

Leveraging Large Language Models for Effective and Explainable Multi-Agent Credit Assignment

Kartik Nagpal

University of California Berkeley
Berkeley, CA, United States of America
kartiknagpal@berkeley.edu

Jean-Baptiste Bouvier

University of California Berkeley
Berkeley, CA, United States of America
bouvier3@berkeley.edu

Dayi Dong

University of California Berkeley
Berkeley, CA, United States of America
dayi.dong@berkeley.edu

Negar Mehr

University of California Berkeley
Berkeley, CA, United States of America
negar@berkeley.edu

ABSTRACT

Recent work, spanning from autonomous vehicle coordination to in-space assembly, has shown the importance of learning collaborative behavior for enabling robots to achieve shared goals. A common approach for learning this cooperative behavior is to utilize the centralized-training decentralized-execution paradigm. However, this approach also introduces a new challenge: how do we evaluate the contributions of each agent’s actions to the overall success or failure of the team. This “credit assignment” problem has remained open, and has been extensively studied in the Multi-Agent Reinforcement Learning (MARL) literature. In fact, humans manually inspecting agent behavior often generate better credit evaluations than existing methods. We combine this observation with recent works which show Large Language Models (LLMs) demonstrate human-level performance at many pattern recognition tasks. Our key idea is to reformulate credit assignment to the two pattern recognition problems of sequence improvement and attribution, which motivates our novel Large Language Model Multi-agent Credit Assignment (LLM-MCA) method. Our approach utilizes a centralized LLM reward-critic which numerically decomposes the environment reward based on the individualized contribution of each agent in the scenario. We then update the agents’ policy networks based on this feedback. We also propose an extension LLM-TACA where our LLM critic performs explicit task assignment by passing an intermediary goal directly to each agent policy in the scenario. Both our methods far outperform the state-of-the-art on a variety of benchmarks, including Level-Based Foraging, Robotic Warehouse, and our new “Spaceworld” benchmark which incorporates collision-related safety constraints. As an artifact of our methods, we generate large trajectory datasets with each timestep annotated with per-agent reward information, as sampled from our LLM critics. By making this dataset available, we aim to enable future works which can directly train a set of collaborative, decentralized policies offline.¹

¹Project Page is available at <http://iconlab.negarmehr.com/LLM-MCA/>

This work is licensed under a Creative Commons Attribution International 4.0 License.

Proc. of the 24th International Conference on Autonomous Agents and Multiagent Systems (AAMAS 2025), Y. Vorobeychik, S. Das, A. Nowé (eds.), May 19 – 23, 2025, Detroit, Michigan, USA. © 2025 International Foundation for Autonomous Agents and Multiagent Systems (www.ifaamas.org).

KEYWORDS

Credit Assignment, Task Allocation, Multi-Agent Reinforcement Learning, Large Language Models, Foundation Models

ACM Reference Format:

Kartik Nagpal, Dayi Dong, Jean-Baptiste Bouvier, and Negar Mehr. 2025. Leveraging Large Language Models for Effective and Explainable Multi-Agent Credit Assignment. In *Proc. of the 24th International Conference on Autonomous Agents and Multiagent Systems (AAMAS 2025)*, Detroit, Michigan, USA, May 19 – 23, 2025, IFAAMAS, 11 pages.

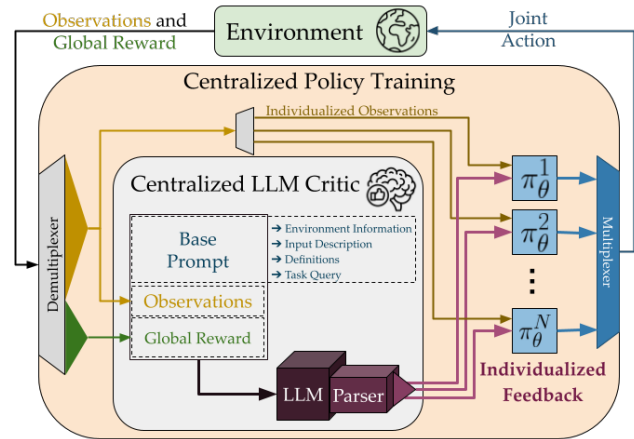


Figure 1: Overall architecture diagram for our LLM-MCA and LLM-TACA methods. Our centralized training architecture utilizes a centralized LLM-critic instantiated with our base prompt (environment description, our definitions, and task query). At each timestep, we update our LLM-critic with the global reward and latest observations from the environment. We then update our agents’ policies with the individualized feedback from our critic.

1 INTRODUCTION

Many real-world scenarios can be posed as multi-agent cooperative problems, where multiple agents are either required to collaborate or benefit from collaborating towards a shared goal. For instance,

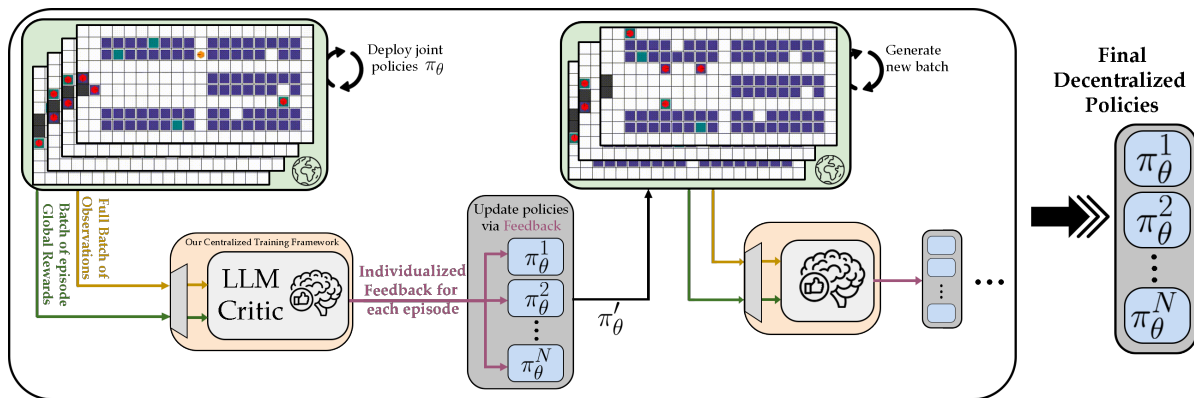


Figure 2: Diagram for our batch-training process with our LLM-MCA method. Our centralized training process allows us to provide entire batches of trajectories at once to our centralized LLM-critic. Our LLM-MCA critic then generates individualized feedback for each agent, which we use to update their policies. After training, we no longer need our LLM-critic, and directly deploy our trained, decentralized agent policies.

moving a heavy object which requires two robots to lift simultaneously, or minimizing city traffic which benefits from vehicles coordinating their routes. While game-theoretic approaches have shown remarkable results for some multi-agent problems [2], they rarely focus on collaborative settings. This has prompted many researchers to use Multi-Agent Reinforcement Learning (MARL) [55, 60] to coordinate large numbers of agents in complex scenarios.

The Centralized-Training Decentralized-Execution (CTDE) paradigm is a popular framework in MARL where agents are trained together via a centralized critic with access to the global states, actions, and rewards [6, 12, 13, 27, 56]. Then during deployment, each autonomous agent only has access to local observations. A primary motivation for the CTDE framework is the nonstationarity of naive decentralized training, as even when a given agent makes no policy changes, another agents’ policy modification can vastly alter the team’s global reward. The CTDE paradigm’s centralized critic can more clearly account for these changes and thus decrease training volatility. Furthermore, the centralized critic can better promote and exploit coordination among agents. Finally, since execution is decentralized, the learned policies for the multi-agent system can also be easily scaled.

A key challenge during the central training phase is how to decouple each policy changes’ effects and assess each agent’s contribution to the overall success or failure of the global task. Traditionally, the environment provides only a collective reward to the agents based on whether they achieved the shared objective or not. From this single reward, a CTDE training algorithm must determine the contributions of each agent and update the agent policies based on their respective actions [55]. This is commonly referred to as the *Structural Credit Assignment* problem [1, 48], not to be mistaken with the temporal credit assignment problem where one must determine the contributions of a succession of actions taken by a single agent receiving a single reward at the end. Note that the centralized critic must receive and process the joint observations and actions of the agents, which are often of large dimension, and must make assignments based on only partial exploration of

the massive joint state-action space, making the contributions difficult to evaluate [38, 58, 61]. Furthermore, this problem is made more challenging by sparse-reward environments, as the agents must learn to solve an entire sub-task or task before getting any reward feedback. As a result credit assignment has remained an open problem, with other works suffering from limitations such as low-quality feedback, low action influence [36], and difficulty with complex interactions [28].

In this work, we reformulate the credit assignment problem as a general pattern recognition problem. More specifically, we view the joint observations and actions of the agents as forming a numerical sequence, which is associated with the global reward signal it earns. Given enough of these sequences, the act of credit assignment is akin to sequence improvement —evolving beyond the original demonstrations to better achieve goals, generalize to new environments, or increase efficiency/robustness—in terms of the global reward [31]. We propose that this sequential pattern-improvement task can be viewed as giving feedback to each agent regarding their performance during the task. In this work, we claim that the underlying pattern seen in collaborative multi-agent scenarios is akin to an “Agreement Problem” where the agents must learn to collectively agree on the strategies needed to effectively collaborate. With this in mind, we label a few common forms of disagreements seen in multi-agent systems and provide them to our centralized critic. Our primary motivation for posing credit assignment in this pattern recognition form is to leverage human-like skills when assigning credit even in multi-step tasks thanks to their pattern recognition skills and complex reasoning capabilities [45].

While these abilities have long been a trademark of human intelligence, recent works have shown that foundational Large Language Models (LLMs) display similar properties [31]. Due to the neuro-morphic characteristics of this learning architecture, some works argue that LLMs can extract the underlying pattern within the text and act as general pattern machines [31] just like humans [26]. Further works have shown how LLMs exhibit human-level intelligence and pattern recognition capabilities [16, 18] when evaluated

on common human benchmarks such as the International Math Olympiad and the Bar Exam [10, 23, 35]. Furthermore, LLMs have already been successfully utilized for temporal credit assignment in the single-agent case [36]. Building on these accomplishments, we propose to leverage the pattern recognition skills of Large Language Models (LLMs) to tackle the multi-agent structural credit assignment problem.

Specifically, we utilize an LLM as our centralized critic, and instantiate it with textual context of the environment, definitions of our agreement problem, and description of its role. We then periodically provide our critic with joint observations, joint actions, and global rewards from the environment. In return, our LLM-critic provides individualized numerical feedback for each agent and explanations for these assignments. We employ a parser function that takes the feedback from the LLM and maps it to the numerical reward value for each agent at each timestep during the trajectory. We further extend this method by adjusting our query to the LLM as well as our parser function such that during training, each agent can also be given an explicit task from the LLM critic. We call this extension our Large Language Models for Task and Credit Assignment (LLM-TACA) method. Both our approaches perform exceptionally well, far surpassing current state-of-the-art methods [38, 58, 61] on common benchmarks [7, 34] while also providing more interpretable agent feedback. Lastly, we introduce "Spaceworld," an environment for modeling the multi-agent in-space assembly scenario with safety constraints [32].

In summary, our main contributions are:

- (1) A novel reformulation of the credit and task assignment problems as a general pattern recognition problem.
- (2) A novel LLM-based centralized training method called "LLM-MCA" that performs credit assignment during training and far surpasses the state-of-the-art at many common MARL benchmarks.
- (3) An extension of our LLM-critic framework to perform explicit task assignment alongside the credit assignment called "LLM-TACA" which further improves performance and explainability.
- (4) A new offline dataset with per-robot task and reward annotations to aid future offline-learning efforts, as well as a new MARL benchmark "Spaceworld" which simulates a cooperative in-space assembly scenario.

2 RELATED WORKS

We start by introducing the literature on the credit assignment problem and how large language models have been used in reinforcement learning.

2.1 Credit Assignment

One of the main challenges of Multi-Agent Reinforcement Learning (MARL) is understanding how each agent should be rewarded for their respective participation toward the fulfillment of their shared goal [55]. This credit assignment problem arises from the difficulty to discern which agent or action deserves credit when dealing with long time horizons or collaborative tasks [39]. These respective *temporal* and *structural* credit assignment problems have

been investigated with reward prediction to assign values to key-actions and adjacent actions [40]. Some work has also used agent-specific utility functions to demonstrate the equivalence between temporal and structural credit assignment problems [1].

Recent works have built on the notion of difference rewards, which are rewards designed specifically to enable agents to estimate their respective contribution to the shared reward [11, 33, 41]. Another line of work uses classical cooperative game-theoretic approaches to assign credit to agents [17, 53] while other works have used counterfactual reasoning by comparing the outcomes to what would have happened if agents had not acted [13, 37]. However, training such a fully centralized state-action value becomes impractical when there are more than a handful of agents.

A more promising subset of work has focused on learning a centralized but factored Q-function, like value decomposition networks (VDN) [47], QMIX [38], and QTRAN [43], where the agents' local Q-functions can be extracted from a centralized Q-function. However, the structural assumptions enabling this factorization limit the complexity of the centralized action-value functions and often leads to subpar performance [29]. This observation led to the development of MAVEN [29] to fight inefficient exploration of Q-DPP [57], which bypasses structural factorization methods by encouraging agents to acquire sufficiently diverse behaviors that can be easily distinguished and then factorized. Similar to QMIX, LICA [61] uses a hyper-network to represent the centralized critic but without the monotony assumption of QMIX [38]. While not explicitly designed to solve the credit assignment problem, multi-agent PPO (MAPPO) [58] has shown strong performance in a variety of MARL tasks and will serve as one of our baselines.

2.2 Large-Language Models in Reinforcement Learning

Large language models are Transformer-based language models [52] trained on immense amounts of text and capable of providing impressive results for a wide array of tasks including text generation [3], question answering [24], and text summarization [59]. In the field of reinforcement learning, LLMs have often been used to generate training priors and for improving agent generalization to new domains [25]. LLMs have also shown promise as zero-shot planners [20], hierarchical planners [22, 49], reward shapers [4], and reward function generators [15] based on natural language instruction. Closely related to this work, LLMs have also been utilized both as binary critics for solving the temporal credit assignment problem [36], and used to help linguistic agents coordinate sub-task plans and waypoint paths through rounds of discussions [30].

LLMs have enabled improvements in reinforcement training, including facilitating inter-agent coordination by dictating high-level task planning [62] and negotiation processes [5, 46]. Additionally, approaches like [42] have used a centralized critic to allow for natural language communication and [19] employs a message pool that agents use to gather relevant messages. While we draw from these works, our approach uses numeric feedback to train neural network policies, as opposed to working with linguistic agents.

3 OUR METHOD

While the centralized training aspect of the centralized-training decentralized-execution (CTDE) paradigm enables agents to learn to coordinate, it also gives rise to the credit assignment problem [12]. Indeed, distributing the shared reward among agents based on their respective contributions towards the overall objective requires a deep understanding of the shared task. To address this problem, we introduce a novel LLM-based centralized critic to perform credit and task assignments during training.

3.1 Problem Formulation

We formalize our scenarios as decentralized partially observable Markov decision processes (Dec-POMDPs) $\mathcal{M} := \langle \mathcal{O}, \mathcal{S}, \mathcal{A}, Z, T, R, \rho \rangle$, where \mathcal{O} , \mathcal{S} , and \mathcal{A} are respectively the joint observation, state, and action spaces, $T : \mathcal{S} \times \mathcal{A} \rightarrow \mathcal{S}$ models the transition probabilities, $Z : \mathcal{S} \rightarrow \mathcal{O}$ is the observation function, $R : \mathcal{O} \times \mathcal{A} \rightarrow \mathbb{R}$ denotes the global reward function, and ρ denote the distribution of initial states $s_0 \in \mathcal{S}$. We consider N agents indexed by the set $\mathcal{N} := \{1, \dots, N\}$. At timestep k , agent $i \in \mathcal{N}$ is in a state $s_k^i \in \mathcal{S}^i$ from which it only observes $o_k^i = Z^i(s_k^i) \in \mathcal{O}^i$ and takes action $a_k^i \sim \pi^i(a_k^i | o_k^i) \in \mathcal{A}^i$ following its decentralized policy $\pi^i : \mathcal{O}^i \rightarrow \mathcal{A}^i$. The joint observation, state, and action spaces are simply the product of their corresponding subspaces for each agent: $\mathcal{O} := \mathcal{O}^1 \times \mathcal{O}^2 \times \dots \times \mathcal{O}^N$, $\mathcal{S} := \mathcal{S}^1 \times \dots \times \mathcal{S}^N$, and $\mathcal{A} := \mathcal{A}^1 \times \dots \times \mathcal{A}^N$.

At any timestep k , the joint action of the agents $a_k := (a_k^0, \dots, a_k^N)$ takes their joint state $s_k := (s_k^0, \dots, s_k^N)$ to $s_{k+1} \sim T(s_{k+1} | s_k, a_k)$. At each time step k , the team of agents receives a single collective scalar reward $r_k := R(o_k, a_k) \in \mathbb{R}$. We denote the joint set of decentralized agent policies by $\pi_\theta := (\pi^1, \dots, \pi^N) : \mathcal{O} \rightarrow \mathcal{A}$, where θ denotes the aggregate of the parameters describing each of the N neural networks modeling the N different policies. Our objective is to maximize the expected total reward over this joint policy:

$$\max_{\theta} \mathbb{E}_{\substack{s_0 \sim \rho \\ a_k \sim \pi_\theta(o_k)}} \left(\sum_{k=0} R(o_k, a_k) \right), \quad (1)$$

where $o_{k+1} \sim Z(T(s_{k+1} | s_k, a_k))$. The challenge in optimizing (1) resides in coordinating policy updates when changing one set of policy parameters affects the reward earned by the whole team, which leads us to our problem of interest.

PROBLEM 1 (CREDIT ASSIGNMENT). *How can we determine the appropriate credit c_k^i to assign to each policy π^i at each timestep k given the shared reward r_k ?*

To address Problem 1 in the CTDE paradigm, a centralized critic distributes the credit c_k^i to each agent i . At the beginning of training, this critic first rewards agents for simply learning to operate under the rules of the environment. Once this knowledge is ingrained in the agents, they can begin to solve simple sub-tasks.

At the next stage of training, collaboration should start arising among the agents which additionally introduces intricate issues regarding whether our decentralized agents are collaborating either too little or too much. We say that a joint set of policies π_θ suffers from *under-collaboration* (respectively *over-collaboration*) if a number j of its decentralized agents coordinate to attempt a sub-task

requiring $m > j$ (resp. $m < j$) agents. A large challenge in cooperative MARL scenarios is overcoming the initial under-collaboration of decentralized agents. However, an excessive push in this direction will lead to the opposite pitfall of over-collaboration resulting in inefficiencies and agent conflicts. The optimal joint policy π_θ strikes the right amount of collaboration to achieve efficiency when the correct number of agents agree to adopt collaborative strategies. We define this challenge as the ‘‘Agreement Problem’’.

It is rather clear that solutions to this ‘‘Agreement Problem’’ are the most efficient outcomes, and thus are associated with the strategy producing the largest global cumulative reward $\sum_{k=1} r_k$ from the environment. Our goal, then, is to map a full sequence of observations, actions, and global reward values from our environment to a set of individualized feedback signals, which reflect how well each agent ‘‘agrees’’ to the same strategy. In essence, we ask our centralized critic to observe this sequence, propose improvements, and then map those improvements back to the agents that can perform them. Under this lens of the ‘‘Agreement Problem’’, we pose the credit assignment problem as a sequence improvement problem, which is then followed by an attribution step, where the critic must choose which among the agents are helping or hurting the overall agreement toward the correct strategy.

Now that we have re-framed the credit assignment problem as an agreement problem, we will introduce our solution.

3.2 Large Language Models for Multi-agent Credit Assignment (LLM-MCA)

Following our reintroduction of Problem 1 into a sequence improvement problem and an attribution step, both of which are pattern recognition tasks, we propose to leverage the strength of Large Language Models (LLMs) as general pattern machines [31], by designing a novel multi-agent LLM-based centralized-critic architecture, as illustrated in Figure 1.

As previously discussed, at the beginning of training, the primary goal of our LLM credit assignment critic is to adjust the sparse rewards from the environment by implicitly generating sub-tasks and rewarding the agents whenever they achieve these sub-goals. In essence, the LLM critic makes the reward signal denser by providing intermediary rewards. This step is similar to the single-agent temporal credit assignment problem in which LLMs have been shown to work well [36]. This process inherently requires an understanding of the overall goal of the scenario, which we provide via our initial text prompt. As current methods are unable to effectively utilize this prior, our approach immediately has an advantage. Furthermore, the LLM gets its own feedback via the global reward signal it receives from the environment. To avoid both of the pitfalls of under- and over-collaboration, we inform the LLM critic of the existence of these pitfalls directly within our base prompt as seen in Figure 3.

These text descriptions of common issues enable the LLM critic to utilize our vocabulary to explain its feedback. We also provide an example trajectory along with these definition prompts to help the LLM understand the issue. Similarly, the agreement problem is explained to the LLM in the form of a concrete example in its current environment, as reproduced below.

DEFINITION PROMPT 1. *Given a set of observations and actions as performed by the agents in the environment, there will be times when the agents will individually accomplish goals, but occasionally, we will need them to collaborate. When two agents correctly agree to collaborate on a task that requires two agents, then they have found a valid solution to the “Agreement Problem”.*

Specifically, we instantiate the LLM critic with a base prompt $p_{\text{base}} \in \mathcal{P}$ where \mathcal{P} is the prompt space. This base prompt is a concatenation of four sub-prompts $p_{\text{base}} := (p_{\text{env}}, p_{\text{desc}}, p_{\text{defn}}, p_{\text{task}})$ tailored for the given environment. The first element p_{env} is a description of this scenario and the rules of the environment. This also includes a description of the scenario’s objective and what will garner the agents better rewards. p_{desc} is an explanation of the kinds of data it will receive from this environment and the agent policies. Following our discussion above, we also provide the LLM critic with a Definition Prompt for under-collaboration, over-collaboration, and the Agreement problem 1. Combining these with traditional definitions of the temporal credit assignment and structural multi-agent credit assignment problems, we have our definition prompt p_{defn} . We also provide the LLM with a task description p_{task} which details its role and objective in the training process, as well as the precise form of its expected output at each timestep. Then, at each training timestep, we provide our LLM-critic C_{LLM} with dynamic data from the environment consisting of past observations, actions, and rewards. See Figure 3 for a more detailed example of our prompt for the “Spaceworld” benchmark.

At each timestep of the training process, we provide the next observation $o_{k+1} = Z(f(s_k, a_k)) \in \mathcal{O}$ and global reward $R(o_k, a_k) \in \mathbb{R}$ in the form of a prompt p_k . Due to the nature of LLM conversations, each new prompt p_{k+1} is effectively augmented with the history of past prompts and feedback, along with the new data, such that $p_{k+1} := (p_k, R(o_k, a_k), o_{k+1}, a_{k+1})$ where $p_0 := (p_{\text{base}}, o_0, a_0)$. The LLM-critic then provides a numerical reward value c_k^i to each agent in the scenario, along with an explanation of how it generated the feedback for that agent. This explanation is often associated with our “Agreement Problem” paradigm and how the LLM’s feedback helps mitigate disagreements among agents. An example of this behavior can be seen in Figure 4

As our LLM-critic $C_{LLM} : \mathcal{P} \times \mathbb{R} \times \mathcal{O} \times \mathcal{A} \rightarrow \mathcal{P}$ maps only to the text space, we utilize a parser function F_{MCA} which takes the entire output of the LLM critic and extracts the array of individualized credit values $c_k \in \mathbb{R}^N$. Because we include requirements for the format of this array within our base’s task prompt p_{task} , our parser function is comprised of a simple regex expressions search and a normalization step. As such, we denote our parsed LLM-critic $F_{MCA} \circ C_{LLM} : \mathcal{P} \times \mathbb{R} \times \mathcal{O} \times \mathcal{A} \rightarrow \mathbb{R}^N$, and explicitly have $F_{MCA}(C_{LLM}(p_k, r_k, o_{k+1}, a_k)) = (c_k^1, \dots, c_k^N)$, where each c_k^i is the credit assigned to each agent $i \in N$ for their contribution to the shared reward r_k .

Then, instead of directly attempting to maximize (1), we optimize the surrogate objective generated by our LLM critic:

$$\max_{\theta} \mathbb{E}_{\substack{s_0 \sim \rho \\ a_k \sim \pi_{\theta}(o_k)}} \left(\sum_{k=0}^{\infty} \sum_{i=1}^N F_{MCA}(C_{LLM}(p_k, R(o_k, a_k), o_{k+1}, a_k)) \right), \quad (2)$$

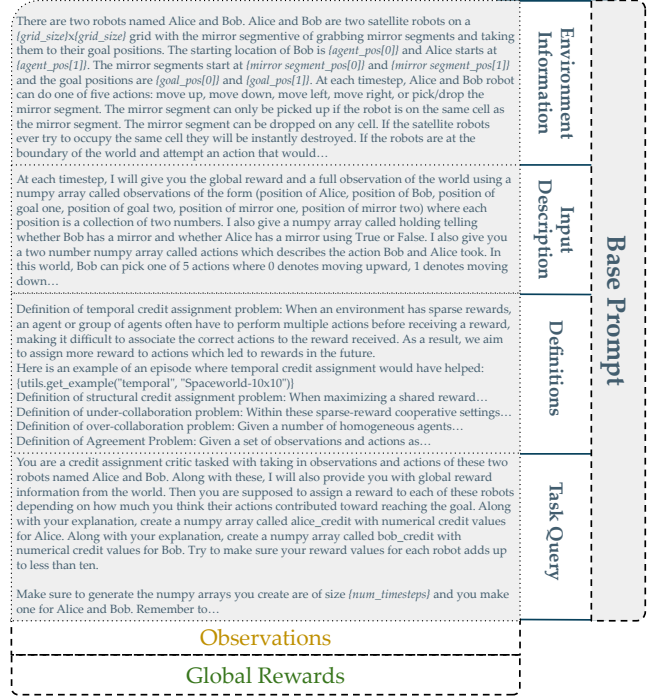


Figure 3: Example prompt for LLM-MCA in the “Spaceworld” benchmark. Our base prompt $p_{\text{base}} := (p_{\text{env}}, p_{\text{desc}}, p_{\text{defn}}, p_{\text{task}})$ is divided into (1) a description of the scenario’s rules and objectives, (2) a description of the kinds of inputs it will now receive from that environment, (3) our agreement problem definitions with examples, and (4) a description of its role as a credit assignment agent along with the formatting requirements of its output.

where $o_{k+1} \sim Z(T(s_{k+1} | s_k, a_k))$. This formulation allows our centralized LLM critic to directly evaluate the complex interplay behavior between agents.

Since we optimize our neural-network policy parameters through batch updates, as opposed to after every timestep, we can instead compute our credit assignment only at these policy update steps. As such, our surrogate objective can be further relaxed, where we compute feedback on an entire trajectory or batch of trajectories at once. We illustrate this batch-training process in Fig. 2. With this in mind, we adjust our base prompt p_{base} to describe this updated task, concatenate all of the observations, actions, and global rewards to a single prompt $p_{\text{batch}} = (p_{\text{base}}, p_0, p_1, \dots)$, and extract a set of credit matrices from the LLM critic at once. By providing an entire set of agent trajectories at once, our LLM-critic can better analyze the agents’ actions in retrospect, which further enables it to notice more intricate patterns within the agents’ behavior. Another major benefit to our batching technique is that it greatly reduces compute times and the data storage burden.

This LLM-MCA method also has several other inherent benefits. Since our LLM-critic is innately language-based, it is also able to explain and justify their feedback strategy, as seen in Figure 4. This explainability contrasts all previous multi-agent credit assignment works, where the crediting strategy often results from an uninterpretable deep neural network optimization, making the process

a complete black-box for the user. Furthermore, since our critics are pre-trained foundation models, they do not require any major hyperparameter tuning.

In summary, our LLM-MCA method can directly map entire trajectories of agents to individualized feedback for each agent, in a computationally efficient and interpretable manner. After extracting the credit values from this feedback, we can directly train our set of decentralized policies.

3.3 Large Language Models for Task and Credit Assignment (LLM-TACA)

Based both upon our interpretation of credit assignment as an ‘‘Agreement Problem’’ and our analysis of the outputs of our LLM-MCA method, we noticed that during the credit assignment process, our critic often implies a task assignment, and disincentivizes the agent from deviating from its internal task allocation strategy. From the point of view of the agent, whenever their learned task assignment differs from the critic’s the credit earned decreases dramatically, causing training volatility. Based on these findings, we introduce our extension LLM-TACA which includes a novel explicit task assignment aspect alongside the centralized credit critic.

We instantiate our LLM-TACA with an adjusted base prompt $p_{\text{base}}^T = (p_{\text{env}}, p_{\text{desc}}^T, p_{\text{defn}}, p_{\text{task}}^T)$, where p_{desc}^T now also includes a description of the task assignment the LLM is allowed to provide and p_{task}^T includes the additional request to generate these individualized task goals. Practically, we allow LLM-TACA to provide additional inputs directly to the agent policies during training time. An example of this assignment is shown in Figure 4. The exact form of the task inputs is a design choice, but we primarily used short integer arrays. At each timestep, on top of the credit feedback, LLM-TACA was given the option to either provide an additional task assignment input t_k directly to the policy of a given agent, or to provide nothing. To extract these features from the overall feedback and explanation from our LLM, we introduce a new parser function

$$F_{TACA}(C_{LLM})(p_k^T, R(o_k, a_k), o_{k+1}, a_k) = (c_k^1, t_k^1, \dots, c_k^N, t_k^N). \quad (3)$$

However, as this task assignment input is only provided during the centralized training and not during the decentralized execution, we must ensure that our agent policies don’t put too much weight on these inputs. To achieve this, we construct our neural network-based policies with a zero-padded input space and controlled dropout, as seen in early neural network literature [14, 44]. More specifically, over the course of the training, we input fewer and fewer task assignments from LLM-TACA by instead overwriting this data with zeros and also increasing the dropout rate of the neurons associated with these inputs.

In summary, our LLM-TACA method allows us to provide credit assignment and explicit task assignment information directly to each agent’s policy during training, while still maintaining the ability to execute them in a decentralized manner.

4 RESULTS

Now that we have introduced our two methods LLM-MCA and LLM-TACA for multi-agent credit assignment, we will evaluate them on a range of discrete, sparse-reward environments against

```

Here's a structured credit assignment based on the given observations and rewards:

1. **Temporal Credit Assignment**: The global reward is sparse but accumulates over time. Since the final reward is significantly higher (11), it likely resulted from earlier good decisions. We should assign more credit to actions that led to this.

2. **Structural Credit Assignment**: We need to ensure that each agent is rewarded based on their actual contributions and avoid freeloading.

3. **Collaboration**: Since some apples require cooperation, if Alice or Bob failed to assist when necessary, they should be penalized. Over-collaboration should also be avoided if unnecessary.

### Credit Assignment:
- The global rewards increase at the end, meaning critical actions were taken.
- As Alice and Bob worked together effectively, they both get positive rewards.
- If one agent was inefficient or blocked the other, they get a penalty.

### Task Assignment for Next Step:
- Based on the apple positions and levels, Alice and Bob should split tasks efficiently.
- If an apple requires collaboration, the agent who failed before should be assisted.
- Otherwise, each agent should move toward the nearest apple they can pick independently.

python
import numpy as np

alice_credit = np.array([3, 1, 1, 1, 1, 1, 1, 1, 1, 5]) # Alice contributed to the final large reward and an intermediate one
bob_credit = np.array([2, 1, 1, 1, 1, 1, 1, 1, 1, 5]) # Bob also contributed similarly

# Task allocation to prevent over-collaboration and promote efficiency
task_alice = np.array([0, 4]) # Alice should go to (0,4) to pick an apple
task_bob = np.array([1, 4]) # Bob should go to (1,4) to pick another apple
...

### Explanation:
- **Alice gets slightly more credit** because she contributed slightly more in the later stages where the major reward was received.
- **Bob still gets a fair amount** because collaboration was required.
- **Future task assignments direct Alice and Bob to specific apples to prevent over-collaboration and increase efficiency**.

This balances the need for individual credit while maintaining overall teamwork efficiency. 🍏🍏

```

Figure 4: Example output from our LLM-TACA method in the level-based foraging benchmark. Our LLM-critic provides individualized credit assignments for the previous timesteps, task assignments, and explanations for its decisions.

state-of-the-art baselines in MARL. In these results, we seek to show the efficacy of our framework rather than that of any specific language foundation model. Thus, we do not compare the performance of different large language models and only use the Gemma-7B model [50] for our trials. In this work, we use an open model to exploit the monetary and ease-of-access benefits, but rudimentary testing suggested that most popular foundation models could be used to generate similar results. Additionally, as all LLMs’ reasoning capabilities improve, the performance of our framework should similarly improve. All of the following results use the double DQN architecture [51] for the agent’s policies, which are trained and run on a NVIDIA 3080 Ti.

Since LLMs generate text, their outputs are interpretable for human users and their choices are thus explainable. This added benefit enables our centralized LLM-critic to explain its decisions when providing reward feedback values, especially how they are aligned with our ‘‘Agreement Problem’’ as discussed in Section 3.2. More specifically, with each alteration our LLM-critic makes, it often provides explanations for each of its actions in terms of how they help mitigate one of the problems as defined in the p_{defn} part of our base prompt, as illustrated in Figure 4. Furthermore, during early points in training, our LLM critics often rewards the agents for accomplishing some internal sub-goal and mention that they are solving the sparse reward or temporal credit assignment problem.

We first test our methods on a fully observable, cooperative climbing matrix game with two agents and dense rewards [9]. This repeated matrix game requires both agents to find an optimal policy that maximizes their rewards as described by the game’s common-payoff matrix illustrated in Figure 5. Each episode is 25 timesteps

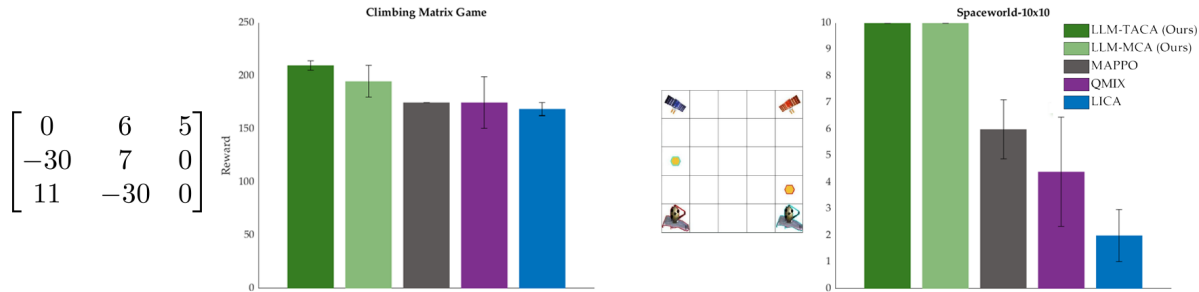


Figure 5: Comparison of our methods with baselines on benchmarks. All values reported are averaged across five separate trainings and the 95% confidence interval is illustrated by the error bars. Spaceworld is our custom benchmark where agents must avoid collision while transporting respective parts to their destinations."

in length and agents are given constant observations. The objective is for agents to learn to maximize their common-payoff by consistently agreeing to play the action profile which provides the highest common reward. It is easy to see that there exists multiple suboptimal equilibria, and so myopic critics may fail to train for the truly optimal solution [8].

Next, we evaluate our approach on the Level-Based Foraging (LBF) environment [7, 34] as seen in Figure 6(a). In this partially observable grid-world environment, multiple agents cooperate to collect apples from trees scattered within the area. Each agent has an associated level and to harvest a given apple the combined levels of the agents next to the tree must be greater than or equal to the level of the food. The environment name lets you customize many elements including the number of agents, the observability, the grid size, cooperativity, and the maximum number of food elements. The exact environment parameters we used for our scenarios can be seen directly in the environment names labeled in Figure 6.

Our third testing environment is Robotic Warehouse (RWARE) [7, 34] as shown in Figure 6(b). In this partially observable environment, robots are tasked with locating and delivering shelves to human workstations before returning them to empty shelf locations. However, as each robot only gets a 3x3 observation centered on themselves it is not always clear which area the agent should traverse. Additionally, there are sparse rewards and agents can get in each others way, causing needless slowdowns and inefficiencies. This environment is available with 3 different overall grid sizes and the number of robots can be specified, both of which are mentioned in the environment labeled in Figure 6.

Finally, we introduce a new MARL benchmark called "Spaceworld" which simulates a cooperative In-Space Servicing, Assembly, and Manufacturing (ISAM) scenario [32]. In this grid-world environment illustrated in Figure 5(c), the agents are satellite robots tasked with servicing the James Webb Space Telescope by replacing damaged mirrors. The blue and red agents are tasked with bringing the correct mirror segment to their target, as illustrated via coloring. However, both agents are allowed to interact and move the mirror segments wherever they like. The primary motivation behind this benchmark was to include safety consequences and unforced collaboration. If the two agents collide (i.e. try to enter the same square) at any point, both agents are destroyed and the episode is

terminated with reward zero. This indirectly incentivizes one agent to stay still and avoid collision while the other agent moves both objects to their required locations. However, this produces suboptimal results as one point is deducted from their global reward for every additional timestep the agents take over the minimal. This reward signal is automatically normalized across grid sizes such that the highest achievable reward is always 10 and the lowest is 0. As such, the centralized critic must balance the safety and efficiency objectives to gain the highest possible reward.

We believe that these benchmarks effectively showcase our methods' capabilities for enabling intricate interactions in sparse-reward environments, handling partial observability, and learning safety-critical behaviors.

4.1 Baselines

To highlight the performance of both of our methods, LLM-MCA and LLM-TACA, we compare them against a number of other Centralized-Training Decentralized-Execution approaches.

First, we compare with LICA (Learning Implicit Credit Assignment for Cooperative Multi-Agent Reinforcement Learning) [61], a multi-agent actor-critic method representing the centralized critic with a hyper-network. This approach is motivated by finding locally optimal cooperative behaviors by directly maximizing a joint action value function.

Our second baseline is MAPPO (Multi-Agent PPO), an extension of PPO to multi-agent scenarios with centralized value function inputs [58]. This multi-agent actor-critic method learns the joint state value function and operates in cooperative, common-reward settings.

Finally, QMIX [38] is an extension of value decomposition networks (VDN) [47] to more complex centralized action-value functions. It uses a mixing network to estimate joint action values as a nonlinear combination of individual agent values.

4.2 Results Analysis

We compiled the results of our method against our 3 baselines in Figures 5 and 6. LICA completely fails to solve the matrix game, while all other methods achieve high scores dominated by our LLM-MCA and LLM-TACA methods. Already with the matrix benchmark,

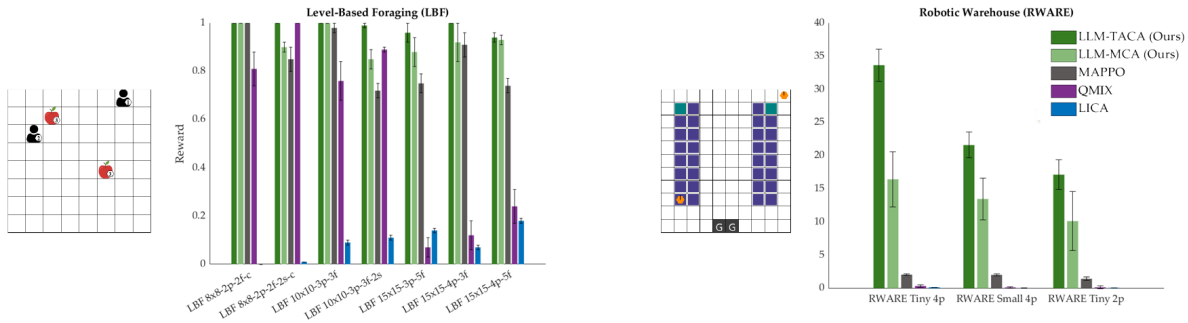


Figure 6: Comparison of our methods with baselines on common MARL benchmarks. All values reported are averaged across five separate trainings and the 95% confidence interval is illustrated by the error bars. We compared our methods on 12 scenarios across two common MARL benchmarks. Level-Based Foraging (LBF) is a common benchmark where agents must coordinate to pick apples based on their levels and the level of the apple. Robotic Warehouse (RWARE) is a partially-observable benchmark where robots must move certain packages to the delivery zone. LBF scenario names follow the format: "{GRID SIZE}-{PLAYER COUNT}p-{FOOD COUNT}f-{SIGHT COUNT}c IF COOPERATIVE MODE". RWARE environment names follow the format "{GRID SIZE} {PLAYER COUNT}p".

we see the benefit of having a critic which easily incorporates prior environment knowledge and retrospective analysis.

On the Spaceworld environment with a 10×10 grid both our methods consistently achieve the highest score whereas the three baselines have significantly lower scores and higher variance over 5 runs as they either succeed or completely fail. Due to any collisions immediately ending an episode, agents must learn to avoid such scenarios consistently and in a sample efficient manner. Additionally, as starting, object, and goal locations are randomized between episodes, our baselines often faced issues with nonstationary policy changes causing poor performance in even previously seen scenarios.

On the Level-Based Foraging environment both our methods consistently score in the 90th percentile with low variance with an advantage for LLM-TACA. MAPPO and QMIX achieve high scores on the smallest grid sizes but their performance drops, drastically for QMIX, as the grid size increases. On the other hand, LICA performs consistently poorly on all sizes.

In the Robotic Warehouse, the results exhibit a similar trend with both our approaches achieving scores an order of magnitude higher than that of MAPPO, which scores an order of magnitude higher than QMIX and LICA.

On the whole, as the tasks become more complex, the impact and efficiency of our method become more evident. Our method was able to harness prior knowledge about the scenarios to better assess agent performance and provide dense feedback quickly, causing our learned policies to be far less myopic than those produced by the baselines.

4.3 Offline Dataset

As a byproduct of our LLM critic method, we have gathered a large number of trajectories in each of these scenarios along with the individualized credit and task assignment for each agent during each timestep of the training process. This novel offline dataset will enable decentralized offline-learning future efforts². While there

have been efforts to learn multi-agent policies from offline data, they often must rely on an underlying centralized structure [21, 54], and so we provide this dataset with the hope that further progress can be made in this direction.

5 CONCLUSIONS AND FUTURE WORK

In this work, we proposed two LLM centralized critics for solving the multi-agent credit assignment problem. Specifically, we reformulated this credit assignment problem as a general sequence improvement problem. Then, building on humans' innate skills at solving general pattern recognition problems, we noticed a similar skill emerging in the latest LLMs. We leveraged this skill by introducing two LLM-based centralized critics tasked with assigning credit to each agent based on their respective participation toward their shared objective.

We evaluated our approach on major MARL benchmarks where we surpassed the state-of-the-art by a significant margin. Thanks to the interpretability and explainability of LLMs, the users of our method can easily understand both the emergent behaviors observed by the LLM critic among the agents and the feedback the LLM is proposing to improve their capabilities.

A limiting factor to this work is the LLM's slowness in generating outputs. We have already worked on mitigating this issue via batch processing of trajectories to minimize the number of LLM queries. Another limitation is the monetary cost of using closed-source LLMs like ChatGPT.

We envision several avenues for future work. We first plan to expand our approach to non-cooperative tasks as they can also suffer from the credit assignment problem. We want to utilize faster LLMs to enable real-time evaluations. This advance would then facilitate the deployment of our method on real robots.

ACKNOWLEDGMENTS

This material is based upon work supported by the National Science Foundation under grants DGE-2146752, CNS-2423130, CCF-2423131, and ECCS-2145134 CAREER Award.

²Dataset instructions at <http://iconlab.negarmehr.com/LLM-MCA/>

REFERENCES

- [1] Adrian K Agogino and Kagan Tumer. 2004. Unifying temporal and structural credit assignment problems. In *Autonomous Agents and Multi-Agent Systems Conference*. 980–987.
- [2] Maulik Bhatt, Yixuan Jia, and Negar Mehr. 2023. Efficient constrained multi-agent trajectory optimization using dynamic potential games. In *2023 IEEE/RSJ International Conference on Intelligent Robots and Systems (IROS)*. IEEE, 7303–7310.
- [3] Tom Brown, Benjamin Mann, Nick Ryder, Melanie Subbiah, Jared Kaplan, Prafulla Dhariwal, Arvind Neelakantan, Pranav Shyam, Girish Sastry, Amanda Askell, Sandhini Agarwal, Ariel Herbert-Voss, Gretchen Krueger, Tom Henighan, Rewon Child, Aditya Ramesh, Daniel M. Ziegler, Jeffrey Wu, Clemens Winter, Christopher Hesse, Mark Chen, Eric Sigler, Mateusz Litwin, Scott Gray, Benjamin Chess, Jack Clark, Christopher Berner, Sam McCandlish, Alec Radford, Ilya Sutskever, and Dario Amodei. 2020. Language models are few-shot learners. arXiv preprint arXiv:2005.14165.
- [4] Thomas Carta, Pierre-Yves Oudeyer, Olivier Sigaud, and Sylvain Lamprier. 2022. Eager: Asking and answering questions for automatic reward shaping in language-guided RL. *Advances in Neural Information Processing Systems* 35 (2022), 12478–12490.
- [5] Huaben Chen, Wenkang Ji, Lufeng Xu, and Shiyu Zhao. 2023. Multi-Agent Consensus Seeking via Large Language Models. *ArXiv abs/2310.20151* (2023). <https://api.semanticscholar.org/CorpusID:264816702>
- [6] Yiqun Chen, Hangyu Mao, Jiaxin Mao, Shiguang Wu, Tianle Zhang, Bin Zhang, Wei Yang, and Hongxing Chang. 2024. PTDE: Personalized training with distilled execution for multi-agent reinforcement learning. In *33rd International Joint Conference on Artificial Intelligence*. International Joint Conferences on Artificial Intelligence Organization, 31–39.
- [7] Filippos Christianos, Lukas Schäfer, and Stefano V Albrecht. 2020. Shared experience actor-critic for multi-agent reinforcement learning. In *34th Conference on Neural Information Processing Systems*. 10707–10717.
- [8] Caroline Claus and Craig Boutilier. 1998. The dynamics of reinforcement learning in cooperative multiagent systems. In *Proceedings of the Fifteenth National/Tenth Conference on Artificial Intelligence/Innovative Applications of Artificial Intelligence (Madison, Wisconsin, USA) (AAAI '98/IAAI '98)*. American Association for Artificial Intelligence, USA, 746–752.
- [9] Caroline Claus and Craig Boutilier. 1998. The dynamics of reinforcement learning in cooperative multiagent systems. *AAAI/IAAI 1998* (1998), 746–752.
- [10] Google DeepMind. 2024. AI achieves silver-medal standard solving International Mathematical Olympiad problems. <https://deepmind.google/discover/blog/ai-solves-imo-problems-at-silver-medal-level>
- [11] Lei Feng, Yuxuan Xie, Bing Liu, and Shuyan Wang. 2022. Multi-level credit assignment for cooperative multi-agent reinforcement learning. *Applied Sciences* 12, 14 (2022), 6938.
- [12] Jakob Foerster, Ioannis Alexandros Assael, Nando De Freitas, and Shimon Whiteson. 2016. Learning to communicate with deep multi-agent reinforcement learning. In *30th International Conference on Neural Information Processing Systems*, Vol. 29. 2145–2153.
- [13] Jakob Foerster, Gregory Farquhar, Triantafyllos Afouras, Nantas Nardelli, and Shimon Whiteson. 2018. Counterfactual multi-agent policy gradients. In *AAAI Conference on Artificial Intelligence*, Vol. 32. 2974–2982.
- [14] Maite Gimenez, Javier Palanca, and Vicent Botti. 2020. Semantic-based padding in convolutional neural networks for improving the performance in natural language processing. A case of study in sentiment analysis. *Neurocomputing* 378 (2020), 315–323.
- [15] Prasoon Goyal, Scott Niekum, and Raymond J Mooney. 2019. Using natural language for reward shaping in reinforcement learning. In *28th International Joint Conference on Artificial Intelligence*. 2385–2391.
- [16] Olle Häggström. 2023. Are large language models intelligent? Are humans?. In *Computer Sciences & Mathematics Forum*, Vol. 8. MDPI, 68.
- [17] Dongge Han, Chris Xiaoxuan Lu, Tomasz Michalak, and Michael Wooldridge. 2022. Multiagent model-based credit assignment for continuous control. In *21st International Conference on Autonomous Agents and Multiagent Systems*. 571–579.
- [18] Zhiwei He, Tian Liang, Wenxiang Jiao, Zhuosheng Zhang, Yujiu Yang, Rui Wang, Zhaopeng Tu, Shuming Shi, and Xing Wang. 2024. Exploring human-like translation strategy with large language models. *Transactions of the Association for Computational Linguistics* 12 (2024), 229–246.
- [19] Sirui Hong, Mingchen Zhuge, Jonathan Chen, Xiawu Zheng, Yuheng Cheng, Jinlin Wang, Ceyao Zhang, Zili Wang, Steven Ka Shing Yau, Zijuan Lin, et al. 2023. MetaGPT: Meta Programming for A Multi-Agent Collaborative Framework. In *The Twelfth International Conference on Learning Representations*.
- [20] Wenlong Huang, Pieter Abbeel, Deepak Pathak, and Igor Mordatch. 2022. Language models as zero-shot planners: Extracting actionable knowledge for embodied agents. In *International Conference on Machine Learning*. PMLR, 9118–9147.
- [21] Jiechuan Jiang and Zongqing Lu. 2023. Offline Decentralized Multi-Agent Reinforcement Learning. In *ECAI*. 1148–1155.
- [22] Shyam Sundar Kannan, Vishnunandan LN Venkatesh, and Byung-Cheol Min. 2023. Smart-LLM: Smart multi-agent robot task planning using large language models. arXiv preprint arXiv:2309.10062.
- [23] Daniel Martin Katz, Michael James Bommarito, Shang Gao, and Pablo Arredondo. 2024. GPT-4 passes the Bar Exam. *Philosophical Transactions of the Royal Society* 382, 2270 (2024), 20230254.
- [24] Jacob Devlin Ming-Wei Chang Kenton and Lee Kristina Toutanova. 2019. BERT: Pre-training of deep bidirectional transformers for language understanding. In *Annual Conference of the North American Chapter of the Association for Computational Linguistics*, Vol. 1. 2.
- [25] Yeseung Kim, Dohyun Kim, Jieun Choi, Jisang Park, Nayoung Oh, and Daehyung Park. 2024. A survey on integration of large language models with intelligent robots. *Intelligent Service Robotics* 17 (2024), 1091–1107.
- [26] Brenden M Lake, Tomer D Ullman, Joshua B Tenenbaum, and Samuel J Gershman. 2017. Building machines that learn and think like people. *Behavioral and brain sciences* 40 (2017), e253.
- [27] Zhiyuan Li, Wenshuai Zhao, Lijun Wu, and Joni Pajarinen. 2024. AgentMixer: Multi-Agent Correlated Policy Factorization. <https://openreview.net/forum?id=tzD9HVgeVx>
- [28] Ryan Lowe, Yi I Wu, Aviv Tamar, Jean Harb, OpenAI Pieter Abbeel, and Igor Mordatch. 2017. Multi-agent actor-critic for mixed cooperative-competitive environments. *Advances in neural information processing systems* 30 (2017).
- [29] Anuj Mahajan, Tabish Rashid, Mikayel Samvelyan, and Shimon Whiteson. 2019. MAVEN: Multi-agent variational exploration. In *33rd International Conference on Neural Information Processing Systems*, Vol. 32. 7613–7624.
- [30] Zhao Mandi, Shreeya Jain, and Shuran Song. 2024. RoCo: Dialectic multi-robot collaboration with large language models. In *International Conference on Robotics and Automation*. IEEE, 286–299.
- [31] Suvir Mirchandani, Fei Xia, Pete Florence, Brian Ichter, Danny Driess, Montserrat Gonzalez Arenas, Kanishka Rao, Dorsa Sadigh, and Andy Zeng. 2023. Large language models as general pattern machines. In *7th Annual Conference on Robot Learning*. PMLR, 2498–2518.
- [32] Kartik Nagpal and Negar Mehr. 2024. Optimal robotic assembly sequence planning: A sequential decision-making approach. In *IEEE/RSJ International Conference on Intelligent Robots and Systems*. IEEE.
- [33] Duc Thien Nguyen, Akshat Kumar, and Hoang Chuin Lau. 2018. Credit assignment for collective multiagent RL with global rewards. In *32nd International Conference on Neural Information Processing Systems*, Vol. 31. 8113–8124.
- [34] Georgios Papoudakis, Filippos Christianos, Lukas Schäfer, and Stefano V. Albrecht. 2021. Benchmarking multi-agent deep reinforcement learning algorithms in cooperative tasks. In *35th Conference on Neural Information Processing Systems Track on Datasets and Benchmarks*.
- [35] Ray Perrault and Jack Clark. 2024. *Artificial Intelligence Index Report*. Technical Report. Stanford University Human-Centered Artificial Intelligence.
- [36] Eduardo Pignatelli, Johan Ferret, Davide Paglieri, Samuel Coward, Tim Rocktäschel, Edward Grefenstette, and Laura Toni. 2024. Assessing the zero-shot capabilities of LLMs for action evaluation in RL. In *Automated Reinforcement Learning: Exploring Meta-Learning, AutoML, and LLMs*. ICML.
- [37] Aida Rahmattalabi, Jen Jen Chung, Mitchell Colby, and Kagan Tumer. 2016. D++: Structural credit assignment in tightly coupled multiagent domains. In *IEEE/RSJ International Conference on Intelligent Robots and Systems*. IEEE, 4424–4429.
- [38] Tabish Rashid, Mikayel Samvelyan, Christian Schroeder De Witt, Gregory Farquhar, Jakob Foerster, and Shimon Whiteson. 2020. Monotonic value function factorisation for deep multi-agent reinforcement learning. *Journal of Machine Learning Research* 21, 178 (2020), 1–51.
- [39] Claude Sammut. 2010. *Credit Assignment*. Springer US, 238–242.
- [40] Minah Seo, Luiz Felipe Vecchietti, Sangkeum Lee, and Dongsoo Har. 2019. Rewards prediction-based credit assignment for reinforcement learning with sparse binary rewards. *IEEE Access* 7 (2019), 118776–118791.
- [41] Arambam James Singh. 2021. *Credit assignment in multiagent reinforcement learning for large agent population*. Ph.D. Dissertation. Singapore Management University.
- [42] Oliver Slumbers, David Henry Mguni, Kun Shao, and Jun Wang. 2024. Leveraging Large Language Models for Optimised Coordination in Textual Multi-Agent Reinforcement Learning. <https://openreview.net/forum?id=1PPjf4wife>
- [43] Kyunghwan Son, Daewoo Kim, Wan Ju Kang, David Earl Hostallero, and Yung Yi. 2019. QTRAN: Learning to factorize with transformation for cooperative multi-agent reinforcement learning. In *International Conference on Machine Learning*. PMLR, 5887–5896.
- [44] Nitish Srivastava, Geoffrey Hinton, Alex Krizhevsky, Ilya Sutskever, and Ruslan Salakhutdinov. 2014. Dropout: a simple way to prevent neural networks from overfitting. *The Journal of Machine Learning Research* 15, 1 (2014), 1929–1958.
- [45] Alexandra Stolyarova. 2018. Solving the credit assignment problem with the prefrontal cortex. *Frontiers in Neuroscience* 12 (2018), 182.
- [46] Chuanneng Sun, Songjun Huang, and Dario Pompili. 2024. LLM-based multi-agent reinforcement learning: Current and future directions. arXiv preprint arXiv:2405.11106.

- [47] Peter Sunehag, Guy Lever, Audrunas Gruslys, Wojciech Marian Czarnecki, Viničius Zambaldi, Max Jaderberg, Marc Lanctot, Nicolas Sonnerat, Joel Z Leibo, Karl Tuyls, et al. 2018. Value-decomposition networks for cooperative multi-agent learning based on team reward. In *17th International Conference on Autonomous Agents and MultiAgent Systems*. 2085–2087.
- [48] Richard Stuart Sutton. 1984. *Temporal credit assignment in reinforcement learning*. University of Massachusetts Amherst.
- [49] Akira Taniguchi, Shuya Ito, and Tadahiro Taniguchi. 2024. Hierarchical path planning from speech instructions with spatial concept-based topometric semantic mapping. *Frontiers in Robotics and AI* 11 (2024), 10.
- [50] Gemma Team, Thomas Mesnard, Cassidy Hardin, Robert Dadashi, Surya Bhupatiraju, Shreya Pathak, Laurent Sifre, Morgane Rivière, Mihir Sanjay Kale, Juliette Love, et al. 2024. Gemma: Open models based on Gemini research and technology. arXiv preprint arXiv:2403.08295.
- [51] Hado Van Hasselt, Arthur Guez, and David Silver. 2016. Deep reinforcement learning with double Q-learning. In *AAAI Conference on Artificial Intelligence*, Vol. 30. 2094–2100.
- [52] Ashish Vaswani, Noam Shazeer, Niki Parmar, Jakob Uszkoreit, Llion Jones, Aidan N Gomez, Lukasz Kaiser, and Illia Polosukhin. 2017. Attention is all you need. In *31st Conference on Neural Information Processing Systems*. 6000–6010.
- [53] Jianhong Wang, Yuan Zhang, Tae-Kyun Kim, and Yunjie Gu. 2020. Shapley Q-value: A local reward approach to solve global reward games. In *AAAI Conference on Artificial Intelligence*, Vol. 34. 7285–7292.
- [54] Xiangsen Wang, Haoran Xu, Yinan Zheng, and Xianyuan Zhan. 2024. Offline multi-agent reinforcement learning with implicit global-to-local value regularization. *Advances in Neural Information Processing Systems* 36 (2024).
- [55] Gerhard Weiß. 1995. Distributed reinforcement learning. In *The Biology and Technology of Intelligent Autonomous Agents*. Springer, 415–428.
- [56] Leixin Xu, Weibin Chen, Xiang Liu, and Yang-Yang Chen. 2021. MADDPG: Multi-agent deep deterministic policy gradient algorithm for formation elliptical encirclement and collision avoidance. In *5th Chinese Conference on Swarm Intelligence and Cooperative Control*. Springer, 239–250.
- [57] Yaodong Yang, Ying Wen, Jun Wang, Liheng Chen, Kun Shao, David Mguni, and Weinan Zhang. 2020. Multi-agent determinantal Q-learning. In *International Conference on Machine Learning*. PMLR, 10757–10766.
- [58] Chao Yu, Akash Velu, Eugene Vinytsky, Jiaxuan Gao, Yu Wang, Alexandre Bayen, and Yi Wu. 2022. The surprising effectiveness of PPO in cooperative multi-agent games. *Advances in Neural Information Processing Systems* 35 (2022), 24611–24624.
- [59] Jingqing Zhang, Yao Zhao, Mohammad Saleh, and Peter Liu. 2020. PEGASUS: Pre-training with extracted gap-sentences for abstractive summarization. In *International Conference on Machine Learning*. PMLR, 11328–11339.
- [60] Kaiqing Zhang, Zhuoran Yang, and Tamer Başar. 2021. *Multi-agent reinforcement learning: A selective overview of theories and algorithms*. Springer, Chapter 12, 321–384.
- [61] Meng Zhou, Ziyu Liu, Pengwei Sui, Yixuan Li, and Yuk Ying Chung. 2020. Learning implicit credit assignment for cooperative multi-agent reinforcement learning. *Advances in Neural Information Processing Systems* 33 (2020), 11853–11864.
- [62] Yuan Zhuang, Yi Shen, Zhili Zhang, Yuxiao Chen, and Fei Miao. 2024. YOLO-MARL: You only LLM once for multi-agent reinforcement learning. arXiv preprint arXiv:2410.03997.

6 APPENDIX

6.1 Implementation Details

Note that when using the Level-Based Foraging (LBF) environment³ [7, 34] the environment name lets you customize many elements including the number of agents, the observability, the grid size, cooperativity, and the maximum number of food elements. These elements are defined in the environment name in 1 with format: "{GRID SIZE}-{PLAYER COUNT}p-{FOOD COUNT}f-{SIGHT COUNT}{-c IF COOPERATIVE MODE}".

Additionally, for our Robotic Warehouse [7, 34] scenarios, this environment is available with 3 different overall grid sizes and the number of robots can be specified, both of which are mentioned in the environment name with the format "{GRID SIZE} {PLAYER COUNT}p".

We summarize our main results in Table 1.

³Note that the Level-Based Foraging environment has since changed the "-c" tag to "-coop" in version 3.

Table 1: Comparison of our method with previous solutions on common benchmarks. All values reported are the means and 95% confidence interval for the collective return of the agents by each algorithm over 5 separate training runs with no common reward or reward scalarization techniques. The bold score denotes the highest value of a row and an asterisk * denotes the second highest.

Environments\Algorithms	LLM-MCA (Ours)	LLM-TACA (Ours)	LICA [61]	MAPPO [58]	QMIX [38]
Climbing Matrix Game	195.0 ± 33.5*	210.0 ± 10.0	169.0 ± 13.4	175.0 ± 0.0	175.0 ± 54.8
Spaceworld-10x10	10.0 ± 0.0	10.0 ± 0.0	2.0 ± 2.2	6.0 ± 2.5*	4.4 ± 4.6
Level Based Foraging					
8x8-2p-2f-c	1.00 ± 0.00	1.00 ± 0.00	0.00 ± 0.00	1.00 ± 0.00	0.81 ± 0.07*
8x8-2p-2f-2s-c	0.90 ± 0.02*	1.00 ± 0.00	0.01 ± 0.00	0.85 ± 0.05	1.00 ± 0.00
10x10-3p-3f	1.00 ± 0.00	1.00 ± 0.00	0.09 ± 0.01	0.98 ± 0.02*	0.76 ± 0.08
10x10-3p-3f-2s	0.85 ± 0.04	0.99 ± 0.01	0.11 ± 0.01	0.72 ± 0.03	0.89 ± 0.01*
15x15-3p-5f	0.88 ± 0.06*	0.96 ± 0.04	0.14 ± 0.01	0.75 ± 0.04	0.07 ± 0.04
15x15-4p-3f	0.92 ± 0.08*	1.00 ± 0.00	0.07 ± 0.01	0.91 ± 0.05	0.12 ± 0.06
15x15-4p-5f	0.93 ± 0.02*	0.94 ± 0.02	0.18 ± 0.01	0.74 ± 0.03	0.24 ± 0.07
Average	0.93 ± 0.03*	0.98 ± 0.01	0.09 ± 0.01	0.85 ± 0.07	0.54 ± 0.05
Robotic Warehouse					
Tiny 4p	16.44 ± 4.16*	33.65 ± 2.44	0.14 ± 0.02	2.05 ± 0.12	0.33 ± 0.18
Small 4p	13.47 ± 3.15*	21.59 ± 1.95	0.03 ± 0.01	2.01 ± 0.18	0.12 ± 0.11
Tiny 2p	10.12 ± 4.43*	17.13 ± 2.27	0.06 ± 0.02	1.43 ± 0.26	0.14 ± 0.19
Average	13.34 ± 3.91*	24.12 ± 2.22	0.08 ± 0.02	1.83 ± 0.19	0.20 ± 0.16