KNOWING WHAT NOT TO DO: LEVERAGE LANGUAGE MODEL INSIGHTS FOR ACTION SPACE PRUNING IN MULTI-AGENT REINFORCEMENT LEARNING

Anonymous authors

Paper under double-blind review

ABSTRACT

Multi-agent reinforcement learning (MARL) is employed to develop autonomous agents that can learn to adopt cooperative or competitive strategies within complex environments. However, the linear increase in the number of agents leads to a combinatorial explosion of the action space, which may result in algorithmic instability, difficulty in convergence, or entrapment in local optima. While researchers have designed a variety of effective algorithms to compress the action space, these methods also introduce new challenges, such as the need for manually designed prior knowledge or reliance on the structure of the problem, which diminishes the applicability of these techniques. In this paper, we introduce Evolutionary action SPAce Reduction with Knowledge (eSpark), an exploration function generation framework driven by large language models (LLMs) to boost exploration and prune unnecessary actions in MARL. Using just a basic prompt that outlines the overall task and setting, eSpark is capable of generating exploration functions in a zeroshot manner, identifying and pruning redundant or irrelevant state-action pairs, and then achieving autonomous improvement from policy feedback. In reinforcement learning tasks involving inventory management and traffic light control encompassing a total of 15 scenarios, eSpark consistently outperforms the combined MARL algorithm in all scenarios, achieving an average performance gain of 34.4% and 9.9% in the two types of tasks respectively. Additionally, eSpark has proven to be capable of managing situations with a large number of agents, securing a 29.7% improvement in scalability challenges that featured over 500 agents. The code can be found in https://anonymous.4open.science/r/OCDH-ODF8/.

006

008 009 010

011

013

014

015

016

017

018

019

021

023

024

025

026

027

028

029

031

1 INTRODUCTION

Multi-agent reinforcement learning (MARL) has emerged as a powerful paradigm for solving complex 037 and dynamic problems that involve multiple decision-makers Zhang et al. (2021); Wang et al. (2021). However, the intricacies of agent interplay and the exponential expansion of state and action spaces render the solution of MARL problems difficult. Researchers have proposed the Centralized Training 040 with Decentralized Execution (CTDE) framework Oliehoek et al. (2008) and parameter sharing 041 methods, decomposing the value or policy functions of a multi-agent system into individual agents 042 and sharing model parameters among all agents. As experimentally verified by many of the most 043 prominent MARL algorithms such as Multi-agent PPO (MAPPO) Yu et al. (2022), QMIX Rashid 044 et al. (2020), QPLEX Wang et al. (2021) or QTRAN Son et al. (2019), these methodologies have been demonstrated to be robust strategies for surmounting the challenges posed by MARL. MARL methods based on parameter sharing and CTDE have achieved notable success in a variety of well-established 046 tasks, including StarCraft Multi-Agent Challenge (SMAC) Li et al. (2023); Wang et al. (2020), the 047 Multi-Agent Particle Environment (MPE) Lowe et al. (2017), and Simulation of Urban MObility 048 (SUMO) Wei et al. (2019); Lu et al. (2023). 049

Despite the great success of parameter-sharing CTDE methods, their practicality dwindles in real world tasks involving large number of agents, such as large-scale traffic signal control Mousavi
 et al. (2017), wireless communication networks Zocca (2019), and humanitarian assistance and
 disaster response Meier (2015), where centralized training becomes impractical due to large problem
 scale Munir et al. (2021). Fully Decentralized Training and Execution (DTDE) methods, such as

Independent PPO (IPPO) de Witt et al. (2020), offer a scalable solution where resource consumption
does not escalate drastically with an increase in the number of agents. However, due to the lack
of consideration for agent interactions, they often struggle to find optimal solutions and fall into
local optima. Current strategies for addressing large-scale MARL tasks involve introducing taskspecific structures to model agent interactions or dividing agents into smaller, independently trained
groups Ying et al. (2023); Chen et al. (2020). These methods, however, are constrained by their
dependence on task-related structuring, limiting their applicability to a narrow range of problems.

061 Additionally, the "curse of dimensionality" presents a significant challenge in multi-agent systems Hao 062 et al. (2022b;a), where agents are required to navigate through an expansive action space saturated with 063 numerous actions that are either irrelevant or markedly suboptimal (relative to states). While humans 064 can deftly employ contextual cues and prior knowledge to sidestep such challenges, MARL algorithms typically engage in the exploration of superfluous and extraneous suboptimal actions Zahavy et al. 065 (2018). Besides, prevailing parameter sharing can exacerbate this exploratory dilemma, as will be 066 elucidated in Proposition 2. The issue occurs primarily because agents with shared parameters often 067 prefer suboptimal policies that present short-term advantages, rather than exploring policies that may 068 potentially deliver higher long-term returns. 069

Exploration is crucial for overcoming local optima, as it encourages agents to discover potentially 071 better states thus refining their policies. While single-agent exploration techniques like the Upper Confidence Bound (UCB) Auer (2002), entropy regularization Haarnoja et al. (2018), and curiosity-072 based exploration Groth et al. (2021); Pathak et al. (2017) have shown promising results, they struggle 073 with the escalated complexity in MARL scenarios, compounded by issues like deceptive rewards 074 and the "Noisy-TV" problem Burda et al. (2018). Integrating domain knowledge into exploration 075 could significantly enhance exploration efficiency, by helping identify critical elements and problem 076 structures, thereby aiding in the selection of optimal actions Simon (1956). However, the integration 077 of this knowledge within a data-driven framework poses significant challenges, particularly when 078 manual input from domain experts is required, thus reducing its practicality. 079

Recently, Large Language Models (LLMs) such as GPT-4 Achiam et al. (2023) have shown formidable skills in language comprehension, strategic planning, and logical reasoning across various tasks Yao 081 et al. (2023); Zhu et al. (2023). Although not always directly solving complex, dynamic problems, 082 their inferential and error-learning abilities facilitate progressively better solutions through iterative 083 feedback Ma et al. (2024). The integration of LLMs with MARL presents a promising new avenue 084 by facilitating exploration through the pruning of redundant actions. In this paper, we introduce 085 Evolutionary action SPAce Reduction with Knowledge (eSpark), a novel approach that utilizes LLMs to improve MARL training via optimized exploration functions, which are used to prune the 087 action space. eSpark begins by using LLMs to generate exploration functions from task descriptions 880 and environmental rules in a zero-shot fashion. It then applies evolutionary search within MARL to pinpoint the best performing policy. Finally, by analyzing the feedback on the performance of this 089 policy, eSpark reflects and proposes a set of new exploration functions, and iteratively optimizes them according to the aforementioned steps. This process enhances the MARL policy by continuously 091 adapting and refining exploration. To summarize, our contributions are as follows: 092

1. We introduce the eSpark framework, which harnesses the intrinsic prior knowledge and encoding capability of LLMs to design exploration functions for action space pruning, thus guiding the exploration and learning process of MARL algorithms. eSpark requires no complex prompt engineering and can be easily combined with MARL algorithms.

094

095

096

098

099

102

- 2. We validate the performance of eSpark across 15 different environments within the inventory management task MABIM Yang et al. (2023) and the traffic signal control task SUMO Behrisch et al. (2011). Combined with IPPO, eSpark outperforms IPPO in all scenarios, realizing an average profit increase of 34.4% in the MABIM and improving multiple metrics in SUMO by an average of 9.9%. Even in the face of scalability challenges where the DTDE methods typically encounter limitations, eSpark elevates the performance of IPPO by 29.7%.
- 3. We conduct controlled experiments and ablation studies to analyze the effectiveness of each component within the eSpark framework. We first validate the advantages of knowledge-based pruning. Subsequently, we conduct ablation studies to demonstrate that both retention training and LLM pruning techniques contribute to the performance of eSpark. These effects are even more pronounced in the more complex MABIM environment.

108 2 RELATED WORKS

110 LLMs for code generation. LLMs have exhibited formidable capabilities in the domain of code 111 generation Roziere et al. (2023); Nejjar et al. (2023). More recently, LLM-based approaches to self-112 improving code generation have been applied to address challenges in combinatorial optimization Liu 113 et al. (2024); Ye et al. (2024), robotic tasks Liang et al. (2023); Wang et al. (2024), reinforcement 114 learning (RL) reward design Ma et al. (2024), and code optimization Romera-Paredes et al. (2024). For the first time, to the best of our knowledge, we propose the use of LLMs to generate code with 115 116 the aim of pruning the redundant action space in MARL environments, and through a process of autonomous reflection and evolution, iteratively enhancing the quality of the generated output. 117

118 LLMs for RL and MARL. The integration of LLMs into RL and MARL has sparked considerable 119 research interest Sharma et al. (2022); Kwon et al. (2023); Li et al. (2024). Some works Hill et al. 120 (2020); Chan et al. (2019) incorporate the goal descriptions of language models that help in enhancing the generalization capabilities of agents designed to follow instructions. Further studies have extended 121 this approach to complex tasks involving reasoning and planning Huang et al. (2023; 2022). Moreover, 122 LLMs have been employed to guide exploration and boost RL efficiency Du et al. (2023); Chang et al. 123 (2023); Hu & Sadigh (2023). However, scaling to high-complexity, real-time, multi-agent settings 124 remains a challenge. Our method mitigates this by generating exploration functions to navigate the 125 policy space, thus facilitating the application to complex MARL scenarios without direct LLM-agent 126 decision-making interaction. 127

Action space pruning in RL and MARL. Pruning the action space has been shown to be effective 128 in guiding agent behaviors in complex environments Lipton et al. (2016); Fulda et al. (2017). 129 Techniques include learning an elimination signal to discard unnecessary actions Zahavy et al. (2018), 130 and employing transfer learning that pre-trains agents to isolate useful experiences for later action 131 refinement Shirali et al. (2023); Lan et al. (2022); Ammanabrolu & Riedl (2018). Some approaches 132 use manually designed data structures based on prior knowledge to filter actions Dulac-Arnold et al. 133 (2015); Padullaparthi et al. (2022); Nagarathinam et al. (2020). However, training pruning signals or 134 applying transfer learning is inherently difficult with many agents, and the need for expert knowledge 135 in manual pruning rules hampers their transferability, limiting the applicability of current methods. 136 We harness the abundant knowledge embedded within LLMs for action space pruning, demonstrating 137 universal applicability across a multitude of scenarios.

138 139

140 141

142

3 PRELIMINARIES

3.1 PROBLEM FORMULATION AND NOTATIONS

At each discrete time step t, the environment is in state $s_t \in S$. Each agent $k \in [1, 2, ..., N]$ receives an observation $o_t^k \in \mathcal{O}^k$ and draws an action from $a_t^k \sim \pi^k(\cdot \mid o_t^k)$, where $\pi^k : \mathcal{O}^k \times \mathcal{A}^k \to [0, 1]$ denotes the policy of agent k, and $\sum_{a^k \in \mathcal{A}^k} \pi^k(a^k \mid o_t^k) = 1$. The joint actions of all agents $\mathbf{a}_t = (a_t^1, a_t^2, ..., a_t^N)$ is drawn from the joint policy $\pi(\cdot \mid s_t) = \prod_{k=1}^N \pi^k(\cdot \mid o_t^k)$. Subsequently, a reward $\mathbf{r}_t = R(s_t, \mathbf{a}_t)$ is given based on the current state and joint action. The state transition is determined by $s_{t+1} \sim P(\cdot \mid s_t, \mathbf{a}_t)$.

In this paper, we focus on a fully cooperative scenario where all agents share a common reward signal. The collective objective is to maximize the expected cumulative reward, starting from an initial state distribution ρ^0 . This collaborative approach emphasizes the alignment of individual agent strategies towards maximizing a unified reward $J(\pi)$:

$$J(\boldsymbol{\pi}) = \mathbb{E}_{s_0 \sim \rho^0, \mathbf{a}_{0:\infty} \sim \boldsymbol{\pi}, s_{1:\infty} \sim P} \left[\sum_{t=0}^{\infty} \gamma^t \mathbf{r}_t \right].$$

162 Policy with exploration function. In the configuration of our study, we introduce the exploration 163 function $E: \mathcal{O}^k \times \mathcal{A}^k \to \{0, 1\}$, indicating whether an action is selectable by agent k. For a given 164 policy π^k of agent k and an exploration function E, we define a new policy π^k_E as follows: 165

$$\pi_E^k(\cdot \mid o_t^k) = \frac{\pi^k(\cdot \mid o_t^k) \cdot E(o_t^k, \cdot)}{\sum_{a^k \in \mathcal{A}^k} \pi^k(a^k \mid o_t^k) \cdot E(o_t^k, a^k)}$$

if $\sum_{a^k \in \mathcal{A}^k} \pi^k(a^k \mid o_t^k) \cdot E(o_t^k, a^k) > 0$; otherwise, $\pi_E^k(\cdot \mid o_t^k) = \pi^k(\cdot \mid o_t^k)$. Consequently, the joint policy for all agents under the guidance of E is defined as: 168 169

$$\boldsymbol{\pi}_E(\cdot \mid s_t) = \prod_{k=1}^N \pi_E^k(\cdot \mid o_t^k)$$

173 We define the set of all joint policies as $\{\pi\}$ and the set of all exploration functions as $\{E\}$. Let $\{\pi_E\}$ 174 denote the set of joint policies when subjected to an exploration function $E \in \{E\}$. An exploration 175 function E is non-trivial if it assigns a zero probability to at least one observation-action pair. The 176 following proposition naturally arises from the definition: 177

Proposition 1.

1. For any $E \in \{E\}$, $\{\pi_E\} \subseteq \{\pi\}$. If E is non-trivial, then $\{\pi_E\} \subset \{\pi\}$. 2. For any $\pi \in \{\pi\}$, there exists a non-trivial $E \in \{E\}$ such that $J(\pi_E) \ge J(\pi)$.

An intelligent choice of exploration functions does not diminish our ability to discover optimal policies; instead, it allows us to refine the policy space, thereby enhancing the efficiency of the learning process. The proof of this proposition can be found in Appendix A.

3.2 CHALLENGES AND MOTIVATIONS

The intricate relationships among multiple agents make it extremely difficult to search for the optimal 188 solution in MARL. Without powerful exploration methods, it is nearly impossible to avoid suboptimal 189 outcomes. We will elaborate on this with an example from the following proposition: 190

Proposition 2. Let's consider a fully cooperative game with N agents, one state, and the joint action 191 space $\{0,1\}^N$, where the reward is given by $r(\mathbf{0}^0,\mathbf{1}^N) = r_1$ and $r(\mathbf{0}^{N-m},\mathbf{1}^m) = -mr_2$ for all 192 $m \neq N$, r_1 , r_2 are positive real numbers. We suppose the initial policy is randomly and uniformly 193 initialized, and the policy is optimized in the form of gradient descent. Let p be the probability that 194 the shared policy converges to the best policy, then: 195

166 167

170 171 172

178

179

181

182

183

185

186 187

197

 $p = 1 - \sqrt[N-1]{\frac{r_2}{r_1 + Nr_2}}.$

Detailed proof is provided in Appendix A. In the above example, we show that the increase in 199 the number of agents makes it more difficult for MARL algorithms to reach the optimal solution. 200 However, based on the problem context, humans can understand problems from a high-level semantic 201 perspective, and quickly find optimal solutions. As LLMs have demonstrated surprising abilities in 202 semantic understanding, reasoning, and planning across various tasks Yuan et al. (2023); Wang et al. 203 (2023), we conduct a simple experiment to test GPT-4's capability for the issue in Proposition 2, and 204 here is GPT-4's response: 205

206 In a fully cooperative game, all agents work together to maximize the total reward. There are two distinct reward conditions: 207 208 1. When all agents choose action 1, the reward is r_1 , a positive real number. 209 210 2. When there is any number of agents m (where 0 < m < N) choosing action 1, the reward is $-mr_2$, where r_2 is a positive real number. 211 212 All agents should act in a way that avoids the negative reward scenario. 213 The negative reward scenario happens anytime there is a mix of 0's and 214 1's in the action space, which means some agents are choosing 1 and others are choosing 0. Therefore, the optimal joint action for all 215 agents is to all choose 1.

GPT-4 exhibits reasoning abilities on par with those of humans and directly solves the problem in
 Proposition 2. Propositions 1 has already shown that an intelligent exploration function can not
 only reduce the searching space but also improve the final performance. This makes us think about
 applying the powerful LLMs to prune the redundant action space and thereby guide the exploration
 in MARL. In the following sections, we propose the eSpark framework, which integrates the prior
 knowledge and inferential capability of LLMs to boost the exploration in MARL.

4 Method

222 223

224 225

226

227

228

229

230 231

232

237

238

239

In this section, we introduce a novel framework, eSpark, which integrates robust prior knowledge encapsulated in LLMs. It improves iteratively through a cycle of trial and error, leveraging the capability of LLMs. Figure 1 illustrates the overall training procedure. eSpark is composed of three components: (i) zero-shot generation of exploration functions, (ii) evolutionary search for best performing MARL policy, and (iii) detailed feedback on the performance of the policy to improve the generation of exploration functions. We show the pseudocode of eSpark in Appendix B.

4.1 EXPLORATION FUNCTION GENERATION

LLMs have been demonstrated to possess exceptional capabilities in both code comprehension and generation. To this end, we employ a LLM as **LLM code generator**, denoted as LLM_g , whose task is to understand the objectives of the current environment, and output an exploration function:

 $E_1, \dots, E_K \sim \text{LLM}_g(\text{prom}),$ (1)

where prom is the prompt for 240 LLM_q , and the generation of K ex-241 ploration functions is to circum-242 vent the suboptimality that may 243 arise from single-sample genera-244 tion. The initial prom includes 245 an *RL* formulation describing the 246 reward system, state items, transi-247 tions, and the action space, along-248 side a *task description* that speci-249 fies the task objectives, expected 250 outputs, and formatting rules. Details on the initial prom are pro-251 vided in Appendix I. We use code 252 for the RL formulation as it effec-253 tively captures the physical transi-254 tion dynamics crucial to RL prob-255 lems, which are always difficult to 256 express precisely through the text 257 alone, especially when environ-258 mental complexity increases. Code 259 contexts also improve code gener-260 ation and clarify environmental se-261 mantics and variable roles Ma et al.



Figure 1: eSpark firstly generates K exploration functions via zero-shot creation. Each exploration function is then used to guide an independent policy, and the evolutionary search is performed to find the best-performing policy. Finally, eSpark reflects on the feedback from the best performance policy, refines and regenerates the exploration functions for the next iteration.

(2024). In Appendix F, we discuss the impact of different forms of *RL formulation* on the final
 performance of eSpark when the environment code is unavailable.

During the code generation, however, LLM_g may incorrectly interpret variables and produce logically
flawed code. This kind of flawed logic could persist if it is added to the prompt context for the
next generation. As research has shown that collaboration among multiple LLMs can enhance the
quality and efficacy of the generated contents Chen et al. (2023); Zhang et al. (2023), we introduce
the LLM checker denoted as LLM_c, which reviews LLM_g's output to pursue an enhanced generation.
LLM_c uses the same prompt as LLM_g but is prompted to focus on verifying the accuracy of code
relative to environmental transitions and variable specifications. If inconsistencies are found, LLM_c

signals the error, prompting LLM_g to regenerate the code. The reasons for introducing LLM_c are further discussed in Section 5.5. Finally, exploration functions are generated by:

- 273
- 274 275

276

277

$$E_1, \ldots, E_K \sim \text{LLM}_c(\text{prom}, \text{LLM}_g(\text{prom})).$$
 (2)

Exploration functions are applied only during the training phase to guide the exploration of MARL. During the execution phase, all exploration functions are removed.

4.2 EVOLUTIONARY SEARCH

280 During the generation, however, it should be noted that the initially generated exploration function 281 may not always guarantee executability and effectiveness. To address this, eSpark performs an 282 evolutionary search paradigm that selects the best-performing policy in one iteration and uses its 283 feedback for subsequent generation Ma et al. (2024). Specifically, eSpark samples a batch of K 284 exploration functions in each generation to ensure there are enough candidates for successful code execution. Performance is assessed at regular checkpoints within an iteration, with the final evaluation 285 based on the average of the last few checkpoints. The policy achieving the highest performance 286 is selected, and the feedback obtained from this policy is integrated to optimize the exploration 287 functions in the following steps. 288

Due to the dynamic nature of exploration, the exploration function generated based on feedback from the best-performing policy may not be applicable to other policies. As the proof of Proposition 1 demonstrates, when an exploration function is incapable of intelligently pruning, it may even impair the performance of the policy. To this end, we utilize **retention training** to maintain continuity of exploration. Let ϕ_{best}^{i-1} represent best-performing policy from the (i-1)-th iteration. For the *i*-th iteration except for the first, at the start of the iteration, we set:

- 295
- 296 297

298 299 $\phi_1^i, \phi_2^i, \dots, \phi_K^i \leftarrow \phi_{\mathsf{best}}^{i-1}. \tag{3}$

This allows us to match exploration functions with their corresponding policies, subsequently refining performance incrementally. We will verify the impact of retention training in Section 5.4.

300 4.3 Reflection and feedback

301 302

303

304

305

306

307

308

309 310 Feedback from the environment can significantly enhance the quality of the generated output by LLMs Nascimento et al. (2023); Du et al. (2023). In eSpark, we leverage **policy feedback**, which contains the evaluation of policy performance from various aspects, to enhance the generation of LLMs. This policy feedback may either come from experts or be automatically constructed from the environment, as long as it encompasses insights into the aspects where the current algorithm performs well and areas where it requires improvement. As illustrated in Equation 4, by correlating the best-performing policy feedback F_{best} and the most effective exploration function E_{best} , LLMs introspect, update the prompt prom, and gear up for the ensuing evolutionary cycle.

$$prom \leftarrow prom : Reflection(E_{best}, F_{best}).$$
(4)

In our experiments, we generate automated policy feedback from environmental reward signals, as domain experts in relevant fields are not available. We acknowledge that obtaining feedback from human experts can be expensive. Nevertheless, it is important to note that within our framework, the number of rounds for feedback collection is specified by a predefined hyperparameter, which is typically kept low (in our experiments, it is set to 10). Therefore, in scenarios where human experts are accessible, incorporating their insights is feasible and can potentially enhance performance.

317 318

319

5 EXPERIMENTS

320 5.1 EXPERIMENT SETTINGS

321

For a comprehensive evaluation of eSpark's capabilities, we perform detailed validations within two distinct industrial environments: the inventory management environment MABIM and the traffic signal control environment SUMO.

- 324 • MABIM setting: MABIM simulates multi-echelon inventory management by modeling 325 each stock-keeping unit (SKU) as an agent, mirroring real-world operations and profits 326 within the MARL framework. The total reward is composed of multiple reward components. 327 We utilize the total reward to identify the best-performing policy, while those components 328 evaluate the policy's multifaceted performance to generate policy feedback. We focus on three key challenges within inventory management: multiple echelons, capacity constraints and scalability, selecting corresponding scenarios for experiments. 330
- 331 • SUMO setting: SUMO is a traffic signal control environment in which each intersection is 332 represented as an agent. It offers a variety of reward functions, and we use "the number of 333 stopped vehicles" as the reward for evolutionary search, while other rewards are for policy feedback. The Average Delay, Average Trip Time, and Average Waiting Time metrics are 334 chosen for evaluation Lu et al. (2023). We employ GESA Jiang et al. (2024) to standardize 335 intersections into 4-arm configurations. Each simulation spans 3600 seconds, with decisions 336 at 15-second intervals.
 - Model setting: We use IPPO as the base MARL framework for eSpark due to its DTDE structure, which is suitable for large-scale challenges. But note that our approach can also be applied as a plugin in other MARL methods. We select GPT-4 for the LLM_c and LLM_q due to its superior comprehension and generation abilities. For each scenario, we conduct three runs with a batch size of K = 16. espark has the same number of training steps as the compared MARL baselines, with 10 iterations evenly selected throughout the training process for feedback, reflection, and exploration function editing.

345 All training jobs are executed with an Intel(R) Xeon(R) Gold 6348 CPU and 4 NVIDIA RTX A6000 346 GPUs. In Appendix C, we provide a detailed introduction and setting for the environments and model. 347 In Appendix D, we give hyperparameter configurations and descriptions of each baseline method.

349 350

348

337

338

339

340

341

342

343

344

5.2 EXPERIMENT RESULTS

351 In this section, we present the key findings for eSpark in the MABIM and SUMO, highlighting the 352 best and second best results in **bold** and <u>underline</u>. More detailed results and computational costs are presented in Appendix E. 353

354 355

364

PERFORMANCE ON MABIM 5.2.1

356 Table 1 shows the results in the 100 SKUs scenarios in terms of capability constraints and multiple 357 echelon challenges. With IPPO as the base MARL algorithm, eSpark not only outperforms IPPO in 358 all scenarios but also exceeds the performance of all compared baselines in 4 out of 5 scenarios. For 359 an in-depth analysis, we discuss the policy differences between IPPO and espark in Appendix G, 360 along with eSpark's reflective mechanism and exploration function adjustments in Appendix J. 361 While IPPO struggles to learn the intricate interplay among SKUs, eSpark excels particularly in 362 navigating cooperation among SKUs and refining its search in a broad space, leading to marked 363 improvements in managing capacity constraints and multi-echelon coordination.

Table 1: Performance in MABIM, a higher profit indicates a better performance. The "Standard" 365 scenario features a single echelon with sufficient capacity. The "2/3 echelons" involves challenges of 366 multi-echelon cooperation. The "Lower/Lowest" scenarios are the challenges where SKUs compete 367 for insufficient capacity, while "500 SKUs scenarios" assess scalability. The '-' indicates CTDE 368 algorithms are not researched in the scalability challenges. 369

					Avg. pr	ofits (K)				
Method		100 S	KUs scenario	s	01	. ,	500 S	KUs scenario	s	
	Standard	2 echelons	3 echelons	Lower	Lowest	Standard	2 echelons	3 echelons	Lower	Lowest
IPPO	690.6	1412.5	1502.9	431.1	287.6	3021.2	7052.0	7945.7	3535.9	2347.4
QTRAN	529.6	1595.3	2012.2	70.1	19.5	-	-	-	-	-
QPLEX	358.9	1580.7	704.2	379.8	259.3	-	-	-	-	-
MAPPO	719.8	1513.8	1905.4	478.3	265.8	-	-	-	-	-
BS static	563.7	1666.6	2338.9	390.7	-1757.5	3818.5	8151.2	11926.3	3115.1	-9063.8
BS dynamic	684.2	1554.2	2378.2	660.6	-97.1	4015.7	8399.3	11611.1	3957.5	2008.6
(S,s)	<u>737.8</u>	1660.8	1725.2	556.9	203.7	<u>4439.4</u>	9952.1	10935.7	3769.3	2212.4
eSpark	823.7	1811.4	2598.7	<u>579.5</u>	405.0	4468.6	<u>9437.3</u>	12134.2	<u>3775.7</u>	2519.5

In Table 1, we also present the performance outcomes of the eSpark algorithm in the scaling-up 500 SKUs scenarios. Due to the centralized nature of the CTDE methods, they struggle to scale to large-scale problems and therefore are not presented in the table. Despite IPPO's markedly inferior performance on scenarios when problems scale up, eSpark exhibits significant enhancements and consistently achieves optimal results across multiple scenarios. We attribute this improvement to eSpark's action space pruning strategy, which effectively addresses the heightened exploration needs in scenarios with many agents, providing a clear advantage in such complex environments.

385 386

387

388

389

390

391

392

393

394

396

397

405

5.2.2 PERFORMANCE ON SUMO

To further assess eSpark's capabilities across different tasks, we have compiled a summary of results in Table 2 based on the SUMO environment. Similar to the outcomes in MABIM, eSpark consistently enhances the performance of the IPPO in all scenarios, and it has outperformed the CTDE baselines as well as domain-specific MARL baselines to achieve the best performance. Notably, even when IPPO alone is capable of good results (as seen in scenarios such as Grid 4×4 and Cologne8), the pruning method designed in eSpark does not compromise the effectiveness of IPPO. We will delve further into the analysis of exploration functions produced by eSpark in Section 5.3.

Table 2: Performance in SUMO, including the mean and standard deviation (in parentheses). A lower time indicates a better performance.

Mathad		Av	g. delay (second	s)		Avg. trip time (seconds)				
Method	Grid 4×4	Arterial 4×4	Grid 5×5	Cologne8	Ingolstadt21	Grid 4×4	Arterial 4×4	Grid 5×5	Cologne8	Ingolstadt21
FTC	161.14 (3.77)	1229.68 (16.79)	820.88 (24.36)	85.27 (1.21)	224.96 (11.91)	291.48 (4.45)	676.77 (13.70)	584.54 (24.17)	145.93 (0.84)	352.06 (9.29)
MaxPressure	63.39 (1.34)	995.23 (77.02)	242.85 (16.04)	31.63 (0.61)	228.64 (15.83)	174.68 (2.05)	702.09 (23.61)	269.35 (9.62)	95.67 (0.62)	352.30 (15.06)
IPPO	48.40 (0.45)	884.73 (38.94)	228.78 (11.59)	27.60 (1.70)	342.97 (43.61)	160.12 (0.60)	506.18 (10.39)	238.03 (7.10)	91.41 (1.60)	464.50 (43.30)
MAPPO	51.25 (0.58)	958.43 (181.72)	958.43 (181.72)	32.55 (4.66)	347.59 (47.59)	160.01 (0.54)	757.40 (242.00)	247.56 (3.71)	94.31 (1.77)	480.66 (49.46)
CoLight	53.40 (1.89)	820.67 (58.65)	339.66 (48.55)	27.48 (1.30)	296.47 (106.82)	165.77 (1.89)	470.33 (12.34)	305.41 (44.43)	91.66 (1.29)	410.59 (97.29)
MPLight	63.51 (0.64)	1142.98 (79.65)	223.44 (16.18)	37.93 (0.45)	196.74 (9.88)	172.47 (0.89)	583.21 (45.84)	255.49 (6.26)	110.56 (1.18)	331.42 (11.79)
eSpark	48.36 (0.32)	854.22 (68.21)	209.49 (13.98)	25.39 (1.27)	243.92 (15.81)	159.74 (0.44)	484.87 (58.21)	235.20 (6.80)	89.50 (1.36)	367.57 (15.03)

5.3 ESPARK LEARNS INTELLIGENT PRUNING METHODS

Given that eSpark employs the prior knowledge of the LLMs to craft its exploration function, our study aimed to investigate two critical aspects: (1) the validity of action space pruning via prior knowledge, and (2) the potential advantages of this method over rule-based heuristic pruning.

To address the questions raised, we devise two pruning strategies. First, we implement a **random** 410 **pruning** method, wherein agents randomly exclude a portion of actions during decision-making to 411 test the validity of knowledge-based pruning. Secondly, we utilize domain-related OR algorithms to 412 implement heuristic pruning methods. For MABIM, actions are pruned using the (S, s) policy and 413 unbound limit, while for SUMO, pruning relies on MaxPressure to keep only a few actions with the 414 highest pressure. The details of these methods are presented in Appendix D.3. Just like eSpark, 415 these pruning strategies are integrated with IPPO during training but not execution. We conducted 416 experiments under the same setting in Section 5.1, with results presented in Tables 3 and Table 4. 417

Table 3: Average performance changes on MABIM. All changes are relative to IPPO.

Table 4: Average performance changes on SUMO. All changes are relative to IPPO.

	inges are ren		Solvio. All changes are relative to if 10.				
Method	Avg. profits 100 SKUs	change ratio (%) 500 SKUs	Method	Avg. 1 Delay	ime change Trip time	ratio (%) Wait time	
Random pruning	2.1	-0.5	Random pruning	-0.1	2.2	-2.5	
(S, s) pruning	-25.9	15.5	MaxPressure pruning	-0.5	1.5	-0.1	
Upbound pruning	-23.2	-32.7	eSpark	-9.7	-5.7	-14.3	
eSpark	39.1	29.7					

⁴²⁴ 425

418

419

426 As shown in the tables, random pruning marginally affects performance by merely altering exploration 427 rates without providing new insights. Heuristic pruning's impact varies with its design and context. In 428 MABIM, (S, s) pruning is less effective in the 100 SKUs scenario, as it restricts the already effective 429 IPPO's exploration in smaller scales. However, it proves beneficial in the 500 SKUs scenario, where it 430 guides the exploration and leads to better results. Upbound pruning consistently underperforms due to 431 its overly simplistic heuristic. For SUMO, pressure-based pruning does not offer significant benefits. 432 Nevertheless, eSpark demonstrates remarkable adaptability across all testing tasks, adeptly selecting pruning methods that substantially enhance results. Its knowledge-based generative technique and
 evolution capability enable it to master intelligent pruning strategies.



Figure 2: Action selection frequency for IPPO and various pruning methods on the 100 SKUs Lowest scenario. "Actions" represents the replenishment quantity is a multiple of the mean demand within the sliding window. eSpark learns not only to minimize restocking but also to diversify with small purchases below the mean demand, balancing demand fulfillment and overflow prevention.

Figure 2 presents a frequency heatmap of action selection for IPPO and various pruning methods in the 100 SKUs Lowest scenario. IPPO learns a minimally restocking strategy, risking unmet demand. Random pruning chooses actions more uniformly yet mirrors IPPO's pattern. (S, s) pruning excessively exceeds mean demand, ignoring no-restock actions and leading to significant overflow. Upbound pruning typically avoids restocking, but prefers to purchase near the mean demand, which could result in overflow costs. In contrast, eSpark adopts a balanced policy, avoiding overstocking while diversifying its minor restocking strategies to meet demand without causing overflow.

5.4 ESPARK BENEFITS FROM RETENTION TRAINING AND ACTION SPACE REDUCTION

Extensive research has underscored the importance of reflection in LLM-driven content generation Ma
et al. (2024); Nascimento et al. (2023). Herein, we focus on the effects of retention training and
action pruning on eSpark's performance.

We first design an ablation experiment, which we refer to as the eSpark w/o retention. The model parameters are initialized when an iteration is finished, and the newly generated exploration functions are equipped, after which the training starts from scratch. Given that the initialized model needs a more extensive number of steps to converge, we accordingly triple the training steps per iteration in comparison to the standard eSpark. Another ablation retains the retention training, while the only difference is that the LLMs and reflection are removed. We name this experiment eSpark w/o LLM. The comparative analysis of these two ablations is delineated in Table 5 and Table 6.

Table 5: Average performance change across
100 SKU scenarios in the MABIM environ-
ment All changes are relative to IPPO

Table 6: Average performance change in the SUMO environment. All changes are relative to IPPO.

mente i m enunges u	<u>to ii i o.</u>				
Method	Avg. profits change ratio (%)	Method	Avg. t	ime change	ratio (%)
eSpark	39.1	- C1-	Delay		
eSpark w/o retention eSpark w/o LLM	24.0 -2.8	eSpark eSpark w/o retention	-9.7 -9.6	- 5. 7 -4.6	-14.3 -11.2
	2.0	eSpark w/o LLM	-9.1	-5.0	-12.8

The removal of retention training and LLMs both result in a decline in the performance of the eSpark. In the SUMO scenario, the performance gap between the two ablations and the complete eSpark is relatively small, whereas it is more pronounced in the MABIM scenarios. This can be attributed to the fact that MABIM involves a greater number of agents and a more complex observation space action space, where a superior pruning can significantly enhance the performance of MARL methods. Additionally, we observe that the lack of LLMs leads to a significant decrease in performance on MABIM, emphasizing the central role of knowledge-based action space pruning within the eSpark.

5.5 LLM CHECKER AND DETAILED REWARD FEEDBACK PROMOTES THE PERFORMANCE OF ESPARK

- We go deeper to investigate the impact of feedback prompt design and the LLM checker on the performance of eSpark. eSpark utilizes detailed reward factors to evaluate policy performance

486
 487
 488
 488
 488
 488
 488
 488
 488

The introduction of LLM checker comes from the following observations: LLM code generator occasionally produces flawed exploration functions (e.g., variable misuse, misaligned task logic). Although these are usually eliminated during evolutionary search, when the pool of executable exploration functions is small, flawed yet executable functions may still be selected as editing templates, thus hindering further refinement. We conduct the second ablation **eSpark w/o checker** by removing the LLM checker. Both ablations are performed in the MABIM (100 SKUs) and SUMO environments, with the performance averaged in respective environments and presented in Table 7 and Table 8.

Table 7: Average performance change across 100 SKU scenarios in the MABIM environment. All changes are relative to IPPO.

Table 8: Average performance change in the
SUMO environment. All changes are relative
to IPPO.

			10 11 1 01				
	Method	vg. profits change ratio (%) 39.1 17.8	Method	Avg. time change ratio (%)			
	eSpark	30.1		Delay	mp une	wait time	
	eSpark w/o r factors	17.8	eSpark	-9.7	-5.7	-14.3	
	eSpark w/o checker	26.2	eSpark w/o r factors	-9.4	-4.8	-13.5	
		20.2	eSpark w/o checker	-8.1	-3.1	-10.8	

After removing detailed reward information, eSpark showed a significant performance decline in both the MABIM and SUMO environments, particularly in MABIM, where complex transition 508 logic and variable usage exacerbated the impact. Without detailed reward factors, eSpark struggled 509 to analyze policy performance and propose targeted improvements, leading to a diminished ability 510 to refine the action space. We provide examples of eSpark's responses with and without reward 511 factors in Appendix H to offer more insights into these results. Additionally, the ablation experiment 512 on the LLM checker revealed its critical role in preventing flawed exploration functions from being 513 selected as editing templates, further demonstrating the importance of both detailed reward feedback 514 and the LLM checker in maintaining eSpark's overall performance. 515

516

496 497 498

499

500

501

517 518

519

6 CONCLUSIONS, LIMITATIONS AND FUTURE WORK

520 We present eSpark, a novel framework for generating exploration functions, leveraging the advanced 521 capabilities of LLMs to integrate prior knowledge, generate code and reflect, thereby refining the 522 exploration in MARL. eSpark has surpassed its base MARL algorithm across all scenarios in 523 both MABIM and SUMO environments. In terms of pruning strategies, pruning based on the prior 524 knowledge from LLMs outshines both random and heuristic approaches. Ablation experiments 525 demonstrate the indispensable role of retention training in accurately improving exploration functions based on policy flaws and enhancing sample efficiency. The LLM checker and detailed policy 526 feedback prompt design together ensure the superior performance of eSpark. 527

Nevertheless, eSpark also has certain limitations. First, currently eSpark is only applicable to tasks involving homogeneous agents. For heterogeneous agents, a potential method could be to generate distinct exploration functions for each agent; however, this approach becomes impractical when the number of agents is too large. Moreover, eSpark benefits from policy feedback to refine the exploration functions. When feedback is not informative regarding how to modify the exploration (e.g., in tasks with sparse rewards, end-of-episode feedback alone is too limited to develop automated feedback), eSpark may struggle to improve and need extra expert input for effective reflection.

Future work encompasses numerous potential directions. Existing research advocates for assigning
different roles or categories to agents Christianos et al. (2021); Wang et al. (2020), which could offer
a compromise for the application of eSpark in heterogeneous multi-agent systems. Furthermore,
state-specific feedback for more granular improvement represents an intriguing avenue Subramanian
et al. (2016). Our future endeavors will investigate these questions, striving to develop algorithms
that are robust and exhibit strong generalizability.

540 REFERENCES

542	Josh Achiam, Steven Adler, Sandhini Agarwal, Lama Ahmad, Ilge Akkaya, Florencia Leoni Aleman,
543	Diogo Almeida, Janko Altenschmidt, Sam Altman, Shyamal Anadkat, et al. Gpt-4 technical report.
544	arXiv preprint arXiv:2303.08774, 2023.

- Prithviraj Ammanabrolu and Mark O Riedl. Playing text-adventure games with graph-based deep reinforcement learning. *arXiv preprint arXiv:1812.01628*, 2018.
- Kenneth J Arrow, Theodore Harris, and Jacob Marschak. Optimal inventory policy. *Econometrica: Journal of the Econometric Society*, pp. 250–272, 1951.
- Peter Auer. Using confidence bounds for exploitation-exploration trade-offs. *Journal of Machine Learning Research*, pp. 397–422, 2002.
- Michael Behrisch, Laura Bieker, Jakob Erdmann, and Daniel Krajzewicz. Sumo–simulation of urban
 mobility: an overview. In *SIMUL*, 2011.
- Alan S Blinder. *Inventory theory and consumer behavior*. Harvester Wheatsheaf, 1990.
- Yuri Burda, Harrison Edwards, Amos Storkey, and Oleg Klimov. Exploration by random network distillation. In *ICLR*, pp. 1–17, 2018.
- Harris Chan, Yuhuai Wu, Jamie Kiros, Sanja Fidler, and Jimmy Ba. Actrce: Augmenting experience via teacher's advice for multi-goal reinforcement learning. *arXiv preprint arXiv:1902.04546*, 2019.
- Jonathan D Chang, Kiante Brantley, Rajkumar Ramamurthy, Dipendra Misra, and Wen Sun. Learning to generate better than your llm. *arXiv preprint arXiv:2306.11816*, 2023.
- Chacha Chen, Hua Wei, Nan Xu, Guanjie Zheng, Ming Yang, Yuanhao Xiong, Kai Xu, and Zhenhui
 Li. Toward a thousand lights: Decentralized deep reinforcement learning for large-scale traffic
 signal control. In *AAAI*, pp. 3414–3421, 2020.
- Weize Chen, Yusheng Su, Jingwei Zuo, Cheng Yang, Chenfei Yuan, Chen Qian, Chi-Min Chan, Yujia
 Qin, Yaxi Lu, Ruobing Xie, et al. Agentverse: Facilitating multi-agent collaboration and exploring
 emergent behaviors in agents. *arXiv preprint arXiv:2308.10848*, 2023.
- Filippos Christianos, Georgios Papoudakis, Muhammad A Rahman, and Stefano V Albrecht. Scaling
 multi-agent reinforcement learning with selective parameter sharing. In *ICML*, pp. 1989–1998,
 2021.
- 575
 576
 576
 576
 577
 578
 578
 578
 578
 579
 578
 574
 575
 576
 577
 578
 578
 578
 578
 578
 578
 578
 578
 578
 578
 578
 578
 578
 578
 578
 578
 578
 578
 578
 578
 578
 578
 578
 578
 578
 578
 578
 578
 578
 578
 578
 578
 578
 578
 578
 578
 578
 578
 578
 578
 578
 578
 578
 578
 578
 578
 578
 578
 578
 578
 578
 578
 578
 578
 578
 578
 578
 578
 578
 578
 578
 578
 578
 578
 578
 578
 578
 578
 578
 578
 578
 578
 578
 578
 578
 578
 578
 578
 578
 578
 578
 578
 578
 578
 578
 578
 578
 578
 578
 578
 578
 578
 578
 578
 578
 578
 578
 578
 578
 578
 578
 578
 578
 578
 578
 578
 578
 578
 578
 578
 578
 578
 578
 578
 578
 578
 578
 578
 578
 578
 578
 578
 578
 578
 578
 578
 578
 578
 578
 578
 578
 578
 578
 578
 578
 578
 578
 578
 578
 578
 578
 578
- Yuqing Du, Olivia Watkins, Zihan Wang, Cédric Colas, Trevor Darrell, Pieter Abbeel, Abhishek
 Gupta, and Jacob Andreas. Guiding pretraining in reinforcement learning with large language
 models. In *ICML*, pp. 8657–8677, 2023.
- Gabriel Dulac-Arnold, Richard Evans, Hado van Hasselt, Peter Sunehag, Timothy Lillicrap, Jonathan Hunt, Timothy Mann, Theophane Weber, Thomas Degris, and Ben Coppin. Deep reinforcement learning in large discrete action spaces. *arXiv preprint arXiv:1512.07679*, 2015.
- Nancy Fulda, Daniel Ricks, Ben Murdoch, and David Wingate. What can you do with a rock?
 affordance extraction via word embeddings. *arXiv preprint arXiv:1703.03429*, 2017.
- Oliver Groth, Markus Wulfmeier, Giulia Vezzani, Vibhavari Dasagi, Tim Hertweck, Roland Hafner, Nicolas Heess, and Martin Riedmiller. Is curiosity all you need? on the utility of emergent behaviours from curious exploration. *arXiv preprint arXiv:2109.08603*, 2021.
- Tuomas Haarnoja, Aurick Zhou, Pieter Abbeel, and Sergey Levine. Soft actor-critic: Off-policy maximum entropy deep reinforcement learning with a stochastic actor. In *ICML*, pp. 1861–1870, 2018.

- Jianye Hao, Xiaotian Hao, Hangyu Mao, Weixun Wang, Yaodong Yang, Dong Li, Yan Zheng, and
 Zhen Wang. Boosting multiagent reinforcement learning via permutation invariant and permutation
 equivariant networks. In *ICLR*, 2022a.
- Xiaotian Hao, Hangyu Mao, Weixun Wang, Yaodong Yang, Dong Li, Yan Zheng, Zhen Wang, and
 Jianye Hao. Breaking the curse of dimensionality in multiagent state space: A unified agent
 permutation framework. *arXiv preprint arXiv:2203.05285*, 2022b.
- Felix Hill, Sona Mokra, Nathaniel Wong, and Tim Harley. Human instruction-following with deep
 reinforcement learning via transfer-learning from text. *arXiv preprint arXiv:2005.09382*, 2020.
- Hengyuan Hu and Dorsa Sadigh. Language instructed reinforcement learning for human-ai coordination. In *ICML*, pp. 13584–13598, 2023.
- Wenlong Huang, Pieter Abbeel, Deepak Pathak, and Igor Mordatch. Language models as zero-shot
 planners: Extracting actionable knowledge for embodied agents. In *ICML*, pp. 9118–9147, 2022.
- Wenlong Huang, Fei Xia, Ted Xiao, Harris Chan, Jacky Liang, Pete Florence, Andy Zeng, Jonathan Tompson, Igor Mordatch, Yevgen Chebotar, et al. Inner monologue: Embodied reasoning through planning with language models. In *CoRL*, pp. 1769–1782, 2023.
- Haoyuan Jiang, Ziyue Li, Zhishuai Li, Lei Bai, Hangyu Mao, Wolfgang Ketter, and Rui Zhao. A
 general scenario-agnostic reinforcement learning for traffic signal control. *IEEE Transactions on Intelligent Transportation Systems*, 2024.
- Anastasios Kouvelas, Jennie Lioris, S Alireza Fayazi, and Pravin Varaiya. Maximum pressure
 controller for stabilizing queues in signalized arterial networks. *Transportation Research Record*, 2421(1):133–141, 2014.
- Minae Kwon, Sang Michael Xie, Kalesha Bullard, and Dorsa Sadigh. Reward design with language
 models. In *ICLR*, 2023.
- Yixing Lan, Xin Xu, Qiang Fang, Yujun Zeng, Xinwang Liu, and Xianjian Zhang. Transfer reinforce ment learning via meta-knowledge extraction using auto-pruned decision trees. *Knowledge-Based Systems*, 242:108221, 2022.
- Dapeng Li, Zhiwei Xu, Bin Zhang, and Guoliang Fan. From explicit communication to tacit
 cooperation: A novel paradigm for cooperative marl. *arXiv preprint arXiv:2304.14656*, 2023.
- Dapeng Li, Hang Dong, Lu Wang, Bo Qiao, Si Qin, Qingwei Lin, Dongmei Zhang, Qi Zhang,
 Zhiwei Xu, Bin Zhang, et al. Verco: Learning coordinated verbal communication for multi-agent
 reinforcement learning. *arXiv preprint arXiv:2404.17780*, 2024.
- Jacky Liang, Wenlong Huang, Fei Xia, Peng Xu, Karol Hausman, Brian Ichter, Pete Florence, and
 Andy Zeng. Code as policies: Language model programs for embodied control. In *ICRA*, pp. 9493–9500, 2023.
- Zachary C Lipton, Kamyar Azizzadenesheli, Abhishek Kumar, Lihong Li, Jianfeng Gao, and Li Deng. Combating reinforcement learning's sisyphean curse with intrinsic fear. *arXiv preprint arXiv:1611.01211*, 2016.
- Fei Liu, Xialiang Tong, Mingxuan Yuan, Xi Lin, Fu Luo, Zhenkun Wang, Zhichao Lu, and Qingfu
 Zhang. An example of evolutionary computation+ large language model beating human: Design
 of efficient guided local search. *arXiv preprint arXiv:2401.02051*, 2024.
- Ryan Lowe, Yi Wu, Aviv Tamar, Jean Harb, Pieter Abbeel, and Igor Mordatch. Multi-agent actor critic for mixed cooperative-competitive environments. In *NeurIPS*, 2017.
- Jiaming Lu, Jingqing Ruan, Haoyuan Jiang, Ziyue Li, Hangyu Mao, and Rui Zhao. Dualight: Enhancing traffic signal control by leveraging scenario-specific and scenario-shared knowledge. *arXiv preprint arXiv:2312.14532*, 2023.
- Yecheng Jason Ma, William Liang, Guanzhi Wang, De-An Huang, Osbert Bastani, Dinesh Jayaraman,
 Yuke Zhu, Linxi Fan, and Anima Anandkumar. Eureka: Human-level reward design via coding large language models. In *ICLR*, pp. 1–39, 2024.

Patrick Meier. <i>Digital humanitarians: how big data is changing the face of humanitarian response</i> . Crc Press, 2015.
Seyed Sajad Mousavi, Michael Schukat, and Enda Howley. Traffic light control using deep policy- gradient and value-function-based reinforcement learning. <i>IET Intelligent Transport Systems</i> , 11 (7):417–423, 2017.
Md Shirajum Munir, Nguyen H Tran, Walid Saad, and Choong Seon Hong. Multi-agent meta- reinforcement learning for self-powered and sustainable edge computing systems. <i>IEEE Transac-</i> <i>tions on Network and Service Management</i> , 18(3):3353–3374, 2021.
Srinarayana Nagarathinam, Vishnu Menon, Arunchandar Vasan, and Anand Sivasubramaniam. Marco - multi-agent reinforcement learning based control of building hvac systems. In ACM e-Energy, 2020.
Nathalia Nascimento, Paulo Alencar, and Donald Cowan. Gpt-in-the-loop: Adaptive decision-making for multiagent systems. <i>arXiv preprint arXiv:2308.10435</i> , 2023.
Mohamed Nejjar, Luca Zacharias, Fabian Stiehle, and Ingo Weber. Llms for science: Usage for code generation and data analysis. <i>arXiv preprint arXiv:2311.16733</i> , 2023.
Frans A Oliehoek, Matthijs TJ Spaan, and Nikos Vlassis. Optimal and approximate Q-value functions for decentralized POMDPs. <i>Journal of Artificial Intelligence Research</i> , 32:289–353, 2008.
Venkata Ramakrishna Padullaparthi, Srinarayana Nagarathinam, Arunchandar Vasan, Vishnu Menon, and Depak Sudarsanam. Falcon-farm level control for wind turbines using multi-agent deep reinforcement learning. <i>Renewable Energy</i> , 181:445–456, 2022.
Deepak Pathak, Pulkit Agrawal, Alexei A Efros, and Trevor Darrell. Curiosity-driven exploration by self-supervised prediction. In <i>ICML</i> , pp. 2778–2787, 2017.
Tabish Rashid, Mikayel Samvelyan, Christian Schroeder De Witt, Gregory Farquhar, Jakob Foerster, and Shimon Whiteson. Monotonic value function factorisation for deep multi-agent reinforcement learning. <i>Journal of Machine Learning Research</i> , 21(1):7234–7284, 2020.
Roger P Roess, Elena S Prassas, and William R McShane. <i>Traffic engineering</i> . Pearson/Prentice Hall, 2004.
Bernardino Romera-Paredes, Mohammadamin Barekatain, Alexander Novikov, Matej Balog, M Pawan Kumar, Emilien Dupont, Francisco JR Ruiz, Jordan S Ellenberg, Pengming Wang, Omar Fawzi, et al. Mathematical discoveries from program search with large language models. <i>Nature</i> , 625(7995):468–475, 2024.
Baptiste Roziere, Jonas Gehring, Fabian Gloeckle, Sten Sootla, Itai Gat, Xiaoqing Ellen Tan, Yossi Adi, Jingyu Liu, Tal Remez, Jérémy Rapin, et al. Code llama: Open foundation models for code. <i>arXiv preprint arXiv:2308.12950</i> , 2023.
Pratyusha Sharma, Antonio Torralba, and Jacob Andreas. Skill induction and planning with latent language. In <i>ACL</i> , pp. 1713–1726, 2022.
Ali Shirali, Alexander Schubert, and Ahmed Alaa. Pruning the way to reliable policies: A multi- objective deep q-learning approach to critical care. <i>arXiv preprint arXiv:2306.08044</i> , 2023.
Herbert A Simon. Rational choice and the structure of the environment. <i>Psychological review</i> , 63(2): 129, 1956.
Kyunghwan Son, Daewoo Kim, Wan Ju Kang, David Earl Hostallero, and Yung Yi. Qtran: Learning to factorize with transformation for cooperative multi-agent reinforcement learning. In <i>ICML</i> , pp. 5887–5896, 2019.

Kaushik Subramanian, Charles L Isbell Jr, and Andrea L Thomaz. Exploration from demonstration 701 for interactive reinforcement learning. In AAMAS, pp. 447-456, 2016.

702 703 704	Guanzhi Wang, Yuqi Xie, Yunfan Jiang, Ajay Mandlekar, Chaowei Xiao, Yuke Zhu, Linxi Fan, and Anima Anandkumar. Voyager: An open-ended embodied agent with large language models. <i>arXiv</i> preprint arXiv:2305.16291, 2023.
705	
706	Jianhao Wang, Zhizhou Ren, Terry Liu, Yang Yu, and Chongjie Zhang. QPLEX: Duplex dueling
707	multi-agent Q-learning. In ICLR, pp. 1–27, 2021.
708	Tonghan Wang, Heng Dong, Victor Lesser, and Chongije Zhang. Roma: Multi-agent reinforcement
709	learning with emergent roles. arXiv preprint arXiv:2003.08039, 2020.
710	
711	Yuki Wang, Gonzalo Gonzalez-Pumariega, Yash Sharma, and Sanjiban Choudhury. Demo2code:
712	From summarizing demonstrations to synthesizing code via extended chain-of-thought. In NeurIPS,
713	2024.
714	Hua Wei, Nan Xu, Huichu Zhang, Guanjie Zheng, Xinshi Zang, Chacha Chen, Weinan Zhang,
715	Yanmin Zhu, Kai Xu, and Zhenhui Li. Colight: Learning network-level cooperation for traffic
716	signal control. In CIKM, pp. 1913–1922, 2019.
717	Tianhao Xie Sibeng Zhao, Chen Henry Wu, Yitao Liu, Oian Luo, Victor Zhong, Yanchao Yang, and
718	Tao Yu. Text2reward: Automated dense reward function generation for reinforcement learning.
719	arXiv preprint arXiv:2309.11489, 2023.
720	
721	Xianliang Yang, Zhihao Liu, Wei Jiang, Chuheng Zhang, Li Zhao, Lei Song, and Jiang Bian. A
722	versatile multi-agent reinforcement learning benchmark for inventory management. arXiv preprint
723	<i>urxiv:2500.07342, 2025.</i>
724	Shunyu Yao, Jeffrey Zhao, Dian Yu, Nan Du, Izhak Shafran, Karthik R Narasimhan, and Yuan Cao.
725	React: Synergizing reasoning and acting in language models. In ICLR, pp. 1–33, 2023.
720	Hooren Va, Jiami Wang, Zhiguang Cao, and Guojia Song, Paavo: Larga languaga models as
728	hyper-heuristics with reflective evolution. arXiv preprint arXiv:2402.01145, 2024.
729	Donghao Ying, Yunkai Zhang, Yuhao Ding, Alec Koppel, and Javad Lavaei. Scalable primal-dual
730	actor-critic method for safe multi-agent RL with general utilities. In NeurIPS, 2023.
731	Cheo Yu Akesh Valu Eugana Vinitsky Jiayuan Cao, Yu Wang Alayandra Payan and Yi Wu, Tha
732	surprising effectiveness of ppo in cooperative multi-agent games. In <i>NeurIPS</i> , pp. 24611–24624
733	2022.
734	
736	Haoqi Yuan, Chi Zhang, Hongcheng Wang, Feiyang Xie, Penglin Cai, Hao Dong, and Zongqing
737	Lu. Plan4mc: Skill reinforcement learning and planning for open-world minecraft tasks. <i>arXiv</i>
738	preprint urxiv:2505.10505, 2025.
739	Tom Zahavy, Matan Haroush, Nadav Merlis, Daniel J Mankowitz, and Shie Mannor. Learn what not
740	to learn: Action elimination with deep reinforcement learning. In NeurIPS, 2018.
741	Bin Zhang, Hangyu Mao, Jingging Ruan, Ying Wen, Yang Li, Shao Zhang, Zhiwei Xu, Dapeng
742	Li, Ziyue Li, Rui Zhao, et al. Controlling large language model-based agents for large-scale
743	decision-making: An actor-critic approach. arXiv preprint arXiv:2311.13884, 2023.
744	Vising Zhang Zhang Wang and Tama Dava M 14' and a' Gamma (1997). A 1997
745	Kalqing Zhang, Zhuoran Yang, and Tamer Başar. Multi-agent reinforcement learning: A selective
746	321_384 2021
747	<i>521 501, 2021.</i>
748	Xizhou Zhu, Yuntao Chen, Hao Tian, Chenxin Tao, Weijie Su, Chenyu Yang, Gao Huang, Bin Li,
749	Lewei Lu, Xiaogang Wang, et al. Ghost in the minecraft: Generally capable agents for open-world
750	enviroments via large language models with text-based knowledge and memory. <i>arXiv preprint</i>
751	arxiv:2505.1/144, 2025.
752	Alessandro Zocca. Temporal starvation in multi-channel csma networks: an analytical framework.
103	ACM SIGMETRICS Performance Evaluation Review, 46(3):52–53, 2019.
104	

PROOFS А

Proposition 1.

 1. For any $E \in \{E\}$, $\{\pi_E\} \subseteq \{\pi\}$. If E is non-trivial, then $\{\pi_E\} \subset \{\pi\}$.

2. For any $\pi \in \{\pi\}$, there exists a non-trivial $E \in \{E\}$ such that $J(\pi_E) \ge J(\pi)$.

Proof. We begin by offering the proof of the first statement in Proposition 1. We denote $\{\pi^k\}$ as the set of all possible policies for agent k, with each π^k satisfying the following two conditions:

$$\begin{aligned} & \mathcal{T}^{k} : \mathcal{O}^{k} \times \mathcal{A}^{k} \to [0,1] \\ & \sum_{a^{k} \in \mathcal{A}^{k}} \pi^{k}(a^{k} \mid o^{k}) = 1 \end{aligned}$$
 (5)

Let $\{\pi_E^k\}$ be the set of policies under an exploration function E. For every element in $\pi_E^k \in \{\pi_E^k\}$, one of the two situations exists:

1. If E is trivial, then $\pi_E^k(\cdot \mid o^k) = \pi^k(\cdot \mid o^k)$, hence $\pi_E^k \in \{\pi^k\}$.

2. If E is non-trivial, then
$$\pi_E^k(\cdot \mid o^k) = \frac{\pi^k(\cdot \mid o^k) \cdot E(o^k, \cdot)}{\sum_{a^k \in \mathcal{A}^k} \pi^k(a^k \mid o^k) \cdot E(o^k, a^k)}$$
. It is clear that $0 \le \pi_E^k(\cdot \mid o^k) \le 1$ and $\sum_{a^k \in \mathcal{A}^k} \pi_E^k(a^k \mid o^k) = 1$, thus $\pi_E^k \in \{\pi^k\}$.

Therefore, we know that:

$$\pi_E^k\} \subseteq \{\pi_k\}.\tag{6}$$

Given that for $\forall \pi(\cdot \mid s_t) = \prod_{k=1}^N \pi^k(\cdot \mid o^k)$, we have $\pi \in \{\pi\}$, where each π^k belongs to $\{\pi^k\}$. According to Formula 6, it is known that for $\forall \pi_E(\cdot \mid s_t) = \prod_{k=1}^N \pi_E^k(\cdot \mid o^k)$, it belongs to $\{\pi\}$. Then:

{

$$\{\boldsymbol{\pi}_E\} \subseteq \{\boldsymbol{\pi}\}.\tag{7}$$

When E is non-trivial, $\pi_E^k \in \{\pi^k\}$ still holds, but $\pi^k \in \{\pi_E^k\}$ may not be true (i.e., when $\pi^k(a^k \mid a^k)$) o^k) > 0 for $\forall a^k \in \mathcal{A}, \pi^k \notin \{\pi_E^k\}$). Hence we can get:

$$\{\pi_E^k\} \subset \{\pi_k\}.\tag{8}$$

In a similar manner, we can deduce that if there exists a $k \in [1, 2, ..., N]$ such that $\pi^k \notin \{\pi_E^k\}$, then $\pi \notin \{\pi_E\}$, which means:

$$\boldsymbol{\pi}_E\} \subset \{\boldsymbol{\pi}\}.\tag{9}$$

Therefore, we finish the proof of the first statement.

To proof the second statement, it is necessary to introduce a series of variables. We define the value function and state-action function for π as follows: $V_{\pi}(s) = \mathbb{E}_{\mathbf{a}_{0:\infty} \sim \pi, s_{1:\infty} \sim P} \left[\sum_{t=0}^{\infty} \gamma^t \mathbf{r}_t \mid \mathbf{s}_0 = \mathbf{s}\right]$ and $Q_{\pi}(s, a) = \mathbb{E}_{\mathbf{a}_{1:\infty} \sim \pi, s_{1:\infty} \sim P} \left[\sum_{t=0}^{\infty} \gamma^t \mathbf{r}_t \mid \mathbf{s}_0 = \mathbf{s}, \mathbf{a}_0 = \mathbf{a}\right]$. The advantage function is defined as $A_{\pi}(s, \mathbf{a}) = Q_{\pi}(s, \mathbf{a}) - V_{\pi}(s)$. The joint exploration function is introduced as $\mathbf{E}(s, \cdot) = \mathbf{T}^N$. $\prod_{k=1}^{N} E(o^k, \cdot)$. The relationship between $V_{\pi}(s)$ and $Q_{\pi}(s, \mathbf{a})$ can be formulated as:

$$V_{\boldsymbol{\pi}}(s) = \sum_{\mathbf{a} \in \boldsymbol{\mathcal{A}}^k} \boldsymbol{\pi}(\mathbf{a} \mid s) Q_{\boldsymbol{\pi}}(s, \mathbf{a})$$
(10)

For a non-trivial E, the value function of π_E can be written as:

$$V_{\boldsymbol{\pi}_{E}}(s) = \sum_{\mathbf{a} \in \boldsymbol{\mathcal{A}}^{k}} \frac{\mathbf{E}(s, \mathbf{a})\boldsymbol{\pi}(\mathbf{a} \mid s)}{\sum_{\mathbf{a} \in \boldsymbol{\mathcal{A}}^{k}} \mathbf{E}(s, \mathbf{a})\boldsymbol{\pi}(\mathbf{a} \mid s)} Q_{\boldsymbol{\pi}}(s, \mathbf{a})$$
$$= \frac{1}{\sum_{\mathbf{a} \in \boldsymbol{\mathcal{A}}^{k}} \mathbf{E}(s, \mathbf{a})\boldsymbol{\pi}(\mathbf{a} \mid s)} \sum_{\mathbf{a} \in \boldsymbol{\mathcal{A}}^{k}} \mathbf{E}(s, \mathbf{a})\boldsymbol{\pi}(\mathbf{a} \mid s)Q_{\boldsymbol{\pi}}(s, \mathbf{a})$$
(11)

$$= \frac{1}{\sum_{\mathbf{a}\in\mathcal{A}^k} \mathbf{E}(s,\mathbf{a})\pi(\mathbf{a}\mid s)} \left[V_{\pi}(s) - \sum_{\mathbf{a}\in\mathcal{A}^k} \left(1 - \mathbf{E}(s,\mathbf{a})\right)\pi(\mathbf{a}\mid s)Q_{\pi}(s,\mathbf{a}) \right].$$

Thus, we have:

$$V_{\boldsymbol{\pi}_{E}}(s) - V_{\boldsymbol{\pi}}(s) = \frac{V_{\boldsymbol{\pi}}(s) - \sum_{\mathbf{a} \in \mathcal{A}^{k}} (1 - \mathbf{E}(s, \mathbf{a})) \, \boldsymbol{\pi}(\mathbf{a} \mid s) Q_{\boldsymbol{\pi}}(s, \mathbf{a})}{\sum_{\mathbf{a} \in \mathcal{A}^{k}} \mathbf{E}(s, \mathbf{a}) \boldsymbol{\pi}(\mathbf{a} \mid s)} - V_{\boldsymbol{\pi}}(s)$$

$$= \frac{(1 - \sum_{\mathbf{a} \in \mathcal{A}^{k}} \mathbf{E}(s, \mathbf{a}) \boldsymbol{\pi}(\mathbf{a} \mid s)) V_{\boldsymbol{\pi}}(s) - \sum_{\mathbf{a} \in \mathcal{A}^{k}} (1 - \mathbf{E}(s, \mathbf{a})) \, \boldsymbol{\pi}(\mathbf{a} \mid s) Q_{\boldsymbol{\pi}}(s, \mathbf{a})}{\sum_{\mathbf{a} \in \mathcal{A}^{k}} \mathbf{E}(s, \mathbf{a}) \boldsymbol{\pi}(\mathbf{a} \mid s)}$$

$$= \frac{\sum_{\mathbf{a} \in \mathcal{A}^{k}} (1 - \mathbf{E}(s, \mathbf{a})) \, \boldsymbol{\pi}(\mathbf{a} \mid s) V_{\boldsymbol{\pi}}(s) - \sum_{\mathbf{a} \in \mathcal{A}^{k}} (1 - \mathbf{E}(s, \mathbf{a})) \, \boldsymbol{\pi}(\mathbf{a} \mid s) Q_{\boldsymbol{\pi}}(s, \mathbf{a})}{\sum_{\mathbf{a} \in \mathcal{A}^{k}} \mathbf{E}(s, \mathbf{a}) \boldsymbol{\pi}(\mathbf{a} \mid s)}$$

$$= -\frac{\sum_{\mathbf{a} \in \mathcal{A}^{k}} (1 - \mathbf{E}(s, \mathbf{a})) \, \boldsymbol{\pi}(\mathbf{a} \mid s) A_{\boldsymbol{\pi}}(s, \mathbf{a})}{\sum_{\mathbf{a} \in \mathcal{A}^{k}} \mathbf{E}(s, \mathbf{a}) \boldsymbol{\pi}(\mathbf{a} \mid s)}.$$
(12)

When $-\sum_{\mathbf{a}\in\mathcal{A}^k} (1 - \mathbf{E}(s, \mathbf{a})) \pi(\mathbf{a} \mid s) A_{\pi}(s, \mathbf{a}) \ge 0$, which means the expectation of the advantage value for the pruned actions is less than or equal to 0, then $V_{\pi_E}(s) \ge V_{\pi}(s)$. Because for every $s \in S$, $\sum_{\mathbf{a} \in \mathcal{A}^k} A_{\pi}(s, \mathbf{a}) = \sum_{\mathbf{a} \in \mathcal{A}^k} Q_{\pi}(s, \mathbf{a}) - V_{\pi}(s) = 0$, there always exist actions for which the advantage function values are less than or equal to zero.

As $J(\pi_E) - J(\pi) = \mathbb{E}_{s_0 \sim \rho^0} [V_{\pi_E}(s_0) - V_{\pi}(s_0)]$, if an exploration function E can satisfy the condition that for all states $s \in S$, the inequality $-\sum_{\mathbf{a} \in \mathcal{A}^k} (1 - \mathbf{E}(s, \mathbf{a})) \pi(\mathbf{a} \mid s) A_{\pi}(s, \mathbf{a}) \ge 0$ holds, then it can be guaranteed that $J(\boldsymbol{\pi}_E) \geq J(\boldsymbol{\pi})$.

Therefore, we finish the proof of the second statement.

Proposition 2. Let's consider a fully cooperative game with N agents, one state, and the joint action space $\{0,1\}^N$, where the reward is given by $r(\mathbf{0}^0, \mathbf{1}^N) = r_1$ and $r(\mathbf{0}^{N-m}, \mathbf{1}^m) = -mr_2$ for all $m \neq N$, r_1 , r_2 are positive real numbers. We suppose the initial policy is randomly and uniformly initialized, and the policy is optimized in the form of gradient descent. Let p be the probability that the shared policy converges to the best policy, then:

$$p = 1 - \sqrt[N-1]{\frac{r_2}{r_1 + Nr_2}}.$$
(13)

Proof. Clearly, the best policy is the deterministic policy with joint action $(\mathbf{0}^0, \mathbf{1}^N)$.

Now, let the shared policy be $(1 - \theta, \theta)$, where θ is the probability that an agent takes action 1. The expected reward can be written as:

$$J(\theta) = \mathbf{Pr}\left(\boldsymbol{a}^{1:N} = (\mathbf{0}^0, \mathbf{1}^N)\right) \cdot r_1 - \sum_{t=0}^{N-1} \mathbf{Pr}\left(\boldsymbol{a}^{1:N} = (\mathbf{0}^{N-t}, \mathbf{1}^t)\right) \cdot t \cdot r_2$$

(

$$= \theta^N \cdot r_1 - \sum_{t=0}^{N-1} t \cdot C_N^t \theta^t (1-\theta)^{N-t} \cdot r_2$$

$$= \theta^N \cdot r_1 - \sum_{t=0}^N t \cdot C_N^t \theta^t (1-\theta)^{N-t} \cdot r_2 + N \cdot \theta^N \cdot r_2,$$
(14)

where C_N^t is a combinatorial number. We need to simplify $\sum_{t=0}^N t \cdot C_N^t \theta^t (1-\theta)^{N-t}$ for further analysis. Notice the structural similarity between the results and the binomial theorem:

$$(1-\theta) + \theta)^{N} = \sum_{t=0}^{N} C_{N}^{t} \theta^{t} (1-\theta)^{N-t}.$$
(15)

We take the derivative of θ on both sides of Formula 15. Because the left side is constant, its derivative is 0. Then:

 $0 = \frac{d \sum_{t=0}^{N} C_{N}^{t} \theta^{t} (1-\theta)^{N-t}}{d\theta}$ $= \sum_{t=0}^{N} C_{N}^{t} t \cdot \theta^{t-1} \cdot (1-\theta)^{N-t} + C_{N}^{t} (N-t) \cdot (-1) \cdot (1-\theta)^{N-t-1} \cdot \theta^{t}$ $= \sum_{t=0}^{N} C_{N}^{t} (1-\theta)^{N-t-1} \theta^{t-1} ((1-\theta)t - (N-t)\theta)$ $= \sum_{t=0}^{N} C_{N}^{t} (1-\theta)^{N-t-1} \theta^{t-1} (t-N\theta).$ (16)

Thus, we have:

$$N\theta \sum_{t=0}^{N} C_{N}^{t} (1-\theta)^{N-t-1} \theta^{t-1} = \sum_{t=0}^{N} t C_{N}^{t} (1-\theta)^{N-t-1} \theta^{t-1}$$

$$N\theta \sum_{t=0}^{N} C_{N}^{t} (1-\theta)^{N-t} \theta^{t} = \sum_{t=0}^{N} t C_{N}^{t} (1-\theta)^{N-t} \theta^{t}.$$
(17)

Notice that the left side of the equation is the expansion form of Formula 15, and the right side of the equation is the desired Formula, we can get:

$$\sum_{t=0}^{N} t C_N^t (1-\theta)^{N-t} \theta^t = N\theta.$$
(18)

Bring Formula 18 back to Formula 14, we get:

$$J(\theta) = \theta^N \cdot r_1 - N\theta r_2 + N \cdot \theta^N \cdot r_2.$$
⁽¹⁹⁾

In order to maximise $J(\theta)$, we must maximise $\theta^N \cdot (r_1 + Nr_2) - N\theta r_2$. Since the policy optimization usually adopts a gradient manner, we calculate the derivative of Formula 19 with respect to θ as:

$$\frac{dJ(\theta)}{d\theta} = N\theta^{N-1}(r_1 + Nr_2) - Nr_2.$$
(20)

the point
$$\theta^* = \sqrt[N-1]{\frac{r_2}{r_1 + Nr_2}}$$
 is the only zero of $\frac{dJ(\theta)}{d\theta}$. When $\theta \le \theta^*$, $\frac{dJ(\theta)}{d\theta} \le 0$; $\theta \ge \theta^*$, $\frac{dJ(\theta)}{d\theta} \ge 0$.

Remember we are trying to maximize $J(\theta)$ through a gradient way, and then the policy improves the parameters in the direction of the gradient. As the initial policy is randomly and uniformly initialized, the θ is uniformly distributed in the interval [0,1], then the probability that the shared policy converges to the best policy is:

$$p = 1 - \sqrt[N-1]{\frac{r_2}{r_1 + Nr_2}}.$$
(21)

Therefore, we finish the proof of Proposition 2.

B PSEUDOCODE

In this Section, we give the pseudocode of our proposed eSpark. We denote the performance of policy *i* as G_i , and the pseudocode of eSpark is shown in Algorithm 1.

- C DETAILED SETTINGS
- 915 C.1 MABIM DETAILS
- 917 MABIM is a simulation environment dedicated to leveraging MARL to tackle the challenges inherent in inventory management problems. Within MABIM, each stock SKU at every echelon is represented



as an autonomous agent. The decision-making process of each agent reflects the procurement requirements for the specific SKU at its corresponding echelon.

Each time step involves the agent making decisions regarding replenishment quantities for SKUs and subsequently transitioning the environment to a new state. Let $M \in \mathbb{Z}^+$ be the total warehouses, with the first one being closest to customers, and $N \in \mathbb{Z}^+$ the total SKUs. Given a variable $X \in \{D, S, L...\}, X_{i,j}^t$ represents its value for the *j*-th SKU in the *i*-th echelon at step *t*, with $0 \le i < M$ and $0 \le j < N$. Given the above notations, the main progression of a step can be described as follows:

$$D_{i+1,j}^{t+1} = R_{i,j}^t$$
(Replenish)

$$S_{i,j}^t = \min(D_{i,j}^t, I_{i,j}^t)$$
(Sell)

$$A_{i,j}^t = \sum_{i=1}^{t-1} \mathbb{I}(k + L_{i,j}^k ==t) \cdot S_{i+1,j}^t$$
(Arrive)

$${}^{t}_{i,j} = \sum_{k=0}^{t} \mathbb{I}(k + L^{k}_{i,j} == t) \cdot S^{t}_{i+1,j}$$
 (And

$$\gamma_{i}^{t} = \min\left(\frac{W_{i} - \sum_{j} I_{i,j}^{t}}{\sum_{j} A_{i,j}^{t}}, 1\right), B_{i,j}^{t} = \lfloor A_{i,j}^{t} \cdot \gamma_{i}^{t} \rfloor \quad (\text{Receive})$$
$$I_{i,j}^{t+1} = I_{i,j}^{t} - S_{i,j}^{t} + B_{i,j}^{t} \quad (\text{Update})$$

945

946

963 Here, $D, R, S, I, A, B \in \mathbb{Z}^+$ and $\mathbb{I}(\text{condition})$ is an indicator function that returns 1 if the condition 964 is true, and 0 otherwise. For the topmost echelon, orders are channeled to a super vendor capable of 965 fulfilling all order demands at that level. Orders from other echelons are directed to their immediate 966 upstream echelons, where the demands are satisfied based on the inventory levels of the upper 967 echelons. The demand at the bottom echelon is derived from actual customer orders captured within 968 real-world data sets. The reward function within MABIM is meticulously calibrated based on the economic realities of inventory management, integrating five fundamental elements: sales profit, 969 order cost, holding cost, backlog cost, and excess cost. The summation of these elements constitutes 970 the reward value, thereby incentivizing agents to optimize inventory control for enhanced profitability 971 and operational efficacy.

MABIM incorporates challenges across five key categories: Scaling up, Cooperation, Competition, Generalization, and Robustness. We concentrate on the challenges associated with Scaling up, Cooperation, and Competition, as these challenges not only manifest in inventory management problems but also exist in a broad range of MARL tasks. We catalog the number of agents, challenges and degrees of difficulty within all the experimental scenarios in Table 9. The specific setting of each scenario is given in Table 10:

Table 9: Tasks and corresponding challenges. '+' denotes the extent of the challenges.					
Task name	Agents number	Scaling up	Challenge Cooperation	Competition	
Standard (100 SKUs)	100				
2 echelons (100 SKUs)	200		+		
3 echelons (100 SKUs)	300		++		
Lower capacity(100 SKUs)	100			+	
Lowest capacity (100 SKUs)	100			++	
Standard (500 SKUs)	500	+			
2 echelons (500 SKUs)	1000	+	+		
3 echelons (500 SKUs)	1500	+	++		
Lower capacity (500 SKUs)	500	+		+	
Lowest capacity (500 SKUs)	500	+		++	

Table 10: Experiments settings. "#SKU * N" indicates N times the number of SKUs.

Task name	#Echelon	#SKU	Capacity per echelon
Standard (100 SKUs)	1	100	#SKU*100
2 echelons (100 SKUs)	2	100	#SKU*100
3 echelons (100 SKUs)	3	100	#SKU*100
Lower capacity(100 SKUs)	1	100	#SKU*50
Lowest capacity (100 SKUs)	1	100	#SKU*25
Standard (500 SKUs)	1	500	#SKU*100
2 echelons (500 SKUs)	2	500	#SKU*100
3 echelons (500 SKUs)	3	500	#SKU*100
Lower capacity (500 SKUs)	1	500	#SKU*50
Lowest capacity (500 SKUs)	1	500	#SKU*25

For each scenario, we carry out three independent runs. Performance is reported as average test set profits from the top model in each run.

1012 C.2 SUMO DETAILS

SUMO is an open-source, highly portable, microscopic and continuous road traffic simulation package
designed to handle large road networks. In the SUMO simulation environment, each intersection is
conceptualized as an autonomous agent equipped with an array of predefined traffic signal phases.
These phases orchestrate the traffic flow across the intersection's multiple approaches. The selection
of these phases, driven by the assessment of live traffic conditions, is aimed at attenuating road
congestion and enhancing the fluidity of vehicular movement through the network, thus contributing
to the overall efficiency of urban traffic management.

To conduct a thorough evaluation of each algorithm, we select a total of five scenarios from both
synthetic and real-world datasets. These datasets encompass a diverse array of intersections, varying
in number and type. The intersections are classified according to their configuration into three
categories: bi-directional (2-arm), tri-directional (3-arm), and quadri-directional (4-arm), indicating
the number of exit points at each junction. We summarize the type of each dataset, the number of
intersections included, and the classification of these intersections in Table 11.

1028	included.						
1029	Dataset	Category	Intersections number	2-arm	3-arm	4-arm	
1030	Grid 4×4	synthesis	16	0	0	16	
1031	Arterial 4×4	synthesis	16	0	0	16	
1032	Grid 5×5 Lu et al. (2023)	synthesis	25	0	0	25	
1033	Cologne8	real-world	8	1	3	4	

real-world

Table 11: The categories of each SUMO dataset, along with the number and types of intersections

To facilitate a homogeneous observation and action space conducive to the deployment of various MARL algorithms, we employ the GEneral Scenario-Agnostic (GESA) framework to parse each intersection into a standardized 4-arm intersection with eight potential actions.

C.3 MODEL DETAILS

Ingolstadt21

We employ IPPO as the base MARL algorithm for eSpark due to its ability to scale to large-scale MARL challenges. We select GPT-4 as our LLM_c and LLM_q , specifically opting for the 2023-09-01-preview version. The temperature of GPT-4 is set to 0.7, with no frequency penalty and presence penalty. For each scenario, we conduct three runs, setting the batch size for each generation of exploration functions to K = 16. This batch size is chosen because it guarantees that the initial generation contains at least one executable exploration function for our environment. We limit the number of training iterations to 10, as we observe that the performance for most scenarios tends to converge within this number of iterations.

D **BASELINE DETAILS**

 In our experiments, we employed three categories of baselines: OR baselines, MARL baselines and pruning baselines. In Table D, we list the characteristics and environment of all the baselines utilized in our study. In the rest of this section, we will elucidate the underlying principles of each OR baseline, articulate the design of the pruning baselines, and present the hyperparameter for the MARL baselines.

Table 12: All the	baselines	used in	the experiments
-------------------	-----------	---------	-----------------

Algorithm name	OR baseline	MARL CTDE	baseline DTDE	Pruning baseline	Used envi MABIM	ironment SUMO
Base stock (BS) Arrow et al. (1951)	\checkmark				\checkmark	
(S, s) Blinder (1990)	\checkmark				\checkmark	
FTC Roess et al. (2004)	\checkmark					\checkmark
MaxPressure Kouvelas et al. (2014)	\checkmark					\checkmark
IPPO			\checkmark		\checkmark	\checkmark
QTRAN		\checkmark			\checkmark	
QPLEX		\checkmark			\checkmark	
MAPPO		\checkmark			\checkmark	\checkmark
MPLight Chen et al. (2020)			\checkmark			\checkmark
CoLight Wei et al. (2019)		\checkmark				\checkmark
Ramdom pruning				\checkmark	\checkmark	\checkmark
(S,s) pruning				\checkmark	\checkmark	
Upbound pruning				\checkmark	\checkmark	
MaxPressure pruning				\checkmark		\checkmark

1080 D.1 OR BASELINES

1082 D.1.1 BASE STOCK ALGORITHM

The base stock algorithm constitutes a streamlined and efficacious approach for inventory control, whereby replenishment orders are initiated upon inventory below a predefined threshold level. This policy is traditionally acknowledged as a fundamental benchmark, favored for its straightforwardness and rapid implementation. The computation of the base stock level is facilitated through a programmatic methodology, as explicated in Equation 22:

1089

1090

1093 1094

1095

$$\begin{split} \max \quad o_{i,j}^{t} &= \bar{p}_{i,j} \cdot S_{i,j}^{t} - \bar{c}_{i,j} \cdot S_{i+1,j}^{t} - \bar{h}_{i,j} \cdot I_{i,j}^{t+1} - \bar{c}_{i,j} \cdot T_{i,j}^{0} - \bar{c}_{i,j} \cdot I_{i,j}^{0} \\ \text{s.t} \quad I_{i,j}^{t+1} &= I_{i,j}^{t} + S_{i+1,j}^{t-\bar{L}_{i,j}} - S_{i,j}^{t} \\ T_{i,j}^{t+1} &= T_{i,j}^{t} - S_{i+1,j}^{t-\bar{L}_{i,j}} + S_{i+1,j}^{t} \\ S_{i,j}^{t} &= \min(I_{i,j}^{t}, R_{i,j}^{t}) \\ T_{i,j}^{0} &= \sum_{t=-\bar{L}_{i,j}}^{-1} S_{i+1,j}^{t} \\ z_{i,j} &= I_{i,j}^{t+1} + S_{i+1,j}^{t} + T_{i,j}^{t} \end{split}$$
(22)

- 1100 $z_{i,j} \in \mathbb{R}^+.$
- 1102

1099

In the above equations, i, j, and t are indexes for the warehouse, SKU, and discrete time, respectively. The indicators $\bar{p}, \bar{c}, \bar{h}$, and \bar{L} represent the average selling price, cost of procurement, cost of holding, and lead time. The variables S, R, I, and T denote the quantities associated with sales, orders for replenishment, inventory in stock and inventory in transit. $o_{i,j}^t$ describes the profit objective, while $z_{i,j}$ is indicative of the base stock level.

We utilize two approaches for computing stock levels. The first approach, named **BS static**, involves calculating all base stock levels with historical data from the training set, which are then applied consistently to the test set. The levels remain unchanged during the test period. The second approach, termed as **BS dynamic**, computes stock levels directly on the test set relying on historical data and updates on a regular basis.

1113 D.1.2 (S, s) Algorithm

The (S, s) inventory policy serves as a robust framework for managing stock levels. Under this policy, a restocking order is triggered once the inventory count falls below a predefined threshold, identified as s. The objective of this replenishment is to elevate the stock to its upper limit, designated as S. Empirical analyses have substantiated the efficacy of this protocol in optimizing inventory control processes. As a result, it is adopted as a benchmark heuristic, with the aim of algorithmically ascertaining the most efficacious (S, s) parameter for each discrete SKU in the given inventory dataset. In our implementation, we conduct a search on the training set to identify the optimal values of s and S, after which we apply these values consistently to the test set.

1122

1114

1123 D.1.3 FIXED-TIME CONTROL ALGORITHM

The Fixed-Time Control (FTC) algorithm is a traditional traffic signal control strategy predicated on predefined signal plans. These plans are typically designed based on historical traffic flow patterns and do not adapt to real-time traffic conditions. The FTC operates on a static schedule where the signal phases at intersections change at fixed intervals. This approach is straightforward and easy to implement but may not be optimal under variable traffic conditions due to its lack of responsiveness to dynamic traffic demands.

In our implementation, the FTC follows a fixed sequence of signal phases: 'WT-ET', 'NT-ST',
'WL-EL', 'NL-SL', 'WL-WT', 'EL-ET', 'SL-ST', 'NL-NT'. Here, 'W', 'E', 'N', and 'S' denote
westbound, eastbound, northbound, and southbound traffic, respectively, while 'T' indicates through
movement, and 'L' signifies a left turn. Each phase has a duration of 30 seconds

1134 D.1.4 MAXPRESSURE ALGORITHM 1135

1136 The MaxPressure algorithm represents a more advanced traffic signal control strategy that dynamically adjusts signal phases in response to real-time traffic conditions. It calculates the "pressure" at each 1137 intersection, defined as the difference between the number of vehicles on the incoming and outgoing 1138 lanes. The algorithm aims to optimize traffic flow by selecting signal phases that reduce the maximum 1139 pressure across the network, thus alleviating congestion and enhancing network throughput. Unlike 1140 FTC, MaxPressure is adaptive and can continuously optimize signal timing based on the current 1141 traffic state, making it more suitable for managing fluctuating traffic volumes. 1142

1143 1144

D.2 HYPERPARAMETERS SETTINGS FOR MARL BASELINES

1145 The following table enumerates the hyperparameters employed during the training process for all 1146 MARL baselines. For all test scenarios, training is performed with a uniform suite of hyperparameters 1147 that have not undergone specialized fine-tuning.

1148 1149

1150

Table 13: Hyperparameters of MARL Algorithms Used in MABIM and SUMO Environments. '-' indicates that the algorithm is not set or does not contain this hyperparameter.

1151 indicates that the algorithm is not set of does not contain this hyperparameter.										
1152	Uuporporemeter	MABIM environment				SU	SUMO environment			
1152	Hyperparameter	IPPO	QTRAN	QPLEX	MAPPO	IPPO	CoLight	MPLight		
1155	Training steps	5020000	5020000	5020000	5020000	2400000	2400000	2400000		
1154	Discount rate	0.985	0.985	0.985	0.985	0.985	0.9	0.9		
1155	Optimizer	Adam	Adam	Adam	Adam	Adam	RMSProp	RMSProp		
1156	Optimizer alpha	0.99	0.99	0.99	0.99	0.99	0.95	0.95		
1157	Optimizer eps	1e-5	1e-5	1e-5	1e-5	1e-5	1e-7	1e-7		
4450	Learning rate	5e-4	5e-4	5e-4	5e-4	5e-4	1e-3	1e-3		
1158	Grad norm clip	10	10	10	10	10	-	-		
1159	Eps clip	0.2	-	-	0.2	0.2	-	-		
1160	Critic coef	0.5	-	-	0.5	0.5	-	-		
1161	Entropy coef	0	-	-	0	0	-	-		
1101	Accumulated episodes	4	8	8	4	4	10	10		
1162	Number of neighbors	-	-	-	-	-	5	-		

1163 1164

D.3 PRUNING BASELINES 1165

1166 D.3.1 RANDOM PRUNING 1167

1168 Random pruning is implemented by randomly masking a certain percentage of the available actions. 1169 During the action selection process, each agent will have p percent of its available actions randomly masked. To balance the observability of the pruning's impact with the preservation of the algorithm's 1170 capacity to utilize prior experience, we set p = 0.3. 1171

1172

D.3.2 (S, s) pruning 1173

1174 According to the (S, s) algorithm, for a given SKU, a replenishment quantity of $\Delta = S - s$ is 1175 ordered when the current inventory level falls below the threshold s; otherwise, no order is placed. 1176 We extend the reference replenishment quantity Δ to a range $[r_1 \times \Delta, r_2 \times \Delta]$, where r_1, r_2 are both 1177 real numbers and $0 \le r_1 \le 1$ and $r_2 \ge 1$. Actions within this interval are deemed available, while 1178 those outside of this range are masked. In our implementation, we select $r_1 = 0.5$ and $r_2 = 2$.

1179

1180 D.3.3 MAXPRESSURE PRUNING

1181 The MaxPressure pruning method utilizes the heuristic concept of "pressure" at an intersection to 1182 prune actions. We calculate the pressure associated with each action, and these pressures are then 1183 ranked. The actions with the top-k highest pressures are rendered available for selection. Actions not 1184 meeting this threshold are subsequently masked. 1185

A standardized intersection warped through the GESA is modeled as a four-arm intersection compris-1186 ing eight potential actions. We empirically set k = 4 to ensure effective pruning while maintaining a 1187 sufficient number of available actions.

Ε ADDITIONAL RESULTS

E.1 ADDITIONAL MAIN RESULTS

In this section, we present the complete main experimental results for SUMO in Section 5.2, and give the consumption of GPT tokens here.

Table 14: Detailed performance in SUMO, includes the mean and standard deviation.

Tuble 11. Detailed performance in Service, menudes the mean and standard deviation.							
Method	Metric	Grid 4×4	Arterial 4×4	Grid 5×5	Cologne8	Ingolstadt21	
FTC	Delay Trip time Wait time	$\begin{array}{c} 161.14 {\pm}~3.77 \\ 291.48 {\pm}4.45 \\ 155.66 {\pm}3.42 \end{array}$	$\begin{array}{c} 1229.68{\pm}16.79\\ 676.77{\pm}13.70\\ 521.86{\pm}13.33\end{array}$	$\begin{array}{c} 820.88{\pm}24.36\\ 584.54{\pm}24.17\\ 441.63{\pm}21.13\end{array}$	$\begin{array}{c} 85.27{\pm}1.21\\ 145.93{\pm}0.84\\ 58.92{\pm}0.68\end{array}$	224.96±11.91 352.06±9.29 161.22±7.88	
MaxPressure	Delay Trip time Wait time	63.39 ± 1.34 174.68 ± 2.05 37.37 ± 1.06	$\begin{array}{c} 995.23{\pm}77.02 \\ 702.09{\pm}23.61 \\ 511.06{\pm}22.55 \end{array}$	$\begin{array}{c} 242.85{\pm}16.04\\ 269.35{\pm}9.62\\ 114.36{\pm}6.48\end{array}$	$\begin{array}{c} 31.63 \pm 0.61 \\ 95.67 {\pm} 0.62 \\ 11.03 {\pm} 0.28 \end{array}$	$\begin{array}{c} 228.64{\pm}15.83\\ 352.30\pm15.06\\ 159.44{\pm}13.34\end{array}$	
IPPO	Delay Trip time Wait time	$\begin{array}{c} 48.40{\pm}0.45\\ 160.12{\pm}0.60\\ 22.69{\pm}0.38\end{array}$	884.73±38.94 506.18±10.39 435.44±77.54	$\begin{array}{c} 228.78{\pm}11.59\\ 238.03{\pm}7.10\\ 91.84{\pm}6.31 \end{array}$	27.60 ± 1.70 91.41 ± 1.60 7.70 ± 0.82	$\begin{array}{r} 342.97{\pm}43.61 \\ 464.50{\pm}43.30 \\ 267.51{\pm}40.53 \end{array}$	
МАРРО	Delay Trip time Wait time	$\begin{array}{c} 51.25{\pm}0.58\\ 160.01{\pm}0.54\\ 25.41{\pm}0.54\end{array}$	$\begin{array}{c} 958.43 {\pm} 181.72 \\ 757.40 {\pm} 242.00 \\ 609.80 {\pm} 255.22 \end{array}$	$\begin{array}{c} 221.62{\pm}20.73\\ 247.56{\pm}3.71\\ 97.10{\pm}5.22 \end{array}$	32.55 ± 4.66 94.31 ± 1.77 9.39 ± 1.53	$\begin{array}{r} 347.59{\pm}47.59\\ 480.66{\pm}49.46\\ 283.59{\pm}43.20 \end{array}$	
MPLight	Delay Trip time Wait time	$\begin{array}{c} 63.51{\pm}0.64\\ 172.47{\pm}0.89\\ 40.32{\pm}0.96\end{array}$	$\begin{array}{c} 1142.98{\pm}79.65\\ 583.21{\pm}45.84\\ 366.27{\pm}58.03\end{array}$	$\begin{array}{c} 223.44{\pm}16.18\\ 255.49{\pm}6.26\\ 126.42{\pm}5.31\end{array}$	37.93 ± 0.45 110.56 ± 1.18 12.98 ± 0.57	$\begin{array}{c} 196.74{\pm}9.88\\ 331.42{\pm}11.79\\ 126.09{\pm}13.60\end{array}$	
CoLight	Delay Trip time Wait time	$\begin{array}{c} 53.40{\pm}1.89\\ 165.77{\pm}1.89\\ 27.25{\pm}1.64\end{array}$	$\begin{array}{c} 820.67{\pm}58.65\\ 470.33{\pm}12.34\\ 312.47{\pm}16.63\end{array}$	$\begin{array}{c} 339.66{\pm}48.55\\ 305.41{\pm}44.43\\ 157.65{\pm}35.69 \end{array}$	27.48 ± 1.30 91.66 ± 1.29 9.35 ± 1.09	296.47±106.82 410.59±97.29 215.98±90.62	
eSpark	Delay Trip time Wait time	$\begin{array}{c} 48.36 {\pm} 0.32 \\ 159.74 {\pm} 0.44 \\ 22.58 {\pm} 0.29 \end{array}$	$\begin{array}{c} 854.22{\pm}68.21\\ 484.87{\pm}58.21\\ 328.82{\pm}61.70 \end{array}$	$\begin{array}{c} 209.49 \pm 13.98 \\ 235.20 \pm 6.80 \\ 88.38 \pm 4.41 \end{array}$	25.39 ± 1.27 89.50 ± 1.36 6.94 ± 0.38	$\begin{array}{c} 243.92{\pm}15.81\\ 367.57{\pm}15.03\\ 180.09{\pm}13.84 \end{array}$	

Since we do not design specific prompts for different scenario tasks within the same environment, we calculate the average token consumption for all scenarios with each environment and display it in Table 15.

Table 15: Average token assumption for MABIM and SUMO.

Environment	Token assumption (M)
MABIM	3.0
SUMO	2.6

The training time and GPU memory usage for eSpark and the baselines across different scenarios are presented in Table 16 and Table 17.

1235	Table 16: GPU memory usage of eSpark and
1236	MARL baselines.

Table 17: Running time of eSpark and MARL baselines

1026	MAKL Dase	mes.			basennes.					
1230	Method	GPU	memory usa	ige (G)	Method	R	Running time (h)			
1238	memou	Standard	2 echelons	3 echelons	witchiou	Standard	2 echelons	3 echelons		
1230	eSpark	27.2	33.9	42.4	eSpark	18	25	30		
1239	IPPO	2.3	3.4	4.4	IPPO	6	8	12		
1240	QTRAN	4.2	6.5	8.7	QTRAN	9	15	21		
1241	-				-					

E.2 DETAILED RESULTS OF THE PRUNING METHODS

In this section, we provide the detailed results for multiple pruning baselines as discussed in Sec-tion 5.3, along with results of eSpark for comparison.

Table 18: Detailed performance of various pruning methods in MABIM.

Avg. profits (K)										
Method		100 SKUs scenarios				500 SKUs scenarios				
	Standard	2 echelons	3 echelons	Lower	Lowest	Standard	2 echelons	3 echelons	Lower	Lowest
Random pruning	733.0	1407.6	1426.6	511.6	262.0	2718.0	8667.4	9464.1	2535.5	2202.1
(S, s) pruning	394.4	832.3	933.4	441.1	258.3	3884.3	9248.3	10282.1	3517.9	2085.6
Upbound pruning	745.0	630.2	-2.2	557.7	294.7	3261.3	2473.6	1657.6	2833.0	2167.7
eSpark	823.7	1811.4	2598.7	579.5	405.0	4468.6	9437.3	12134.2	3775.7	2519.5

Table 19: Detailed performance of various pruning methods in SUMO, includes the mean and standard deviation.

Method	Metric	Grid 4×4	Arterial 4×4	Grid 5×5	Cologne8	Ingolstadt21
Ramdom pruning	Delay Trip time Wait time	$\substack{49.07 \pm 0.36 \\ 160.13 \pm 0.58 \\ 22.66 \pm 0.14 }$	$\begin{array}{c} 858.33 {\pm} 48.20 \\ 548.08 {\pm} 61.84 \\ 387.25 {\pm} 43.93 \end{array}$	$\begin{array}{c} 238.57 {\pm} 9.45 \\ 241.92 {\pm} 9.60 \\ 93.49 {\pm} 8.09 \end{array}$	$\begin{array}{c} 25.89{\pm}1.34\\ 89.75{\pm}1.26\\ 7.03{\pm}0.37\end{array}$	$\begin{array}{c} 353.38{\pm}24.39\\ 478.53{\pm}22.54\\ 281.97{\pm}21.90 \end{array}$
MaxPressure pruning	Delay Trip time Wait time	$\substack{48.78 \pm 0.37 \\ 160.72 \pm 0.17 \\ 23.03 \pm 0.56 }$	$\begin{array}{c} 890.04{\pm}121.50\\ 533.36{\pm}78.08\\ 391.96{\pm}78.34\end{array}$	$\begin{array}{c} 234.27{\pm}14.28\\ 253.68{\pm}17.68\\ 102.29{\pm}10.34\end{array}$	$\begin{array}{c} 26.26{\pm}0.36\\ 90.36{\pm}0.88\\ 7.33{\pm}0.12\end{array}$	$\begin{array}{c} 337.02{\pm}62.18\\ 448.11{\pm}65.32\\ 257.98{\pm}61.91 \end{array}$
eSpark	Delay Trip time Wait time	$\substack{48.36 \pm 0.32 \\ 159.74 \pm 0.44 \\ 22.58 \pm 0.29}$	854.22±68.21 484.87±58.21 328.82±61.70	$\begin{array}{c} 209.49 \pm 13.98 \\ 235.20 {\pm} 6.80 \\ 88.38 {\pm} 4.41 \end{array}$	$\begin{array}{c} 25.39{\pm}1.27\\ 89.50{\pm}1.36\\ 6.94{\pm}0.38\end{array}$	$\begin{array}{c} 243.92{\pm}15.81\\ 367.57{\pm}15.03\\ 180.09{\pm}13.84 \end{array}$

E.3 DETAILED RESULTS OF THE ABLATIONS

In this section, we provide the detailed results for ablations in Table 20 and Table 21, along with results of eSpark for comparison.

Mathad	Avg. profits (K)								
Methoa	Standard	2 echelons	3 echelons	Lower	Lowest				
eSpark w/o retention	719.0	1806.1	2388.6	547.7	294.1				
eSpark w/o LLM	754.7	1538.9	1109.9	536.7	198.5				
eSpark w/o checker	780.7	1741.6	2037.6	494	295.3				
eSpark w/o r factors	758.2	1688.1	2375.1	498.6	368.9				
eSpark	823.7	1811.4	2598.7	579.5	405.0				

Table 20: Detailed performance of ablations in MABIM.

F ESPARK'S PERFORMANCE WITH DIFFERENT RL FORMULATION

When the environment code is unavailable, we consider manually adding descriptions for transitions and rewards. We first introduce the ablation eSpark w/ lang, which replaces the environment codes with natural language descriptions. We conduct experiments on selected scenarios in MABIM with 100 SKUs, and the results are shown in Table 22 and Table 23. While the SUMO tasks remain largely unaffected, the MABIM tasks experience noticeable performance declines. This is because the transition logic and variable usage in SUMO are relatively simple, whereas MABIM involves more complex environment transitions and variable interactions. While environment codes provide precise transition dynamics and variable meanings, natural language often lacks this level of detail.

To further explore effective representation for improving performance in complex environments, we also adopt a strategy similar to Text2reward Xie et al. (2023), where experts create a simplified Python-style representation, referred to as eSpark w/ pyrep. As shown in Table 22, the kind of representation obtain performance on par with eSpark, which indicates that as long as accurate details and sufficient information are included, the language model can effectively understand the environment and eSpark can work well.

Table 21: Detailed performance of ablations in SUMO, includes the mean and standard deviation.

-						
Method	Metric	Grid 4×4	Arterial 4×4	Grid 5×5	Cologne8	Ingolstadt21
eSpark w/o retention	Delay Trip time Wait time	$\substack{48.75 \pm 0.49 \\ 160.42 \pm 0.54 \\ 22.97 \pm 0.39}$	$\begin{array}{c} 851.56{\pm}37.98\\ 487.05{\pm}65.66\\ 338.64{\pm}67.07 \end{array}$	211.07±22.01 248.96±15.04 97.77±10.25	$\begin{array}{c} 25.18{\pm}0.51\\ 89.23{\pm}0.49\\ 7.14{\pm}0.15\end{array}$	$\begin{array}{c} 246.05{\pm}14.88\\ 363.28{\pm}14.83\\ 175.79{\pm}12.47\end{array}$
eSpark w/o LLM	Delay Trip time Wait time	$\begin{array}{c} 48.67{\pm}0.48\\ 159.90{\pm}0.56\\ 22.99{\pm}0.43\end{array}$	$\begin{array}{c} 854.10{\pm}63.38\\ 491.49{\pm}67.52\\ 342.62{\pm}73.88\end{array}$	$\begin{array}{c} 212.25{\pm}16.13\\ 238.33{\pm}9.41\\ 90.89{\pm}6.60\end{array}$	$\begin{array}{c} 24.70{\pm}0.56\\ 88.57{\pm}0.60\\ 6.69{\pm}0.23\end{array}$	$\begin{array}{c} 257.14{\pm}40.50\\ 376.46{\pm}40.01\\ 188.14{\pm}37.03 \end{array}$
eSpark w/o checker	Delay Trip time Wait time	$\substack{48.29 \pm 0.53 \\ 159.62 \pm 0.48 \\ 22.59 \pm 0.48 }$	$\begin{array}{c} 872.6{\pm}94.89\\ 511.45{\pm}79.67\\ 347.43{\pm}68.39\end{array}$	$\begin{array}{c} 215.58{\pm}11.29\\ 246.56{\pm}9.38\\ 95.95{\pm}6.33\end{array}$	$\begin{array}{c} 25.22{\pm}0.69\\ 89.20{\pm}0.68\\ 6.92{\pm}0.28\end{array}$	$\begin{array}{c} 258.39{\pm}14.44\\ 383.33{\pm}17.07\\ 192.96{\pm}15.37 \end{array}$
eSpark w/o r factors	Delay Trip time Wait time	$\substack{48.54 \pm 0.41 \\ 159.89 \pm 0.63 \\ 22.85 \pm 0.36}$	$\begin{array}{c} 849.82{\pm}44.81\\ 479.29{\pm}57.12\\ 331.69{\pm}74.99 \end{array}$	$\begin{array}{c} 216.27{\pm}20.01\\ 247.28{\pm}8.30\\ 90.14{\pm}6.95\end{array}$	$\begin{array}{c} 25.29{\pm}0.61\\ 89.25{\pm}0.62\\ 6.88{\pm}0.19\end{array}$	$\begin{array}{c} 240.47{\pm}19.69\\ 373.32{\pm}21.19\\ 180.90{\pm}13.42 \end{array}$
eSpark	Delay Trip time Wait time	$\substack{48.36 \pm 0.32 \\ 159.74 \pm 0.44 \\ 22.58 \pm 0.29}$	854.22±68.21 484.87±58.21 328.82±61.70	$\begin{array}{c} 209.49 \pm \! 13.98 \\ 235.20 {\pm} 6.80 \\ 88.38 {\pm} 4.41 \end{array}$	$\begin{array}{c} 25.39{\pm}1.27\\ 89.50{\pm}1.36\\ 6.94{\pm}0.38\end{array}$	$\begin{array}{c} 243.92{\pm}15.81\\ 367.57{\pm}15.03\\ 180.09{\pm}13.84 \end{array}$

Table 22: Performance of eSpark and its ablations in the 100 SKUs setting of the MABIM environ-ment.

Mathad	Profits (K)			
Methou	Standard	3 echelons	Lowest	
eSpark	823.7	2598.7	405.0	
eSpark w/ lang	777.5	752.7	347.0	
eSpark w/ pyrep	817.4	2522.2	409.7	

G POLICY PERFORMANCE ANALYSIS

To gain a deeper understanding of the policy difference between eSpark and IPPO, we select the capacity limit and multiple echelons challenges within the 100 SKUs scenario as representative cases, presenting in Figure 4 the daily profit of eSpark and IPPO on the test dataset challenged with capacity limit and multiple echelons. In the capacity limit challenges, a high daily overflow ratio and low fulfillment ratio suggest that IPPO falls short in mastering the adjustment of restocking quantities for individual agents when capacity is limited, leading to overstocking and substantial overflow. Concurrently, this prevents SKUs required by consumers from being accommodated, culminating in an exceedingly low fulfillment ratio. In multiple echelon challenges, the fulfillment ratio at each echelon gradually decreases over time, indicating that IPPO struggles to comprehend and learn the intricate interplay required for cooperation among various echelons, thereby inadequately fulfilling the demands of each echelon. Such shortcomings not only diminish potential profits but also subject the system to considerable backlog expenses. However, through action space pruning, evolutionary search, and reflection, eSpark manages to reduce the search within the vast space, selecting the most effective exploration functions and continuously improving. This approach significantly reduces overflow in the capacity limit scenario and successfully learns suitable cooperation methods for multiple echelons.

Table 23: Performance of eSpark and its ablations in SUMO, includes the mean and standard deviation.

Method	Motric	Time usage (s)	
Method	Mente	Arterial 4×4	Ingolstadt21
	Delay	851.56 ± 37.98	$246.05{\pm}14.88$
eSpark	Trip time	$484.84{\pm}58.21$	$367.57 {\pm} 15.03$
	Wait time	$328.82{\pm}61.70$	180.09 ± 13.84
	Delay	830.05±75.95	243.32±19.43
eSpark w/ lang	Trip time	$495.18{\pm}18.60$	$365.65{\pm}21.48$
	Wait time	$351.32{\pm}27.07$	178.31 ± 19.31



Figure 4: The performance comparison between eSpark and IPPO in 100 SKUs scenarios. In capacity-limited scenarios, eSpark strives to meet the demands while minimizing overflow costs, boasting a lower overflow ratio and a higher fulfillment ratio. In the multiple-echelon challenge, eSpark achieves nuanced collaboration across different echelons, ensuring high fulfillment ratios.

H ESPARK'S RESPONSE WITH AND WITHOUT REWARD FACTORS

1367 1368

1361

1362

1363

1364 1365 1366

We select the policy feedback from one iteration in the Lower scenario of the MABIM environment with 100 SKUs as an example. The complete policy feedback is shown in Figure 5. In the ablation experiment, all reward factors (highlighted in red) are removed. Figure 6 and Figure 7 illustrate the differences in GPT-4's reflections with and without the reward components. After removing the reward feedback, GPT-4 struggles to identify the reasons for poor performance and to thoroughly analyze ways to improve it.

1375	
1376	We trained a RL policy using the provided exploration function code and tracked the values of the individual
1377	training stage, mean in the late training stage, mean in all the training stage, minimum values for reward and its components after every (enoch fred) enochs. Each element is a one-dimensional array of length n warehouse
1378	representing the value of that component on different warehouse:
1379	metric_name: reward, Max: [435592.087], Min: [-1331093.355], Mean in the early training stage: [-212101.5832], Mean in the late training stage: [338683.3766], Mean in all the training stage: [71512.5884375]
1380	stage: [564756.2], Mean in all the training stage: [527295.125]
1381	metric_name: excess_cost, Max: [1648683.5], Min: [1/213.5], Mean in the early training stage: [596130.9], Mean in the late training stage: [123902.4], Mean in all the training stage: [349431.25]
1382	metric_name: order_cost, Max: [60580.], Min: [56490.], Mean in the early training stage: [58864.], Mean in the late training stage: [58758.], Mean in all the training stage: [58650.625]
1383	<pre>metric_name: holding_cost, Max: [456.135], Min: [272.157], Mean in the early training stage: [390.6432], Mean in the late training stage: [399.6234], Mean in all the training stage: [394.1803125]</pre>
1384	<pre>metric_name: backlog_cost, Max: [56796.1], Min: [40408.7], Mean in the early training stage: [50364.04], Mean in the late training stage: [43012.8], Mean in all the training stage: [47306.48125]</pre>
1385	Please carefully analyze the policy feedback and provide a new, improved exploration function that can better solve
1386	the task. Some helpful tips for analyzing the policy feedback: (1) You can start with "let's think step by step", and then look at each reward individually and think about how
1387	can you improve it (2) If the total reward maintains in the same level or even reduce, then you must rewrite the entire exploration
1388	(3) If the values for a certain reward component are near identical throughout, then this means PL is not able
1389	to optimize this component as it is written. You may consider: (a) Changing the temperature scale or value of the related mask component so that more action can be
1390	explored
1391	 (b) Re-writing the related mask component (c) Discarding the mask component or add a new mask component



1393 1394

1392

1395

1396	
1397	From the policy feedback, we can identify a few key issues:
1398	1. The total reward is negative in the early training stage and only becomes positive in the later training stages. This suggests that the agent is not learning effective strategies early on, potentially due to a lack of exploration.
1399	The excess cost is quite high, especially in the early training stages. This suggests that the agent is overstocking the warehouse.
1400	The holding cost remains high throughout the training, indicating that the agent is holding onto inventory for too long.
1401	
1402	Given these issues, we need to revise the exploration function to encourage the agent to explore actions that would result in less overstocking. Here's an improved exploration function:
1400	

Figure 6: eSpark's reflection with reward factors.

From	the policy foodback, we can identify a few key issues:
1 т	the policy recuback, we can identify a rew key issues.
sign	he policy feedback indicates that the reward function has a wide range, with the minimum values being ificantly negative and the maximum being positive. The mean reward in the early stage is negative, while it mean reward in the late stage. This suggests that the policy improves over time but has considerable variability.
in o	mes positive in the tate stage. This suggests that the policy improves over time but has considerable variability utcomes. He high variance in reward components suggests that the evoloration may be too aggressive leading to
over 3. T opti	he high fortained in Found components suggests that the experiment may be used signestre, treating to stocking (incurring high holding costs and overflow costs) or understocking (incurring backlog costs). he exploration function may be too restrictive or lenient, preventing the RL agent from effectively learning the mal policy.
We w expl rang	ill modify the mask to include a temperature parameter, which will initially be set high to allow for broader oration and gradually decrease it as the training progresses. This change will help the agent explore a wider e of actions early on and focus more on promising actions as it learns.
	Figure 7: eSpark's reflection without reward factors.
F	ULL PROMPTS
/e re cplo inct /e pi	eference the prompt design outlined in the Eureka Ma et al. (2024) and adapt it specifically fration function generation. Our prompt provides general guidance on the design of explorations, specific code formatting suggestions, feedback, and recommendations for improveme resent our prompts for MABIM below.
You a funci	re an expert in both inventory management learning and reinforcement learning. You will get some exploration in
Your any i help and t (1) / (2) /	goal is to evaluate whether the given exploration function matches the task description and whether it contair llogical errors in the code content, and evaluate whether it's possible to avoid some unreasonable actions, the exploration of reinforcement learning. You need to pay special attention to the meaning of each state item the logic of the task, making sure to detect all incorrect use of variables in code all the parts that don't follow logic.
The e {tasl	<pre>:xploration function signature can be: (_exploration_signature_string}</pre>
Your impor Under	advice can be text or snippets of code, but it should not be the full exploration function code. Most 'tantly, remember that your response must begin with either "Code passes check." or "Code fails to pass check." ' no circumstance can you begin your answer with other content.
	Figure 8: System prompt for LLM_c .
	Figure 8: System prompt for LLM _c .
	Figure 8: System prompt for LLM _c . You are an expert in both inventory management learning and reinforcement learning. You are trying to write some exploration functions, which helps mask some actions that are logically unlikely to be selected, to help exploration in reinforcement learning tasks as effective as possible. Your goal is to write a exploration function for the agent that will mask the actions that's almost impossible to be chosen in the task described in text
	Figure 8: System prompt for LLM _c . You are an expert in both inventory management learning and reinforcement learning. You are trying to write some exploration functions, which helps mask some actions that are logically unlikely to be selected, to help exploration in reinforcement learning tasks as effective as possible. Your goal is to write a exploration function for the agent that will mask the actions that's almost impossible to be chosen in the task described in text. The exploration function signature can be: {task exploration function signature string}
	Figure 8: System prompt for LLM _c . You are an expert in both inventory management learning and reinforcement learning. You are trying to write some exploration functions, which helps mask some actions that are logically unlikely to be selected, to help exploration in reinforcement learning tasks as effective as possible. Your goal is to write a exploration function for the agent that will mask the actions that's almost impossible to be chosen in the task described in text. The exploration function signature can be: {task_exploration_signature_string} Your exploration function should only use the variables from the argument list. Please just give only the exploration function and don't put it in a class. Please make sure that the code is compatible with numpy (e.g., use numpy array instead of torch tensor).
	Figure 8: System prompt for LLM _c . You are an expert in both inventory management learning and reinforcement learning. You are trying to write some exploration functions, which helps mask some actions that are logically unlikely to be selected, to help exploration in reinforcement learning tasks as effective as possible. Your goal is to write a exploration function for the agent that will mask the actions that's almost impossible to be chosen in the task described in text. The exploration function signature can be: {task_exploration_signature_string} Your exploration function should only use the variables from the argument list. Please just give only the exploration function and don't put it in a class. Please make sure that the code is compatible with numpy (e.g., use numpy array instead of torch tensor).
	Figure 8: System prompt for LLM _c . You are an expert in both inventory management learning and reinforcement learning. You are trying to write some exploration functions, which helps mask some actions that are logically unlikely to be selected, to help exploration in reinforcement learning task as effective as possible. Your goal is to write a exploration function for the agent that will mask the actions that's almost impossible to be chosen in the task described in text. The exploration function signature can be: {task_exploration_signature_string} Your exploration function should only use the variables from the argument list. Please give only the exploration function and on't put if in a class. Please make sure that the code is compatible with numpy (e.g., use numpy array instead of torch tensor). Figure 9: System prompt for LLM _g .
Vrite {tack	<pre>Figure 8: System prompt for LLM_c. You are an expert in both inventory management learning and reinforcement learning. You are trying to write some exploration functions, which helps mask some actions that are logically unlikely to be selected, to help exploration in reinforcement learning tasks as effective as possible. Your goal is to write a exploration function for the agent that will mask the actions that's almost impossible to be chosen in the task described in text. The exploration function signature can be: {task_exploration_signature_string} Your exploration function should only use the variables from the argument list. Please just give only the exploration function and don't put it in a class. Please make sure that the code is compatible with numpy (e.g., use numpy array instead of torch tensor). Figure 9: System prompt for LLM_g. a exploration function for the following task : introduction] </pre>
√rite {task	<pre>Figure 8: System prompt for LLM_c. You are an expert in both inventory management learning and reinforcement learning. You are trying to write some exploration functions, which helps mask some actions that are logically unlikely to be selected, to help exploration in reinforcement learning tasks as effective as possible. Your goal is to write a exploration function for the agent that will mask the actions that's almost impossible to be chosen in the task described in text. The exploration function signature can be: task_exploration function should only use the variables from the argument list. Please just give only the exploration function and don't put it in a class. Please make sure that the code is compatible with numpy (e.g., use numpy array instead of torch tensor). Figure 9: System prompt for LLM_g. a exploration function for the following task :introduction} </pre>
√rite {task [he d {tran [he a	<pre>Figure 8: System prompt for LLM_c. You are an expert in both inventory management learning and reinforcement learning. You are trying to write some exploration functions, which helps mask some actions that are logically unlikely to be selected, to help exploration in reinforcement learning tasks as effective as possible. Your goal is to write a exploration function for the agent that will mask the actions that's almost impossible to be chosen in the task described in text. The exploration function signature can be: {task_exploration function sould only use the variables from the argument list. Please just give only the exploration function and don't put it in a class. Please make sure that the code is compatible with numpy (e.g., use numpy array instead of torch tensor). figure 9: System prompt for LLM_g. setup of environment and transition are : sition_definition; gent and state definition is : edefinition? </pre>
Vrite {task The d {tran Che a (stat	<pre>Figure 8: System prompt for LLM_c. Vou are an expert in both inventory management learning and reinforcement learning. You are trying to write some exploration functions, which helps mask some actions that are logically unlikely to be selected, to help exploration in reinforcement learning tasks as effective as possible. Vour goal is to write a exploration function for the agent that will mask the actions that's almost impossible to be the exploration function signature can be: {task_exploration_signature_string} Vour exploration function should only use the variables from the argument list. Please just give only the exploration function and don't put it in a class. Please make sure that the code is compatible with numpy (e.g., use numpy array instead of torch tensor). f gure 9: System prompt for LLM_g. sequeration function for the following task :</pre>

Figure 10: Initial prompt for LLM_g .

1/59	
1458	Please carefully check the exploration function for the following task : {task_introduction}
1460	The definition of environment and transition are :
1461	The agent and state definition is t
1462	{state_definition}
1463	The reward definition is : {reward definition}
1465	The requirements of code :
1466	{code_output_tip}
1467	
1468	Figure 11: Initial prompt for LLM_c .
1469	
1470	
1471	
1472	
1473	Here is the latest exploration function and it's description: {gpt_response}
1474	Please carefully review this code, check whether it matches the task description and whether it contains any illegical errors in the code content, and evaluate whether it's possible to improve the evaluation and the
1470	performance.
1477	
1478	Figure 12: LLM_g 's feedback to LLM_c .
1479	
1480	
1481	
1482	
1483	We discuss this code with experts, and the code is not approved by experts, and the comments of experts on this code are as follows:
1484	{checker_feedback}
1485	Please refer to expert advice and combine your own knowledge, fix the problems and provide a new, improved
1486	exploration function!
1487	
1489	Figure 13: LLM_c 's feedback to LLM_g .
1490	
1491	
1492	
1493	<pre>def compute_mask(agent_states: AgentStates, supply_chain: SupplyChain, action_space: list):</pre>
1494	# Here are some code you can refer to when you generate your exploration function.
1495	return total_mask, {}
1496	
1497	Figure 14: Signature of exploration function.
1499	0 · · · · 0 · · · · · · · · · · · · · ·
1500	
1501	
1502	
1503	Please carefully analyze the policy feedback and provide a new, improved exploration function that can better solve the task. Some helpful tips for analyzing the policy feedback:
1504	(1) You can start with "let's think step by step", and then look at each reward individually and think about how can you improve it.
1505	 (2) If the total reward maintains in the same level or even reduce, then you must rewrite the entire exploration
1506	(3) If the values for a certain reward component are near identical throughout, then this means RL is not able
1507	 (a) Changing the temperature scale or value of the related mask component so that more action can be
1508	(b) Re-writing the related mask component
1509	(c) viscarding the mask component or add a new mask component
1510	
1311	

Figure 15: Output and improvement tips for LLM_g .

The output of the exploration function should be a total mask (1 denote action is not masked and 0 otherwise). The code output should be formatted as a python code string: "```python ... ```". Some helpful tips for writing the exploration function code: (1) Your total mask and its component should be a 3-D numpy array in [warehouse_name, sku_type, action_mask]. The masks of the actions that are available are set to 1, otherwise 0. No other elements are allowed to appear in total mask. (2) If you choose to transform a component mask, then you must also introduce a temperature parameter inside the transformation function; this parameter must be a named variable in the mask function and it must not be an input variable. Each transformed mask component should have its own temperature variable (3) Make sure the type of each input variable is correctly specified; For example, a float input variable should not be specified as torch.Tensor (4) Most importantly, the exploration code's can only use variables defined in its arguments. Under no circumstance can you introduce new input variables. You only need to give the definition of exploration function and no other function or class should be defined. Figure 16: Output format for LLM_q . Your output should contain two parts: (1) Your response must begin with either "Code passes check." or "Code fails to pass check.". "Code passes check." means you believe that there are no logical errors in the code and the variables are taken in accordance with the description of the task. "Code fails to pass check." indicates the provided exploration function contains logical errors, or you think the code is obviously flawed and you can point out how to facilitate more effective exploration. (2) If you begin with "Code fails to pass check.", you have to explain why the code fails to pass the check and give your advice on fixing the problems; If you begin with "Code passes check.", you also have to state why each part is logical and reasonable. Some common logical errors include: (1) Misunderstanding the meaning of the state items, or including syntax errors when using variables (2) Illogically handling the state items (3) Using multiple unrelated state items to calculate mask component(4) The way of combining mask components into total mask is unreasonable Figure 17: Output format for LLM_c. We trained a RL policy using the provided exploration function code and tracked the values of the individual components of the reward function as well as global policy metrics. We also compute the maximum, mean in the early training stage, mean in the late training stage, mean in all the training stage, minimum values for reward and its components after every {epoch_freq} epochs. Each element is a one-dimensional array of length n_warehouse, representing the value of that component on different warehouses: <Reward Feedback Here> Figure 18: Reward feedback and action feedback.

1566 J ESPARK'S EXPLORATION FUNCTION EDITING

In this section, we demonstrate the reward editing capabilities of eSpark. eSpark is capable of reflecting upon feedback to optimize the exploration for subsequent iterations.

We trained a RL policy using the provided exploration function code and tracked the values of the individual components of the reward function as well as global policy metrics. We also compute the maximum, mean in the early training stage, mean in the late training stage, mean in all the training stage, minimum values for reward and its components after every (epoch_freq) epochs. Each element is a one-dimensional array of length n_warehouse, representing the value of that component on different warehouses: metric_name: reward, Max: [435592.087], Min: [-1331893.335], Mean in the early training stage: [-212101.5832], Mean in the late training stage: [338683.3766], Mean in all the training stage: [71512.5884375] metric_name: profit, Max: [567080.1, Min: [43087.], Mean in the early training stage: [493648.], Mean in the late training stage: [564756.2], Mean in all the training stage: [527295.125] metric_name: excess_cost, Max: [1648683.5], Min: [17213.5], Mean in the early training stage: [596130.9], Mean in the late training stage: [58654.2], Mean in all the training stage: [54758.], Mean in the late training stage: [58650.6], Min: [50490.1], Mean in the early training stage: [58664.], Mean in the late training stage: [58654.], Mean in all the training stage: [58650.625] metric_name: order_cost, Max: [60680.], Min: [272.157], Mean in the early training stage: [390.6432], Mean in the late training stage: [58758.], Mean in all the training stage: [394.1803125] metric_name: holding_cost, Max: [456796.1], Min: [272.157], Mean in the early training stage: [390.6432], Mean in the late training stage: [4306.42], Mean in all the training stage: [394.1803125] metric_name: backlog_cost, Max: [65796.1], Min: [40408.7], Mean in the early training stage: [50364.04], Mean in the late training stage: [4301.28], Mean in all the training stage: [394.6432], Mean in the late training stage: [4306.48125] 1571 1572 1573 1574 1575 1576 1578 1579 late training stage: [43012.8], Mean in all the training stage: [47306.48125] 1580 Please carefully analyze the policy feedback and provide a new, improved exploration function that can better solve the task. Some helpful tips for analyzing the policy feedback: (1) You can start with "let's think step by step", and then look at each reward individually and think about how 1581 can you improve it
 (2) If the total reward maintains in the same level or even reduce, then you must rewrite the entire exploration 1582 function (3) If the values for a certain reward component are near identical throughout, then this means RL is not able to optimize this component as it is written. You may consider: 1584 (a) Changing the temperature scale or value of the related mask component so that more action can be 1585 explored (b) Re-writing the related mask component 1586 (c) Discarding the mask component or add a new mask component 1587

Figure 19: Policy feedback from the last iteration.

From the policy feedback, we can identify a few key issues:

 The total reward is negative in the early training stage and only becomes positive in the later training stages. This suggests that the agent is not learning effective strategies early on, potentially due to a lack of exploration.
 The excess cost is quite high, especially in the early training stages. This suggests that the agent is overstocking the warehouse.
 The holding cost remains high throughout the training, indicating that the agent is holding onto inventory for too long.
 Given these issues, we need to revise the exploration function to encourage the agent to explore actions that would

result in less overstocking. Here's an improved exploration function:

Figure 20: eSpark reflects upon the policy feedback and proposes the modification.

mpute_mask(agent_states: AgentStates, supply_c chain: SupplyChain, action_space: list) -> np.ndarray warehouse_list = supply_chain.get_warehouse_list() num sku = len(agent states) total_mask = np.ones((len(warehouse_list), num_sku, len(action_space))) capacity = supply chain [warehouse name, in_stock = agent_states[warehouse_name, "in_stock", "today", "all_skus"] history_demand_mean = np.average(agent_states[warehouse_name, "demand", "lookback", "all_skus"], axis=1) for a, action in enumerate(action_space): replenishment = history_demand_mean * action if np.sum(sku_volume * (in_stock + replenishment)) > capacity: total_mask[i, :, a] = 0 # Add an additional condition to mask actions that bring in too much replenishment if action > 2: total mask[i. :. a] = 0 if agent_states.current_step < agent_states.durations * 0.3 and action < 0.5: total_mask[i, :, a] = 0



1588 1589 1590

1591

1592

1593

1594

1595

1596

1597 1598 1599

1602

1605

1606

1609

1610

1611

1612

1613

1614

1615

1616 1617

Figure 21: Comparison of exploration functions before and after editing.