# YOLO-MARL: You Only LLM Once for Multi-Agent Reinforcement Learning

Yuan Zhuang[1†], Yi Shen[2†], Zhili Zhang[1], Yuxiao Chen[3] and Fei Miao[1*]

*Abstract*— **Advancements in deep multi-agent reinforcement learning (MARL) have positioned it as a promising approach for decision-making in cooperative games. However, it still remains challenging for MARL agents to learn cooperative strategies for some game environments. Recently, large language models (LLMs) have demonstrated emergent reasoning capabilities, making them promising candidates for enhancing coordination among the agents. However, due to the model size of LLMs, it can be expensive to frequently infer LLMs for actions that agents can take. In this work, we propose You Only LLM Once for MARL (YOLO-MARL), a novel framework that leverages the high-level task planning capabilities of LLMs to improve the policy learning process of multi-agents in cooperative games. Notably, for each game environment, YOLO-MARL only requires one time interaction with LLMs in the proposed strategy generation, state interpretation and planning function generation modules, before the MARL policy training process. This avoids the ongoing costs and computational time associated with frequent LLMs API calls during training. Moreover, trained decentralized policies based on normal-sized neural networks operate independently of the LLM. We evaluate our method across two different environments and demonstrate that YOLO-MARL outperforms traditional MARL algorithms. The Github repository of our code can be found at https://github.com/paulzyzy/YOLO-MARL.**

## I. INTRODUCTION

MARL has become a vital framework for solving decision-making problems in complex multi-agent systems. With the rising applications of multi-agent systems, it is increasingly crucial for individual agents to cooperate or compete in dynamic environments without relying on centralized control [1]. In cooperative Markov games, agents work together to maximize joint rewards. However, existing MARL approaches often struggle with learning effective distributed policies, particularly in tasks characterized by sparse rewards, dynamic environments, and large action spaces.

Concurrently, LLMs have demonstrated remarkable capabilities as high-level semantic planners by leveraging in-context learning and extensive prior knowledge [2]. Recent studies have showcased LLMs deployed as embodied agents [3], [4], as well as their use in guiding reinforcement learning (RL) as ELLM, which leverages LLMs to suggest goals [5], and work focusing on aligning LLM-provided actions with RL policies [6]. Despite these promising developments, integrating LLMs into multi-agent settings remains largely unexplored. Moreover, repeated interactions with LLMs during long-episode or complex tasks can be both time-consuming and computationally prohibitive, particularly when training over tens of millions of steps.

Addressing these challenges, we propose YOLO-MARL, a novel framework combining LLM's high-level planning capabilities with MARL policy training, as shown in Fig. 1. YOLO-MARL uniquely requires just one LLM interaction per environment. Once the strategy generation, state interpretation, and planning function generation modules produce the necessary guidance, the MARL training process proceeds without further LLM involvement. This design substantially reduces the communication overhead and computational cost typically associated with LLM inferences. Moreover, YOLO-MARL demonstrates its strong generalization capability and simplicity for application, since it only requires a basic understanding of the new environments from users. We validate our framework on a sparse reward multi-agent environment: Level-Based Foraging (LBF) environment [1] as well as the Multi-Agent Particle (MPE) environment [7], and our results show that YOLO-MARL outperforms several MARL baselines, such as [7]–[9]. To the best of our knowledge, this work represents one of the first trials to fuse the high-level reasoning and planning abilities of LLMs with MARL, pointing a new direction for scalable and efficient multi-agent coordination [10].

In summary, our proposed method YOLO-MARL has the following advantages:

- This framework synergizes the planning capabilities of LLMs with MARL to enhance the policy learning in challenging cooperative game environments. In particular, our approach exploits the LLM's wide-ranging reasoning ability to generate high-level planning functions to facilitate agents in coordination.
- YOLO-MARL requires minimal LLMs involvement, which significantly reduces computational overhead and mitigates communication connection instability concerns when invoking LLMs during the training process.
- Our approach leverages zero-shot prompting and can be easily adapted to various game environments, with only basic prior knowledge required from users.

An overview of YOLO-MARL is presented in Fig. 1, and all related prompts, environments, and generated planning functions are available in our GitHub repository.

## II. RELATED WORK

### A. Multi-Agent Reinforcement Learning

MARL has attracted significant attention for its potential to address complex, decentralized problems. A popular

---
[†] These authors contributed equally to this work.
[1]University of Connecticut
[2]University of Pennsylvania
[3]NVIDIA
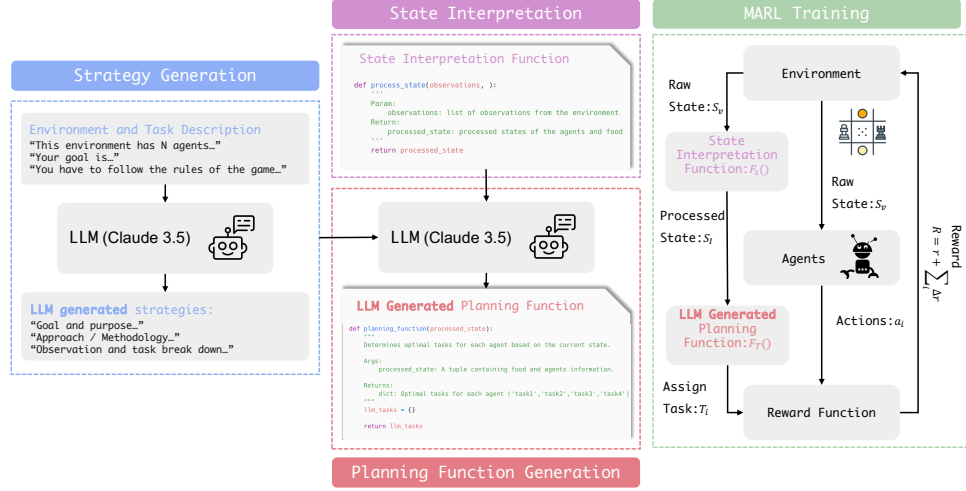[*]Corresponding author: Fei Miao (fei.miao@uconn.edu)

State Interpretation

Strategy Generation

MARL Training

State Interpretation Function

```
def process_state(observations, ):
    '''
    Param:
        observations: list of observations from the environment
    Return:
        processed_state: processed states of the agents and food
    '''
    return processed_state
```

Environment and Task Description
"This environment has N agents…"
"Your goal is…"
"You have to follow the rules of the game…"

LLM (Claude 3.5)

LLM (Claude 3.5)

LLM generated strategies:
"Goal and purpose…"
"Approach / Methodology…"
"Observation and task break down…"

LLM Generated Planning Function

```
def planning_function(processed_state):
    '''
    Determines optimal tasks for each agent based on the current state.

    Args:
        processed_state: A tuple containing food and agents information.

    Returns:
        dict: Optimal tasks for each agent {'task1','task2','task3','task4'}
    '''
    llm_tasks = {}
    return llm_tasks
```

Planning Function Generation

Environment

Raw State: $S_v$

State Interpretation Function: $F_s()$

Raw State: $S_v$

Processed State: $S_I$

Agents

LLM Generated Planning Function: $F_T()$

Actions: $a_i$

Reward $R = r + \sum_i \Delta r$

Assign Task: $T_i$

Reward Function

Fig. 1: Depiction of our framework YOLO-MARL. (a). Strategy Generation: We pass basic environment and task description into the LLM to get generated strategies for this specific environment. (b). State Interpretation: We process the global states so that the format of global states will be more structured and organized for better comprehension by the LLM. (c). Planning Function Generation: We chain together the environment and task description, LLM generated strategies and state interpretation function. These prompts are then fed into the LLM to generate a Python planning function for this environment. (d). MARL Training: The state interpretation function and the generated planning function are integrated into the MARL training process. The LLM is no longer required for further interaction after the Planning Function Generation. The more detailed explanation of MARL Training part can be found in Algorithm 1

paradigm is centralized training with decentralized execution. Methods like QMIX [9] and MADDPG [7] employ centralized critics during training to coordinate agents, while allowing independent execution at test time. In cooperative environments, COMA [11] and VDN [12] enable agents to share rewards and maximize joint returns. Despite these advances, many existing MARL algorithms struggle in sparse-reward settings and have difficulty learning fully cooperative policies. Moreover, only limited work has explored the integration of LLMs with MARL [10], leaving open questions about leveraging LLMs for multi-agent decision-making.

### B. Large Language Models for RL and Decision-Making

Recent studies have incorporated LLMs into the RL training process to enhance performance. For instance, [5] improve agent exploration by aligning LLM-suggested goals with observed behaviors, while [13] generate actions conditioned on language-based prompts during online RL. Similarly, [6] uses LLMs to provide scalar rewards that guide training. However, many of these approaches focus on single-agent settings or require extensive interactions with LLMs during training. And [14] highlights the limitations of LLMs in handling complex low-level tasks. Other works, such as [15], studies a multi-task RL problem. SayCan [2] grounds LLMs via value functions of pretrained skills to execute abstract commands on robots. [16] finds that code-writing LLMs can be re-purposed to write robot policy code. Additionally, studies like [17] and [18] exploit LLMs' prior knowledge and code-generation capabilities to produce

reward functions. In contrast, our approach leverages LLMs to generate planning functions, thereby enhancing MARL without continuous LLM involvement.

### C. Large Language Models for Multi-Agent Systems

LLM-based multi-agent systems have been employed in diverse applications requiring varied agent roles and collaborative decision-making [10], [19]. Camel [20] and MetaGPT [21] deploy multiple LLM agents for tasks like brainstorming and software development. SMART-LLM [22] decomposes multi-robot task planning into subgoals, and Co-NavGPT [23] uses LLMs as global planners for cooperative navigation. Other research has focused on applying LLMs in strategic gaming environments [24]–[28]. Unlike these approaches, which use LLMs directly as agents or decision-makers, our method harnesses the planning capabilities of LLMs to train compact, efficient MARL policies.

### III. PROBLEM FORMULATION

**Markov game** is defined as a multi-agent decision-making problem when the interaction between multiple agents affect the state dynamics of the entire system and the reward of each agent under certain conditions [29]. In this work, we consider a Markov game, or a stochastic game [30] defined as a tuple $G := (\mathcal{N}, S, A, \{r^i\}_{i \in \mathcal{N}}, p, \gamma)$, where $\mathcal{N}$ is a set of $N$ agents, $S = S^1 \times \cdots \times S^N$ is the joint state space, $A = A^1 \times \cdots \times A^N$ is the joint action space, with $(S^i, A^i)$ as the state space and action space of agent $i$, respectively, $\gamma \in [0, 1)$ is the discounting factor [29], [30]. The state transition $p : S \times A \rightarrow$

$\Delta(S)$ is controlled by the current state and joint action, where $\Delta(S)$ represents the set of all probability distributions over the joint state space $S$. Each agent has a reward function, $r^i : S \times A \to \mathbb{R}$. At time $t$, agent $i$ chooses its action $a_t^i$ according to a policy $\pi^i : S \to \Delta(A^i)$.

For each agent $i$, it attempts to maximize its expected sum of discounted rewards, i.e. its objective function $J^i(s, \pi) = \mathbb{E}\left[\sum_{t=1}^{\infty} \gamma^{t-1} r_t^i(s_t, a_t) | s_1 = s, a_t \sim \pi(\cdot | s_t)\right]$. In the literature, MARL algorithms [7]–[9] have been designed to train neural network-based policies $\pi_i(\theta_i)$. For a cooperative game, one shared reward function for all agents is widely used during the training process, which is also considered in this work.

## IV. METHODOLOGY

In this section, we introduce YOLO-MARL, a framework that leverages the high-level planning capabilities of LLMs to enhance MARL. YOLO-MARL integrates four key components: Strategy Generation, State Interpretation, Planning Function Generation, and MARL training process with the LLM generated Planning Function incorporated throughout.

---

**Algorithm 1** YOLO-MARL Training Process

---

**Require:** Large Language Model *LLM*, State Interpretation function $F_S$, MARL actor $\mathscr{A}$, MARL algorithm $MARL_{alg}$, Initial Prompts $P_{init}$
1: **Hyperparameters:** reward signal $r'$, penalty signal $p'$
2: $P_{Strategy} \sim LLM(P_{init})$ // Strategy Generation
3: $P = P_{init} + P_{Strategy} + F_S$ // Chaining all the prompt for Planning Function Generation
4: $\mathscr{F}_{\mathscr{T}} \sim LLM(P)$ // Planning Function Generation: Sample functions code from the LLM
5: **MARL training with generated planning function**
6: **for** each training step **do**
7:     $S_I \leftarrow F_S(S_v)$ // State Interpretation: Get processed global observation $S_I$ from $F_S$
8:     $\mathscr{T}_1, \mathscr{T}_2, \ldots \leftarrow \mathscr{F}_{\mathscr{T}}(S_I)$ // Assign tasks $\mathscr{T}$ to each agent
9:     $a_1, a_2, \ldots \leftarrow \mathscr{A}(S_v)$ // Output actions from the actor
10:     **for** each agent $i$ **do**
11:         **if** $a_i \in \mathscr{T}_i$ **then**
12:             $\Delta r_i \leftarrow r'$
13:         **else**
14:             $\Delta r_i \leftarrow p'$
15:         **end if**
16:     **end for**
17:     $R \leftarrow r + \sum_i \Delta r_i$ // Compute final reward for critic: More details are in equation 1, 2
18:     $\pi(\theta) = MARL_{alg}(R)$ // Use $R$ as the final reward for MARL training
19: **end for**
20: **return** Trained MARL policy

---

### A. Strategy Generation

To ensure applicability across diverse environments—even for users with limited domain knowledge—we incorporate a *Strategy Generation Module*, as shown in the blue box of Figure 1(a). In this module, the LLM receives basic environment details (e.g., task descriptions, rules, and constraints) along with a general guideline, and then autonomously outputs a detailed strategy in a prescribed format without requiring extensive human input or expertise.

The Strategy Generation is crucial for several reasons:

- Reducing User Burden: It alleviates the need for users to comprehensively understand new environments, saving time and effort.
- Enhancing Generalization: It enables the framework to adapt to different environments with minimal prompt modifications.
- Facilitating Planning Function Generation: The strategies serve as vital components in the prompts used for the Planning Function Generation Module. The results of using YOLO-MARL but without Strategy Generation Module are shown in ablation study VI-A .

The LLM-generated strategies are incorporated into the prompt alongside other necessary information to facilitate the subsequent planning function generation.

### B. State Interpretation

In many simulation environments, observations are provided as vectors, with each component encoded in a non-semantic format. Although such representations are effective for training deep RL models, they pose challenges for LLMs that require context to interpret each component correctly.

We propose the *State Interpretation Module* to assist the LLM in interpreting the environment state. By providing a semantically meaningful representation of the state, the LLM can successfully generate executable planning functions for training. Formally, given the current environment state in vector form $S_v$, we define an interpretation function $F_S$ such that $F_S(S_v) \to S_I$, where $S_I$ provides more explicit and meaningful information about each state component.

Recent works like [17] and [18] have demonstrated the success of enhancing LLMs performance by providing relevant environment code. In the same manner, we include the interpretation function $F_S$ in the prompting pipeline, formatted as Python environment code as shown in the purple box in Figure 1(b). The State Interpretation Module significantly reduces the risk of the LLM generating erroneous functions with outputs incompatible with the training procedures. An ablation study on the effectiveness of this module can be found in Sec VI-B.

### C. Planning Function Generation

A crucial component of our method is leveraging the LLM to perform high-level planning instead of handling low-level actions. We combine all the prompts from the previous modules and input them into the LLM. The LLM then generates a reasonable and executable planning function that can be directly utilized in the subsequent training process.

To be more concise, given any processed state $S_I$, we define an assignment planning function as $\mathscr{F}_{\mathscr{T}}(S_I) \to \mathscr{T}_i \in \mathscr{T}$, where $\mathscr{T} = \{\mathscr{T}_1, ..., \mathscr{T}_n\}$ is a set of target assignments that each agent can take. We define the assignment set $\mathscr{T}$ over the action space such that an action can belong to multiple assignments and vice versa. For example, if the assignment space is defined as $\mathscr{T} = \{Landmark\_0, Landmark\_1\}$, and landmark 0 and landmark 1 are located at the top right and top left positions relative to the agent respectively, then

taking the action "UP" can be associated with both assignments. Conversely, we can have multiple actions correspond to an assignment. For instance, moving towards "Landmark 0" may involve actions like "UP" and "RIGHT".

The *Planning Function Generation* will only be required once for each new environment you try to use. After you interact with the LLM to get generated planning function, you can directly use it in the later training process with different MARL algorithms. This is referred to the red module in Fig. 1(c).

### D. MARL training with Planning function incorporation

To incorporate the planning function into MARL training, we add an extra reward term to the original reward provided by environments. Specifically, we define the final reward $R$ used by the critic as:

$$R = r + \sum_i \Delta r_i. \tag{1}$$

Here, $r$ is the original reward from the environment. For each agent $i$, $\Delta r_i$ is an additional reward or penalty that determined based on whether the action taken by the agent aligns with the task assigned by the planning function. Specifically:

$$\Delta r_i = \begin{cases} r', \text{if } a_i \in \mathscr{T}_i \\ p', \text{if } a_i \notin \mathscr{T}_i \end{cases} \tag{2}$$

Notably, we don't need to interact with the LLM during the entire training process, nor do we need to call the planning function after the policy has been trained. The training process $MARL_{alg}(R)$ takes $R$ as the reward function, uses the same state and action space. We follow the standard MARL algorithms and evaluation metrics within the literature, such as [8], [9], and [7]. Our method, as shown in the green box in Fig. 1(d), is highly efficient compared to approaches that interact with LLMs throughout the whole training process or directly use LLMs as agents. In practice, using the LLM's API to generate the planning function incurs minimal cost—less than a dollar per environment—even when using the most advanced LLM APIs.

## V. EXPERIMENTS

In this section, we evaluate our method across two different environments: MPE and LBF.

### A. Setup

**Baselines.** In our experiments, we compare the MARL algorithm MADDPG [7], MAPPO [8] and QMIX [9] and set default hyper-parameters according to the well-tuned performance of human-written reward, and fix that in all experiments on this task to do MARL training. Experiment hyper parameters are listed in Appendix.

**LLM** As one of the most advanced LLMs available at the time of experimentation, Claude 3.5 Sonnet exhibited enhanced reasoning and code generation abilities critical for our framework's strategy formulation and planning functions. This performance advantage led us to select

`claude-3-5-sonnet-20240620` as our primary LLM for all reported experiments.[1]

**Metrics.** To assess the performance of our method, we use the mean return in evaluation for all other environments. During evaluation, we rely solely on the default return values provided by the environments for both the baseline and our method, ensuring a fair comparison.

### B. Results

**Level-Based Foraging.** LBF [1] is a challenging sparse reward environment designed for MARL training. In this environment, agents must learn to navigate a path and successfully collect food, with rewards only being given upon task completion. To evaluate our framework in a cooperative setting, we selected the 2-player, 2-food fully cooperative scenario. In this setting, all agents must work together and coordinate their actions to collect the food simultaneously. The environment offers an action space consisting of [NONE, NORTH, SOUTH, WEST, EAST, LOAD], and we define the task set as [NONE, Food i, ..., LOAD]. Using the relative positions of agents and food items, we map assigned tasks to the corresponding actions in the action space and calculate the reward based on this alignment. We evaluated our framework over 3 different seeds, with the results shown in Figure 2 and Table I. LLM assist the MARL algorithm by providing reward signals, our framework significantly outperformed the baseline, achieving a maximum improvement of **105 %** in mean return and a **2x faster** convergence rate among all tested MARL algorithms. According to the results, our framework is effective across all the baseline algorithms, with particularly large improvements observed in QMIX and MADDPG, and a faster convergence rate for MAPPO. To assess the variability in the quality of our generated functions, we present the results of three different generated functions in Figure 6 and Table II in Appendix. The results demonstrate that our framework consistently generates high-quality functions, with each achieving similar improvements across all baseline algorithms.

**Multi-Agent Particle Environment.** We evaluate our framework in MPE [7] simple spread environment which is a fully cooperative game. This environment has N agents, N landmarks. At a high level, agents must learn to cover all the landmarks while avoiding collisions. It's action space consists of [no_action, move_left, move_right, move_down, move_up]. We define the assignment for each agent to take to be [Landmark_i,...,No action]. During training, based on the global observation, we obtain the relative position of each agent with respect to the landmarks. Similar to LBF, we map each assignment of agent back to the corresponding action space and then reward the action of policy in action space level. We evaluate our approach on 3-agent and 4-agent scenarios using QMIX and MADDPG as baselines. As shown in Figure 3, our framework(colored line) outperforms the baseline(black line) algorithm in mean returns by **7.66%** and **8.8%** for 3-agent scenario, and **2.4%** and **18.09%** for

---

[1]`https://www.anthropic.com/news/claude-3-5-sonnet`

TABLE I: Comparison between YOLO-MARL and MARL in the LBF environment across three seeds. The highest evaluation return means during training are highlighted in bold. The corresponding results can be found in Figure 2. The M means one million training steps. We run all the experiments on the same machine.

| | Mean Return after 0.2M / 0.4M / 1.5M / 2M Steps | | |
| --- | --- | --- | --- |
| | QMIX | MADDPG | MAPPO |
| MARL | 0.00/ 0.01/ 0.25/ 0.38 | 0.09/ 0.33/ 0.26/ 0.32 | 0.31/ 0.72/ 0.99/ 0.99 |
| YOLO-MARL | **0.01/ 0.02 / 0.60/ 0.78** | **0.13/ 0.38/ 0.39/ 0.44** | **0.93/ 0.98/ 0.99/ 0.99** |



(a) MADDPG      (b) MAPPO      (c) QMIX

Fig. 2: **Results for LBF environment across 3 seeds:** The solid lines indicate the mean performance, and the shaded areas represent the range (minimum to maximum) across 3 different seeds.

4-agent scenario with QMIX and MADDPG respectively. These improvements demonstrate the effectiveness of our framework in enhancing coordination among agents to cover up all the landmarks.

## VI. ABLATION STUDY

In this section, we conduct the ablation studies mainly in LBF 2 players 2 food fully cooperative environment since rewards in LBF are sparser compared to MPE [1]. We refer to V-B for more information about the environment.

### A. Comparison between YOLO-MARL with and without Strategy Generation

In this section, we examine the impact of the Strategy Generation Module on the performance of the YOLO-MARL framework. Specifically, we compare the standard YOLO-MARL with a variant that excludes the Strategy Generation Module to assess its significance.

According to our tests, the Strategy Generation Module plays an important role in the YOLO-MARL method. As shown in Figure 4 , without the LLM generated strategy, we obtain a worse-performing planning function. Interestingly, the mean returns of evaluations for the functions without the LLM generated strategy are not always close to zero, indicating that the generated planning functions are not entirely incorrect. Based on this, we could confirm that the Strategy Generation Module would help Planning Function Generation Module provides better solutions to this game. Moreover, giving the strategy also helps stabilize the quality of the generated code. We observe a higher risk of obtaining erroneous functions without supplying the strategy.

### B. Comparison between YOLO-MARL with and without State Interpretation

To demonstrate how the State Interpretation Module enhances our framework, we present two failure case snippets:
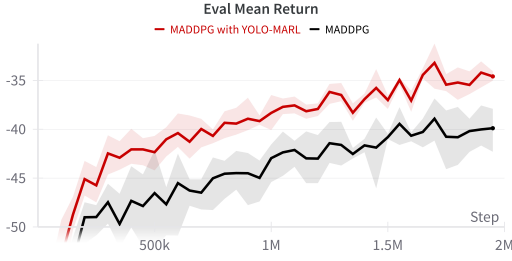
- Without the Interpretation Function: The interpretation function is omitted entirely from the prompting pipeline.
- Providing Raw Environment Code Directly: Raw environment source code is fed directly to the LLM.

The LLM is unable to infer the type of state and attempts to fetch environment information via a non-existent key if no preprocessing code provided. And if environment code is provided without dimensional context for each component, the LLM is likely to make random guesses. In both scenarios, the absence of explicit state interpretation hinders the LLM's ability to generate accurate and executable planning functions. These failures underscore the importance of the State Interpretation Module in bridging the gap between vectorized observations and the LLM's requirement for semantically meaningful input.
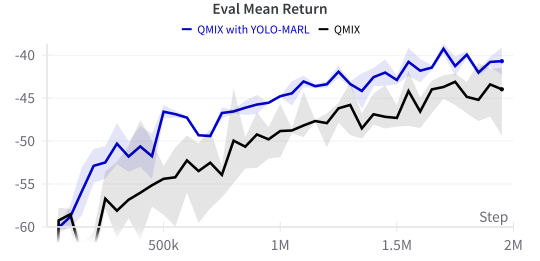
By incorporating the State Interpretation Module, we enable the LLM to understand the environment's state representation effectively. This results in the generation of reliable planning functions that significantly enhance the performance of our YOLO-MARL framework.

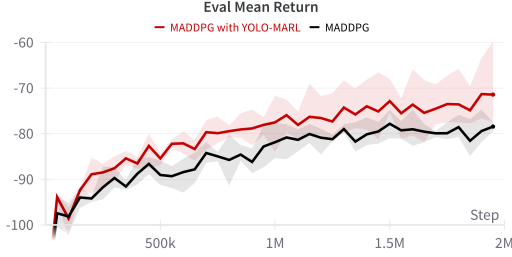### C. Comparison between YOLO-MARL and reward generation

In this section, we compare our YOLO-MARL method with approaches that utilize the LLM for reward generation without reward function template. We explore two scenarios: reward generation without feedback and reward generation with feedback. For the reward generation without feedback, the reward function is generated at the same stage as the planning function for fair comparison. This means that we generate the reward function before all the training process for each new environment. For the reward generation with feedback, we first generate a reward function just like the reward generation without feedback. And then, iteratively, we will run a whole training process on this environment
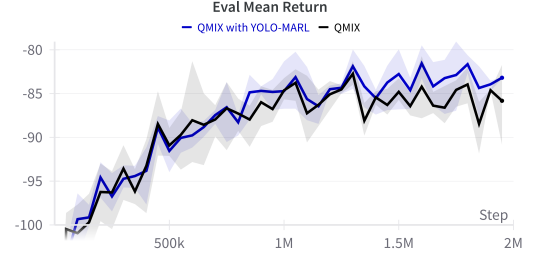
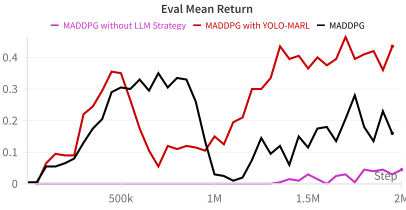(a) 3 Agents – MADDPG

(b) 3 Agents – QMIX
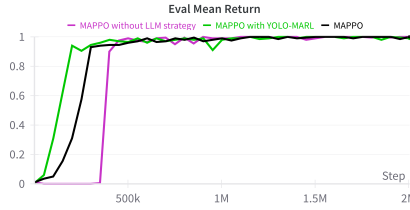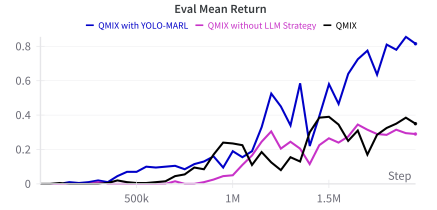
(c) 4 Agents – MADDPG

(d) 4 Agents – QMIX

Fig. 3: **Results for MPE simple spread environment** with 3 agents (top row a and b) and 4 agents (bottom row c and d). The solid lines indicate the mean performance, and the shaded areas represent the range (minimum to maximum) across 3 different generated planning functions.



(a) MADDPG

(b) MAPPO

(c) QMIX

Fig. 4: Comparison between YOLO-MARL with and without using LLM generated strategies in LBF

and pass the feedback of this training performance to the LLM, combined with previous prompts and ask the LLM to refine the previous generated reward function.

Our experiments show that relying solely on the LLM-generated reward function leads to poor performance. As shown in Figure 5, the mean return for the LLM-generated reward function pair consistently falls below the performance of all three MARL algorithms. This indicates that agents are not learning effectively under the LLM-generated reward function. However, we do observe a slight positive return. This suggest the potential of using this framework for reward shaping tasks, particularly in situations where standard MARL algorithms struggle to learn in sparse reward scenarios. To investigate whether iterative refinement could improve the LLM generated reward function, we supply the LLM with the generated reward function from the prior iteration and feedback on its performance. Despite this iterative process, the LLM still fails to output a suitable reward function for the LBF environment. The mean return of evaluations remains close to zero, as shown in figure 7.

## VII. CONCLUSION

We propose YOLO-MARL, a novel framework that enhances MARL policy training by integrating LLMs' high-level planning capabilities with a single interaction per environment. This design reduces computational overhead and mitigates instability issues associated with frequent LLM interactions during training. This approach not only outperforms traditional MARL algorithms but also operates independently of the LLM during execution, demonstrating strong generalization capabilities across various environments.

We evaluate YOLO-MARL across two different environments: the MPE environment and the LBF environment. Our experiments showed that YOLO-MARL outperforms or achieve competitive results compared to baseline MARL methods. The integration of LLM-generated high-level assignment planning functions facilitated improved policy learning in challenging cooperative tasks. Finally, we mention a possible way to incorporate reward generation to our framework and we will step further.
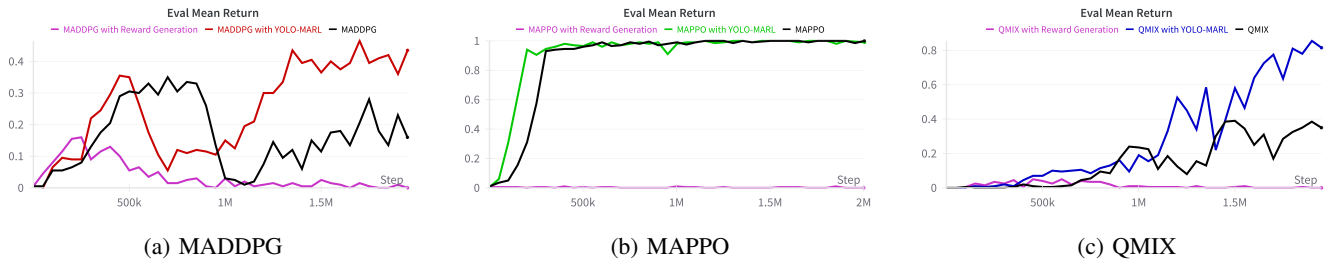
| (a) MADDPG | (b) MAPPO | (c) QMIX |

Fig. 5: Comparison between YOLO-MARL and reward generation without feedback in LBF

## REFERENCES

[1] G. Papoudakis and L. Schäfer, "Benchmarking Multi-Agent Deep Reinforcement Learning Algorithms in Cooperative Tasks," *arXiv preprint arXiv:2006.07869*, 2021.

[2] M. Ahn, A. Brohan, N. Brown, Y. Chebotar, O. Cortes, B. David, C. Finn, C. Fu, K. Gopalakrishnan, K. Hausman *et al.*, "Do as i can, not as i say: Grounding language in robotic affordances," *arXiv preprint arXiv:2204.01691*, 2022.

[3] H. Zhang, W. Du, J. Shan, Q. Zhou, Y. Du, J. B. Tenenbaum, T. Shu, and C. Gan, "Building cooperative embodied agents modularly with large language models," *arXiv preprint arXiv:2307.02485*, 2023.

[4] S. S. Kannan, V. L. N. Venkatesh, and B.-C. Min, "Smart-llm: Smart multi-agent robot task planning using large language models," *arXiv preprint arXiv:2309.10062*, 2024.

[5] Y. Du, O. Watkins, Z. Wang, C. Colas, T. Darrell, P. Abbeel, A. Gupta, and J. Andreas, "Guiding pretraining in reinforcement learning with large language models," *arXiv preprint arXiv:2302.06692*, 2023.

[6] M. Kwon, S. M. Xie, K. Bullard, and D. Sadigh, "Reward Design with Language Models," *arXiv preprint arXiv:2303.00001*, 2023.

[7] R. Lowe, Y. Wu, A. Tamar, J. Harb, P. Abbeel, and I. Mordatch, "Multi-Agent Actor-Critic for Mixed Cooperative-Competitive Environments," *arXiv preprint arXiv:1706.02275*, 2020.

[8] C. Yu, A. Velu, E. Vinitsky, J. Gao, Y. Wang, A. Bayen, and Y. Wu, "The Surprising Effectiveness of PPO in Cooperative, Multi-Agent Games," *arXiv preprint arXiv:2103.01955*, 2022.

[9] T. Rashid, M. Samvelyan, C. S. De Witt, G. Farquhar, J. Foerster, and S. Whiteson, "QMIX: Monotonic Value Function Factorisation for Deep Multi-Agent Reinforcement Learning," *arXiv preprint arXiv:1803.11485*, 2018.

[10] C. Sun, S. Huang, and D. Pompili, "Llm-based multi-agent reinforcement learning: Current and future directions," *arXiv preprint arXiv:2405.11106*, 2024.

[11] J. N. Foerster, G. Farquhar, T. Afouras, N. Nardelli, and S. Whiteson, "Counterfactual multi-agent policy gradients," *arXiv preprint arXiv:1705.08926*, 2017.

[12] P. Sunehag, G. Lever, A. Gruslys, W. M. Czarnecki, V. F. Zambaldi, M. Jaderberg, M. Lanctot, N. Sonnerat, J. Z. Leibo, K. Tuyls, and T. Graepel, "Value-decomposition networks for cooperative multi-agent learning," *arXiv preprint arXiv:1706.05296*, 2017.

[13] T. Carta, C. Romac, T. Wolf, S. Lamprier, O. Sigaud, and P.-Y. Oudeyer, "Grounding large language models in interactive environments with online reinforcement learning," *arXiv preprint arXiv:2302.02662*, 2023.

[14] F. Lin, E. L. Malfa, V. Hofmann, E. M. Yang, A. Cohn, and J. B. Pierrehumbert, "Graph-enhanced large language models in asynchronous plan reasoning," *arXiv preprint arXiv:2402.02805*, 2024.

[15] L. Fan, G. Wang, Y. Jiang, A. Mandlekar, Y. Yang, H. Zhu, A. Tang, D.-A. Huang, Y. Zhu, and A. Anandkumar, "Minedojo: Building open-ended embodied agents with internet-scale knowledge," *Advances in Neural Information Processing Systems*, vol. 35, pp. 18 343–18 362, 2022.

[16] J. Liang, W. Huang, F. Xia, P. Xu, K. Hausman, B. Ichter, P. Florence, and A. Zeng, "Code as policies: Language model programs for embodied control," *Proceedings of the IEEE International Conference on Robotics and Automation (ICRA)*, 2023.

[17] Y. J. Ma, W. Liang, G. Wang, D.-A. Huang, O. Bastani, D. Jayaraman, Y. Zhu, and A. Anandkumar, "EUREKA: HUMAN-LEVEL REWARD DESIGN VIA CODING LARGE LANGUAGE MODELS," *arXiv preprint arXiv:2310.12931*, 2024.

[18] T. Xie, S. Zhao, C. H. Wu, Y. Liu, Q. Luo, V. Zhong, Y. Yang, and T. Yu, "Text2Reward: Automated Dense Reward Function Generation for Reinforcement Learning," *arXiv preprint arXiv:2309.11489*, 2023.

[19] T. Guo, X. Chen, Y. Wang, R. Chang, S. Pei, N. V. Chawla, O. Wiest, and X. Zhang, "Large language model based multi-agents: A survey of progress and challenges," *arXiv preprint arXiv:2402.01680*, 2024.

[20] G. Li, H. Hammoud, H. Itani, D. Khizbullin, and B. Ghanem, "Camel: Communicative agents for" mind" exploration of large language model society," *Advances in Neural Information Processing Systems*, vol. 36, 2024.

[21] S. Hong, X. Zheng, J. Chen, Y. Cheng, J. Wang, C. Zhang, Z. Wang, S. K. S. Yau, Z. Lin, L. Zhou *et al.*, "Metagpt: Meta programming for multi-agent collaborative framework," *arXiv preprint arXiv:2308.00352*, 2023.

[22] S. S. Kannan, V. L. Venkatesh, and B.-C. Min, "Smart-llm: Smart multi-agent robot task planning using large language models," *arXiv preprint arXiv:2309.10062*, 2023.

[23] B. Yu, H. Kasaei, and M. Cao, "Co-navgpt: Multi-robot cooperative visual semantic navigation using large language models," *arXiv preprint arXiv:2310.07937*, 2023.

[24] S. Hu, T. Huang, F. Ilhan, S. Tekin, G. Liu, R. Kompella, and L. Liu, "A survey on large language model-based game agents," *arXiv preprint arXiv:2404.02039*, 2024.

[25] R. Gong, Q. Huang, X. Ma, H. Vo, Z. Durante, Y. Noda, Z. Zheng, S.-C. Zhu, D. Terzopoulos, L. Fei-Fei, and J. Gao, "Mindagent: Emergent gaming interaction," *arXiv preprint arXiv:2309.09971*, 2023.

[26] S. Wu, L. Zhu, T. Yang, S. Xu, Q. Fu, Y. Wei, and H. Fu, "Enhance reasoning for large language models in the game werewolf," *arXiv preprint arXiv:2402.02330*, 2024.

[27] H. Li, Y. Q. Chong, S. Stepputtis, J. Campbell, D. Hughes, M. Lewis, and K. Sycara, "Theory of mind for multi-agent collaboration via large language models," *arXiv preprint arXiv:2310.10701*, 2023.

[28] W. Ma, Q. Mi, X. Yan, Y. Wu, R. Lin, H. Zhang, and J. Wang, "Large language models play starcraft ii: Benchmarks and a chain of summarization approach," *arXiv preprint arXiv:2312.11865*, 2023.

[29] M. L. Littman, "Markov games as a framework for multi-agent reinforcement learning," *https://www.sciencedirect.com/science/article/pii/B9781558603356500271*, 1994.

[30] G. Owen, "Game theory," *https://books.google.com/books?id=pusfAQAAIAAJ*, 1982.

## VIII. APPENDIX

Given the page constraints, we present some additional experiments in this section.

TABLE II: Comparison between YOLO-MARL and MARL in the LBF environment across three different generated planning functions. The highest evaluation return means during training are highlighted in bold. The corresponding results can be found in figure 6. The M means one million training steps. We use two different machines to generate planning functions and run MARL and YOLO-MARL on the same machines where the planning functions are generated.

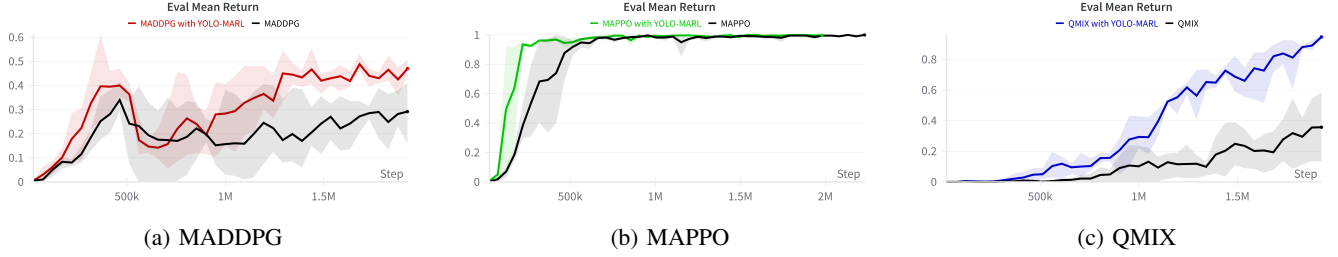| | Mean Return after 0.2M / 0.4M / 1.5M / 2M Steps | | |
| --- | --- | --- | --- |
| | QMIX | MADDPG | MAPPO |
| MARL | 0.00/ 0.01/ 0.25/ 0.36 | 0.08/ 0.28/ 0.24/ 0.29 | 0.38/ 0.74/ 0.99/ 0.99 |
| YOLO-MARL | **0.00/ 0.03/ 0.69/ 0.95** | **0.18/ 0.40/ 0.42/ 0.47** | **0.94/ 0.97/ 0.99/ 0.99** |



(a) MADDPG     (b) MAPPO     (c) QMIX

Fig. 6: **Results for LBF environment across 3 seeds:** The solid lines indicate the mean performance, and the shaded areas represent the range (minimum to maximum) across 3 different seeds.



(a) Iteration one     (b) Iteration two
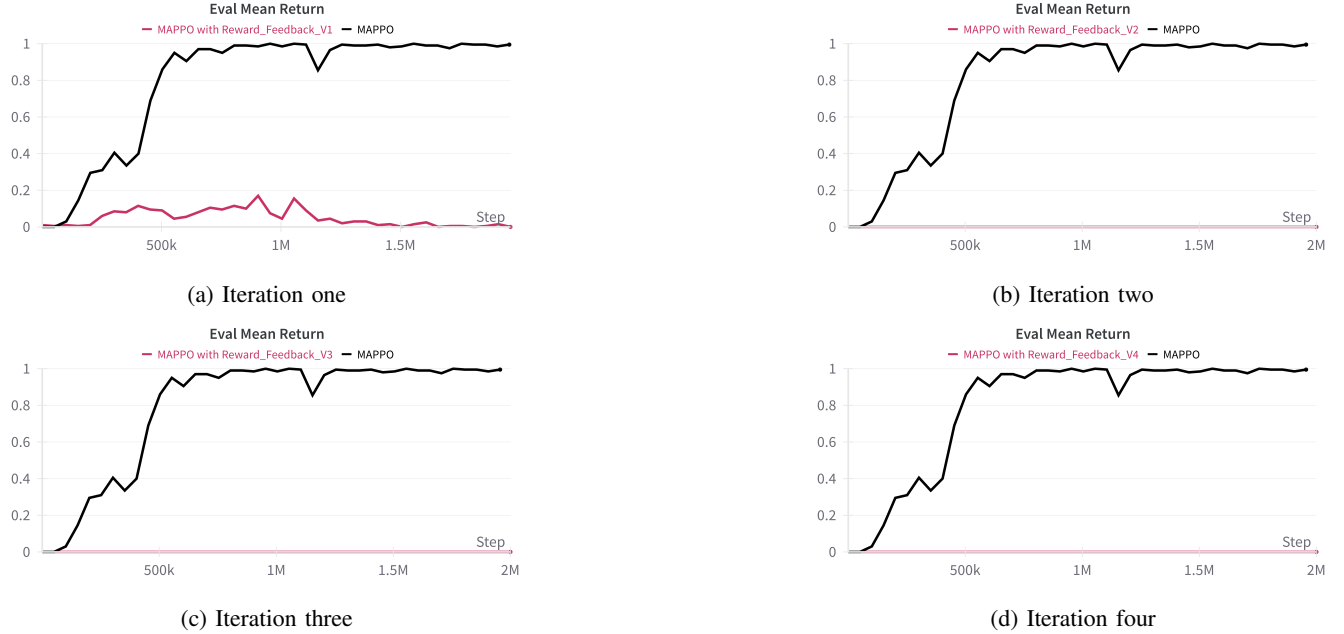
(c) Iteration three     (d) Iteration four

Fig. 7: Results of only reward generation with feedback in the LBF environment. The total number of iteration is 4 and the MARL algorithm we used here is MAPPO.