%)	1102.57 \pm 101.60 (-26.15%)	655.22 \pm 35.08 (+26.01%)
RELED- $\alpha_{0.2}$	513.13 \pm 115.46 (-34.41%)	623.31 \pm 31.06 (+23.38%)	1143.70 \pm 121.82 (-23.39%)	519.97 \pm 19.94 (+12.65%)
RELED- $\alpha_{0.5}$	-75.51 \pm 359.01 (-109.65%)	1690.39 \pm 63.25 (+234.59%)	688.90 \pm 37.48 (-53.86%)	1708.52 \pm 55.19 (+228.58%)
RELED-Logit	443.02 \pm 63.78 (-43.37%)	834.79 \pm 49.16 (+65.24%)	1060.71 \pm 70.87 (-28.95%)	770.90 \pm 41.58 (+48.26%)
IPPO	481.58 \pm 135.40 (-38.44%)	682.86 \pm 39.22 (+35.16%)	865.21 \pm 35.76 (-42.05%)	1241.49 \pm 81.22 (+138.76%)

Reported as mean \pm SD over runs. The \uparrow indicates higher value is better, while the \downarrow indicates lower value is better.

efficiency; an overly large expert weight ($\alpha = 0.5$) leads to insufficient exploration and a sharp late-stage reward drop, especially in (a) Orlando. This validates the advantage of DTW-based adaptive weighting for scheduling the transition from imitation to exploration. RELED-Logit shows higher early sample efficiency but overall lags behind RELED (Full), suggesting that LLM-derived trajectories provide higher-quality, more context-aware guidance than logit-stochastic routing.

2) *Converged Performance*: Table IV reports the performance of the convergent model of each method, tested over 20 independent runs. RELED (Full) achieves the highest rewards with minimal travel times in both networks. RELED-Logit attains 43.37% lower reward and 65.24% longer travel time than RELED (Full) in Orlando, with similar gaps in Hong Kong. Fixed-weight configurations show severe performance degradation, with $\alpha = 0.5$ yielding negative rewards in Orlando (-109.65%) and substantial drops in Hong Kong (-53.86%), while $\alpha = 0.2$ also degrades performance in both environments (-34.41% and -23.39% lower rewards). Removing refinement (RELED w/o Refn.) results in modest performance drops across both environments (-19.79% and -26.15% lower rewards). These results quantitatively demonstrate the impact of each RELED component on overall converged performance.

VI. CONCLUSION

We propose RELED, a scalable MARL framework that effectively integrates LLM-driven expert demonstrations with autonomous agent exploration. The SED module leverages theoretical non-stationarity bounds to guide the generation and refinement of high-quality expert trajectories, improving cumulative rewards and training stability. Furthermore, the HPO module adaptively balances learning from both expert- and agent-generated experiences, accelerating policy convergence and enhancing generalization for each agent. Experimental results show that RELED achieves better sample and time efficiency across scenarios, maintains performance advantages with increasing agent numbers in practical traffic management applications. As a potential future direction, we are looking forward to extending our method to improve the performance of various applications such as robot control [68]–[73] and autonomous driving [23], [74]–[76].

REFERENCES

- [1] L. U. Khan, I. Yaqoob, N. H. Tran, S. A. Kazmi, T. N. Dang, and C. S. Hong, "Edge-computing-enabled smart cities: A comprehensive survey," *IEEE Internet Things J.*, vol. 7, no. 10, pp. 10 200–10 232, 2020.
- [2] T. Duan, Z. Zhang, S. Guo, Y. Zhao, Z. Lin, Z. Fang, Y. Liu, D. Luan, D. Huang, H. Cui *et al.*, "Sample efficient experience replay in non-stationary environments," *arXiv preprint arXiv:2509.15032*, 2025.
- [3] Z. Lin, Z. Chen, Z. Fang, X. Chen, X. Wang, and Y. Gao, "Fedsn: A federated learning framework over heterogeneous leo satellite networks," *IEEE Transactions on Mobile Computing*, 2024.
- [4] H. Yuan, Z. Chen, Z. Lin, J. Peng, Z. Fang, Y. Zhong, Z. Song, and Y. Gao, "SatSense: Multi-Satellite Collaborative Framework for Spectrum Sensing," *IEEE Trans. Cogn. Commun. Netw.*, 2025.
- [5] Y. Tang, Z. Chen, A. Li, T. Zheng, Z. Lin, J. Xu, P. Lv, Z. Sun, and Y. Gao, "MERIT: Multimodal Wearable Vital Sign Waveform Monitoring," *arXiv preprint arXiv:2410.00392*, 2024.
- [6] J. Peng, Z. Chen, Z. Lin, H. Yuan, Z. Fang, L. Bao, Z. Song, Y. Li, J. Ren, and Y. Gao, "SUMS: Sniffing Unknown Multiband Signals under Low Sampling Rates," *IEEE Trans. Mobile Comput.*, 2024.
- [7] Z. Fang, Z. Lin, S. Hu, Y. Tao, Y. Deng, X. Chen, and Y. Fang, "Dynamic uncertainty-aware multimodal fusion for outdoor health monitoring," *arXiv preprint arXiv:2508.09085*, 2025.
- [8] H. Yuan, Z. Chen, Z. Lin, J. Peng, Y. Zhong, X. Hu, S. Xue, W. Li, and Y. Gao, "Constructing 4D Radio Map in LEO Satellite Networks with Limited Samples," *IEEE INFOCOM*, 2025.
- [9] Z. Lin, G. Qu, W. Wei, X. Chen, and K. K. Leung, "Adaptsfl: Adaptive Split Federated Learning in Resource-Constrained Edge Networks," *IEEE Trans. Netw.*, 2024.
- [10] H. Yuan, Z. Chen, Z. Lin, J. Peng, Z. Fang, Y. Zhong, Z. Song, X. Wang, and Y. Gao, "Graph Learning for Multi-Satellite Based Spectrum Sensing," in *Proc. IEEE Int. Conf. Commun. Technol. (ICCT)*, 2023, pp. 1112–1116.
- [11] Z. Lin, Y. Zhang, Z. Chen, Z. Fang, C. Wu, X. Chen, Y. Gao, and J. Luo, "LEO-Split: A Semi-Supervised Split Learning Framework over LEO Satellite Networks," *IEEE Trans. Mobile Comput.*, 2025.
- [12] Z. Ning and L. Xie, "A survey on multi-agent reinforcement learning and its application," *J. Autom. Intell.*, vol. 3, no. 2, pp. 73–91, 2024.
- [13] T. Duan, Z. Zhang, S. Guo, D. Huang, Y. Zhao, Z. Lin, Z. Fang, D. Luan, H. Cui, and Y. Cui, "Lead: A highly efficient and scalable llm-empowered expert demonstrations framework for multi-agent reinforcement learning," *arXiv preprint arXiv:2509.14680*, 2025.
- [14] D. Chen, K. Zhang, Y. Wang, X. Yin, Z. Li, and D. Filev, "Communication-efficient decentralized multi-agent reinforcement learning for cooperative adaptive cruise control," *IEEE Trans. Intell. Veh.*, 2024.
- [15] Y. Xiao, W. Tan, and C. Amato, "Asynchronous actor-critic for multi-agent reinforcement learning," *Proc. NeurIPS*, vol. 35, pp. 4385–4400, 2022.
- [16] Z. Gao, L. Yang, and Y. Dai, "Large-scale computation offloading using a multi-agent reinforcement learning in heterogeneous multi-access edge computing," *IEEE Trans. Mobile Comput.*, vol. 22, no. 6, pp. 3425–3443, 2022.
- [17] Z. Lin, W. Wei, Z. Chen, C.-T. Lam, X. Chen, Y. Gao, and J. Luo, "Hierarchical Split Federated Learning: Convergence Analysis and System Optimization," *IEEE Trans. Mobile Comput.*, 2025.
- [18] W. Wei, Z. Lin, X. Liu, H. Du, D. Niyato, and X. Chen, "Optimizing split federated learning with unstable client participation," *arXiv preprint arXiv:2509.17398*, 2025.
- [19] Z. Lin, G. Qu, X. Chen, and K. Huang, "Split Learning in 6G Edge Networks," *IEEE Wirel. Commun.*, 2024.
- [20] T. Du, X. Gui, and T. Sheng, "Multi-agent Deep Reinforcement Learning-Based Hierarchical Scheduling in Heterogeneous UAVs Enabled Vehicular Networks," *IEEE Internet Things J.*, 2025.
- [21] Z. Gao, J. Fu, Z. Jing, Y. Dai, and L. Yang, "MOIPC-MAAC: Communication-assisted multiobjective Marl for trajectory planning and task offloading in multi-UAV-assisted MEC," *IEEE Internet Things J.*, vol. 11, no. 10, pp. 18 483–18 502, 2024.
- [22] Y. Zhang, Z. Yu, J. Zhang, L. Wang, T. H. Luan, B. Guo, and C. Yuen, "Learning decentralized traffic signal controllers with multi-agent graph reinforcement learning," *IEEE Trans. Mobile Comput.*, vol. 23, no. 6, pp. 7180–7195, 2023.
- [23] Y. Wang, T. Xu, X. Niu, C. Tan, E. Chen, and H. Xiong, "STMARL: A spatio-temporal multi-agent reinforcement learning approach for cooperative traffic light control," *IEEE Trans. Mobile Comput.*, vol. 21, no. 6, pp. 2228–2242, 2020.

- [24] K. Zhang, Z. Yang, and T. Başar, “Multi-agent reinforcement learning: A selective overview of theories and algorithms,” *Handb. Reinforcement Learn. Control*, pp. 321–384, 2021.
- [25] C. Yu, A. Velu, E. Vinititsky, J. Gao, Y. Wang, A. Bayen, and Y. Wu, “The surprising effectiveness of ppo in cooperative multi-agent games,” *Proc. NeurIPS*, vol. 35, pp. 24611–24624, 2022.
- [26] C. S. De Witt, T. Gupta, D. Makoviichuk, V. Makoviychuk, P. H. Torr, M. Sun, and S. Whiteson, “Is independent learning all you need in the starcraft multi-agent challenge?” *arXiv preprint arXiv:2011.09533*, 2020.
- [27] T. Rashid, M. Samvelyan, C. S. De Witt, G. Farquhar, J. Foerster, and S. Whiteson, “Monotonic value function factorisation for deep multi-agent reinforcement learning,” *J. Mach. Learn. Res.*, vol. 21, no. 178, pp. 1–51, 2020.
- [28] J. Ackermann, V. Gabler, T. Osa, and M. Sugiyama, “Reducing overestimation bias in multi-agent domains using double centralized critics,” *arXiv preprint arXiv:1910.01465*, 2019.
- [29] L. Wang, Z. Yang, and Z. Wang, “Breaking the curse of many agents: Provable mean embedding Q-iteration for mean-field reinforcement learning,” in *Proc. ICML*, 2020, pp. 10092–10103.
- [30] P. Hernandez-Leal, M. Kaisers, T. Baarslag, and E. M. De Cote, “A survey of learning in multiagent environments: Dealing with non-stationarity,” *arXiv preprint arXiv:1707.09183*, 2017.
- [31] B. Pan, J. Lu, K. Wang, L. Zheng, Z. Wen, Y. Feng, M. Zhu, and W. Chen, “AgentCoord: Visually exploring coordination strategy for llm-based multi-agent collaboration,” *arXiv preprint arXiv:2404.11943*, 2024.
- [32] Z. Lin, G. Qu, Q. Chen, X. Chen, Z. Chen, and K. Huang, “Pushing Large Language Models to the 6G Edge: Vision, Challenges, and Opportunities,” *IEEE Communication Magazine*, 2023.
- [33] Z. Wang, S. Cai, G. Chen, A. Liu, X. Ma, and Y. Liang, “Describe, explain, plan and select: Interactive planning with large language models enables open-world multi-task agents,” *arXiv preprint arXiv:2302.01560*, 2023.
- [34] Z. Lin, Y. Zhang, Z. Chen, Z. Fang, X. Chen, P. Vepakomma, W. Ni, J. Luo, and Y. Gao, “HSplitLoRA: A Heterogeneous Split Parameter-Efficient Fine-Tuning Framework for Large Language Models,” *arXiv preprint arXiv:2505.02795*, 2025.
- [35] Z. Fang, Z. Lin, Z. Chen, X. Chen, Y. Gao, and Y. Fang, “Automated Federated Pipeline for Parameter-Efficient Fine-Tuning of Large Language Models,” *arXiv preprint arXiv:2404.06448*, 2024.
- [36] Z. Lin, X. Hu, Y. Zhang, Z. Chen, Z. Fang, X. Chen, A. Li, P. Vepakomma, and Y. Gao, “SplitLoRA: A Split Parameter-Efficient Fine-Tuning Framework for Large Language Models,” *arXiv preprint arXiv:2407.00952*, 2024.
- [37] M. Ahn, A. Brohan, N. Brown, Y. Chebotar, O. Cortes, B. David, C. Finn, C. Fu, K. Gopalakrishnan, K. Hausman *et al.*, “Do as i can, not as i say: Grounding language in robotic affordances,” *arXiv preprint arXiv:2204.01691*, 2022.
- [38] J. Liu, D. Shen, Y. Zhang, W. B. Dolan, L. Carin, and W. Chen, “What makes good in-context examples for GPT-3?” in *Proc. DeeLLO*, 2022, pp. 100–114.
- [39] G. Qu and N. Li, “Exploiting fast decaying and locality in multi-agent mdp with tree dependence structure,” in *Proc. CDC*, 2019, pp. 6479–6486.
- [40] M. Yuan, H. Huang, Z. Li, C. Zhang, F. Pei, and W. Gu, “A multi-agent double deep-Q-network based on state machine and event stream for flexible job shop scheduling problem,” *Adv. Eng. Inform.*, vol. 58, p. 102230, 2023.
- [41] P. Sunehag, G. Lever, A. Gruslys, W. M. Czarnecki, V. Zambaldi, M. Jaderberg, M. Lanctot, N. Sonnerat, J. Z. Leibo, K. Tuyls *et al.*, “Value-decomposition networks for cooperative multi-agent learning,” *arXiv preprint arXiv:1706.05296*, 2017.
- [42] Y. Bai, Y. Lv, and J. Zhang, “Smart mobile robot fleet management based on hierarchical multi-agent deep Q network towards intelligent manufacturing,” *Eng. Appl. Artif. Intell.*, vol. 124, p. 106534, 2023.
- [43] S. Iqbal and F. Sha, “Actor-attention-critic for multi-agent reinforcement learning,” in *Proc. ICML*. PMLR, 2019, pp. 2961–2970.
- [44] J. G. Kuba, R. Chen, M. Wen, Y. Wen, F. Sun, J. Wang, and Y. Yang, “Trust region policy optimisation in multi-agent reinforcement learning,” *arXiv preprint arXiv:2109.11251*, 2021.
- [45] S. Gu, J. G. Kuba, M. Wen, R. Chen, Z. Wang, Z. Tian, J. Wang, A. Knoll, and Y. Yang, “Multi-agent constrained policy optimisation,” *arXiv preprint arXiv:2110.02793*, 2021.
- [46] R. Lowe, Y. I. Wu, A. Tamar, J. Harb, O. Pieter Abbeel, and I. Mordatch, “Multi-agent actor-critic for mixed cooperative-competitive environments,” *Proc. NeurIPS*, vol. 30, 2017.
- [47] M. Pternea, P. Singh, A. Chakraborty, Y. Oruganti, M. Milletari, S. Bapat, and K. Jiang, “The rl/llm taxonomy tree: Reviewing synergies between reinforcement learning and large language models,” *J. Artif. Intell. Res.*, vol. 80, pp. 1525–1573, 2024.
- [48] Y. Feng, Y. Wang, J. Liu, S. Zheng, and Z. Lu, “Llama rider: Spurring large language models to explore the open world,” *arXiv preprint arXiv:2310.08922*, 2023.
- [49] Y. Wu, S. Y. Min, S. Prabhume, Y. Bisk, R. R. Salakhutdinov, A. Azaria, T. M. Mitchell, and Y. Li, “Spring: Studying papers and reasoning to play games,” *Proc. NeurIPS*, vol. 36, pp. 22383–22687, 2023.
- [50] J. Qiu, M. Xu, W. Han, S. Moon, and D. Zhao, “Embodied executable policy learning with language-based scene summarization,” *arXiv preprint arXiv:2306.05696*, 2023.
- [51] C. L. Shek, X. Wu, W. A. Suttle, C. Busart, E. Zaroukian, D. Manocha, P. Tokekar, and A. S. Bedi, “Lancar: Leveraging language for context-aware robot locomotion in unstructured environments,” in *Proc. IROS*. IEEE, 2024, pp. 9612–9619.
- [52] A. Szot, M. Schwarzer, H. Agrawal, B. Mazouze, R. Metcalfe, W. Talbott, N. Mackraz, R. D. Hjelm, and A. T. Toshev, “Large language models as generalizable policies for embodied tasks,” in *Proc. ICLR*, 2023.
- [53] S. Sontakke, J. Zhang, S. Arnold, K. Pertsch, E. Biyik, D. Sadigh, C. Finn, and L. Itti, “Roboclip: One demonstration is enough to learn robot policies,” *Proc. NeurIPS*, vol. 36, pp. 55681–55693, 2023.
- [54] E. Triantafyllidis, F. Christianos, and Z. Li, “Intrinsic language-guided exploration for complex long-horizon robotic manipulation tasks,” in *Proc. ICRA*, 2024, pp. 7493–7500.
- [55] M. Kwon, S. M. Xie, K. Bullard, and D. Sadigh, “Reward design with language models,” *arXiv preprint arXiv:2303.00001*, 2023.
- [56] J. Song, Z. Zhou, J. Liu, C. Fang, Z. Shu, and L. Ma, “Self-refined large language model as automated reward function designer for deep reinforcement learning in robotics,” *arXiv preprint arXiv:2309.06687*, 2023.
- [57] G. Zhu, R. Zhou, W. Ji, and S. Zhao, “LAMARL: LLM-Aided Multi-Agent Reinforcement Learning for Cooperative Policy Generation,” *IEEE Robot. Autom. Lett.*, 2025.
- [58] S. Liu, Z. Liang, X. Lyu, and C. Amato, “Llm collaboration with multi-agent reinforcement learning,” *arXiv preprint arXiv:2508.04652*, 2025.
- [59] F. A. Oliehoek, C. Amato *et al.*, *A concise introduction to decentralized POMDPs*. Springer, 2016, vol. 1.
- [60] J. Achiam, D. Held, A. Tamar, and P. Abbeel, “Constrained policy optimization,” in *Proc. ICML*, 2017, pp. 22–31.
- [61] P. Yu, M. Mishra, A. Koppel, C. Busart, P. Narayan, D. Manocha, A. Bedi, and P. Tokekar, “Beyond joint demonstrations: Personalized expert guidance for efficient multi-agent reinforcement learning,” *arXiv preprint arXiv:2403.08936*, 2024.
- [62] L. Weiwei, J. Wei, L. Shanqi, R. Yudi, Z. Kexin, Y. Jiang, and L. Yong, “Expert demonstrations guide reward decomposition for multi-agent cooperation,” *Neural Comput. Appl.*, vol. 35, no. 27, pp. 19847–19863, 2023.
- [63] H. Wei, Z. Zhang, S. He, T. Xia, S. Pan, and F. Liu, “Plangen-llms: A modern survey of llm planning capabilities,” *arXiv preprint arXiv:2502.11221*, 2025.
- [64] Z. Zhou, J. Song, K. Yao, Z. Shu, and L. Ma, “Isr-llm: Iterative self-refined large language model for long-horizon sequential task planning,” in *Proc. ICRA*, 2024, pp. 2081–2088.
- [65] J. Schulman, F. Wolski, P. Dhariwal, A. Radford, and O. Klimov, “Proximal policy optimization algorithms,” *arXiv preprint arXiv:1707.06347*, 2017.
- [66] OpenStreetMap contributors, “Planet dump retrieved from <https://planet.osm.org>,” <https://www.openstreetmap.org>, 2017.
- [67] P. A. Lopez, M. Behrisch, L. Bieker-Walz, J. Erdmann, Y.-P. Flötteröd, R. Hilbrich, L. Lücken, J. Rummel, P. Wagner, and E. Wießner, “Microscopic traffic simulation using sumo,” in *Proc. ITSC*, 2018, pp. 2575–2582.
- [68] Z. Zhang, T. Duan, Z. Lin, D. Huang, Z. Fang, Z. Sun, L. Xiong, H. Liang, H. Cui, Y. Cui *et al.*, “Robust deep reinforcement learning in robotics via adaptive gradient-masked adversarial attacks,” *arXiv preprint arXiv:2503.20844*, 2025.
- [69] Z. Zhang, T. Duan, Z. Lin, D. Huang, Z. Fang, Z. Sun, L. Xiong, H. Liang, H. Cui, and Y. Cui, “State-aware perturbation optimization for robust deep reinforcement learning,” *arXiv preprint arXiv:2503.20613*, 2025.
- [70] T. Duan, Z. Zhang, Z. Lin, Y. Gao, L. Xiong, Y. Cui, H. Liang, X. Chen, H. Cui, and D. Huang, “Rethinking adversarial attacks in reinforcement learning from policy distribution perspective,” in *Proc. ICASSP*, 2025, pp. 1–5.

- [71] J. Wang, Z. Sun, X. Guan, T. Shen, Z. Zhang, T. Duan, D. Huang, S. Zhao, and H. Cui, "Agrnav: Efficient and energy-saving autonomous navigation for air-ground robots in occlusion-prone environments," in *Proc. ICRA*, 2024, pp. 11 133–11 139.
- [72] Z. Lin, Z. Chen, X. Chen, W. Ni, and Y. Gao, "HASFL: Heterogeneity-Aware Split Federated Learning over Edge Computing Systems," *arXiv preprint arXiv:2506.08426*, 2025.
- [73] J. Wang, Z. Sun, X. Guan, T. Shen, D. Huang, Z. Zhang, T. Duan, F. Liu, and H. Cui, "He-nav: A high-performance and efficient navigation system for aerial-ground robots in cluttered environments," *IEEE Robot. Autom. Lett.*, 2024.
- [74] Z. Lin, L. Wang, J. Ding, B. Tan, and S. Jin, "Channel Power Gain Estimation for Terahertz Vehicle-to-Infrastructure Networks," *IEEE Commun. Lett.*, vol. 27, no. 1, pp. 155–159, 2022.
- [75] Z. Lin, L. Wang, J. Ding, Y. Xu, and B. Tan, "Tracking and transmission design in terahertz v2i networks," *IEEE Transactions on Wireless Communications*, vol. 22, no. 6, pp. 3586–3598, 2022.
- [76] Z. Fang, Z. Lin, S. Hu, H. Cao, Y. Deng, X. Chen, and Y. Fang, "IC3M: In-Car Multimodal Multi-Object Monitoring for Abnormal Status of Both Driver and Passengers," *arXiv preprint arXiv:2410.02592*, 2024.