# MARLIN: Multi-Agent Reinforcement Learning Guided by Language-Based Inter-Robot Negotiation

Toby Godfrey[1], William Hunt[1] and Mohammad D. Soorati[1]

*Abstract*— **Multi-agent reinforcement learning is a key method for training multi-robot systems over a series of episodes in which robots are rewarded or punished according to their performance; only once the system is trained to a suitable standard is it deployed in the real world. If the system is not trained enough, the task will likely not be completed and could pose a risk to the surrounding environment. We introduce Multi-Agent Reinforcement Learning guided by Language-based Inter-Robot Negotiation (MARLIN), in which the training process requires fewer training episodes to reach peak performance. Robots are equipped with large language models that negotiate and debate a task, producing plans used to guide the policy during training. The approach dynamically switches between using reinforcement learning and large language model-based action negotiation throughout training. This reduces the number of training episodes required, compared to standard multi-agent reinforcement learning, and hence allows the system to be deployed to physical hardware earlier. The performance of this approach is evaluated against multi-agent reinforcement learning, showing that our hybrid method achieves comparable results with significantly reduced training time.**

## I. INTRODUCTION

Mutli-robot systems have emerged as a promising paradigm for tackling complex tasks in various domains. Advanced decision-making algorithms are essential to coordinate and control these systems effectively. Multi-Agent Reinforcement Learning (MARL) is a powerful framework for training agents to collaborate, or compete [1], and has become an integral technique for multi-robot systems' development and deployment. However, the complexity of real-world scenarios often requires sophisticated learning and reasoning capabilities, which can be difficult and time-consuming to train. Recently, pre-trained Large Language Models (LLMs) have demonstrated few-shot performance in a range of tasks [2], [3]. This knowledge distillation and reasoning are useful in enhancing the decision-making capabilities of multi-robot systems.

By taking actions in an environment and learning through punishment and reward, agents learn policies through experience and exploration [1]. Proximal Policy Optimisation (PPO) is a type of reinforcement learning policy that has become popular in robotics applications [4]. Due to the clipped surrogate objective function, the size of policy updates is limited, resulting in more stable training. PPO was extended

to multi-agent learning in [5] by adding a centralised critic to aid in training. Multi-Agent PPO (MAPPO) maintains the clipped surrogate objective function, and hence the stability, but also follows the centralised training, decentralised execution pattern, which has shown promise in overcoming the nonstationarity problem commonly found in MARL [6].

When training multi-robot systems using MARL, significant time and data are often required to train models before they can be deployed onto robots. This paper aims to address the amount of training required to deploy a MARL policy to robots in the real world. This problem can cause delays in the deployment of novel multi-robot systems due to training requirements, or because robots become damaged from prematurely deploying a policy. This work could have wide-reaching impacts by reducing the number of episodes it takes for a policy to perform well enough such that it can reliably be applied to hardware without adding significant risk to the surrounding environment or the integrity of robots.

LLMs are trained on vast corpora of text and other internet-scale data, and so contain prior knowledge and reasoning skills unseen in other approaches. Since their inception, the community has seen an increasing number of LLM-integrated robotic systems, alongside many attempts at mitigating and overcoming common pitfalls.

This work harnesses the reasoning skills and flexibility of LLMs to improve the efficiency at which multi-robot systems are trained using MARL. **M**ulti-**A**gent **R**einforcement **L**earning guided by Language-Based, **I**nter-Agent **N**egotiation (MARLIN) is proposed; a method where agents negotiate how to solve a given task using natural language plans, which are then used to help train the MARL policy. This balances quick, robust MARL with prior-knowledge-based, but slower inter-agent negotiation. To save training time, off-the-shelf language models without any additional fine-tuning are used. This approach leverages encoded reasoning and natural language understanding to generate more performant plans. As a result, policies can be deployed to physical robots after fewer training episodes, compared to a standard MARL implementation.

The primary contribution of this work is the development of a novel hybrid approach which augments MARL with dialogical LLMs that utilise inter-agent negotiation to generate natural-language plans which are then used to improve the in-training performance of cooperative robots, with no loss in accuracy compared to conventional MARL approaches. This approach is evaluated against MAPPO and an LLM baseline in a series of cooperative tasks where robots must navigate a narrow corridor and swap positions, both in simulation and

[1]School of Electronics & Computer Science, University of Southampton, Southampton, SO17 1BJ, United Kingdom
{t.godfrey, w.hunt, m.soorati}@soton.ac.uk

on physical robots.

## II. RELATED WORK

Multi-robot systems have become more common in the landscape due to their increased robustness and flexibility when compared to single robot systems [7]. MARL is used extensively for multi-robot systems to learn a desired behaviour. A more recent development is LLMs, which are continually being improved upon [8], [9], and have recently been integrated into robotics [10]. These models have been applied to multi-robot systems to build plans in natural language [11]. When language models are extended with action tokens and visual processing capabilities, they have been used to develop end-to-end robot control with improved generality and semantic reasoning [12]. Similarly, PaLM-E was developed as an embodied multi-modal language model capable of mapping information from the textual-visual domain to real-world actions [13]. When LLMs are capable of interacting with other LLMs, they have shown the ability to mimic complex human behaviours, such as reflection on past experiences and forming opinions [14]. Interacting LLMs can also increase the accuracy of results through multi-agent debate [15].

The intersection of LLMs and MARL is less developed than that of language models and robotic agents. Sun *et al.* detailed the current landscape of LLM-based MARL and discussed the potential for language models to be used in the training stage of MARL [16]. LLMs may be able to increase the sample efficiency and explainability of MARL [17]. In [18], an LLM was combined with MARL to generate safe MARL policies. Natural language constraints were translated into semantic embeddings for use in the MARL training process. Their framework performed comparably to conventional MARL models but with fewer constraint violations. LLMs have also been used to interpret natural language navigation instructions and pass them to a reinforcement learning or MARL model, resulting in a policy that can act on human instructions [19], [20].

Zhang *et al.* used an LLM to generate an intrinsic reward function to help alleviate the credit assignment problem in MARL environments with sparse rewards [21]. Their integration of models showed improved sample efficiency over a standard Deep Q-Learning implementation. Tarakli *et al.* used a language model to solicit human feedback used to guide MARL training in order to improve performance [22]. Liu *et al.* developed the eSpark system to integrate LLMs (one as a generator and one as a critic) into MARL exploration [23]. The eSpark system directly generates exploration functions using one of the language models and then masks out potential actions. These functions are then applied and the best policy is selected. Their method increased performance, compared to a standard MARL implementation, however, the eSpark system is only applicable to homogeneous agents and dense environments.

Our system differs from the aforementioned works by using inter-agent negotiation to produce plans to guide MARL training. This is less centralised and utilises the prior knowledge encoded into the language models to solve the problem in fewer episodes. Actions are generated through inter-agent negotiation and debate instead of using the approach of eSpark to centrally generate functions for all agents. As our negotiation-based action generation system makes no assumptions regarding robot morphologies, MARL rewards, or the underlying MARL architecture, MARLIN can be used in a wider array of systems.

## III. METHOD

MAPPO is used as the base MARL model for several reasons. It achieves stable training due to the clipping of the surrogate objective function [5]. It also follows the centralised training, decentralised execution paradigm which has been shown to achieve state-of-the-art performance [24]. MAPPO is commonly used in multi-robot systems, allowing MARLIN to be applicable to a wider range of systems. As the agents are cooperative, parameters are shared between agents to decrease the number of trainable parameters. MAPPO is extended with an additional method for action generation to increase performance throughout training. Our system is comprised of two action-generation methods: (1) sampling the MARL action distribution (standard MAPPO), and (2) natural language inter-agent negotiation (LLM planner). At each inference step one of these methods is used, and the generated actions control the robots during training. Situations where LLMs are likely to generate superior plans compared to the action distribution are capitalised on through runtime switching between action-generation functions. Irrespective of how actions are generated, the MARL model benefits by training on the trajectory. Both the standard MAPPO algorithm and our novel approach were implemented using the Ray RLlib library [25]. A key facet of our work is that no fine-tuning is performed on language models. Fine-tuning may yield performance gains, but MARLIN was designed for scenarios where robots are in use during training, resulting in limited opportunities for LLM fine-tuning before training.

The system can be modelled as a partially observable Markov decision process $(\mathbf{S}, \mathbf{A}, \mathbf{O}, T, \Omega, \mathscr{R})$ where $S$ denotes the state space, $\mathbf{A} = \{W, F, B, L, R\}$ is the set of actions available to agents, corresponding to waiting ($W$), moving forwards ($F$), backwards ($B$), left ($L$) or right ($R$). The observation space $\mathbf{O} \subseteq \mathbb{Z}^4$ for each agent is constructed such that $o \in \mathbf{O}$ is a 4-tuple $(x, y, x_g, y_g)$, representing the agent's current position and its goal position. The transition function $T : S \times A \times S \to [0, 1]$ specifies the probability of transitioning from state $s$ to state $s'$ under action $a$ and the observation function $\Omega : S \times A \times O \to [0, 1]$ returns the probability of observing $o \in \mathbf{O}$ for a given agent in state $s$ taking action $a$. The reward function is $\mathscr{R} : \mathbf{S} \times \mathbf{A_i} \to \mathbb{R}$ such that $\mathscr{R}(i) = w \cdot \max(0, \frac{d_i}{D_i}) - \rho$ where $w \in \mathbb{N}$ is a constant to weight the performance over penalties, $d_i$ is the Manhattan distance from agent $i$'s current position to its goal position, $D_i$ is the Manhattan distance from agent $i$'s starting position to its goal position and $\rho$ represents penalties for collisions with other agents or the boundaries.
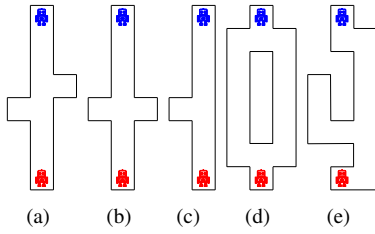
Fig. 1. Diagrams of the scenarios used for evaluation; (a) Asymmetrical Two Slot Corridor, (b) Symmetrical Two Slot Corridor, (c) Single Slot Corridor, (d) Two Path Corridor, (e) Maze Like Corridor

The environment consists of narrow corridors in which two agents cannot move past each other, and a way for the agents to either let the other pass or move around each other. Diagrams of the environments can be seen in Fig. 1. The Symmetrical Two Slot Corridor (Fig. 1(b)) is based on the environment used in [26]. This idea was expanded to other corridors for the rest of the environments. The goal of all environments is for the agents to swap positions from one end of the corridor to the other. The task was kept simple to aid in reproducibility.

A feedforward Neural Network (NN) is trained to generate actions for each robot. Each robot has a NN with a $\mathbb{Z}^4$ input vector, two fully connected hidden layers, each with 256 units and tanh activation functions, and a $\mathbb{R}^5$ output vector representing the action logits.

MAPPO follows the centralised training, decentralised execution paradigm, and so during training there is a centralised value function. For a system with two robots, the feedforward NN to predict the value function has a $\mathbb{Z}^{13}$ input vector (own observation, opponent's observation, and opponent's one-hot encoded action), one fully connected hidden layer with 16 units and tanh activation functions, and an $\mathbb{R}^1$ output.

### A. Generation Methods

A generator function, $G$, is used, such that $G : \mathbf{S} \to \mathbf{A}^n$ where $\mathbf{S}$ is the set of states, $\mathbf{A}$ is the set of all actions and $n$ is the total number of agents. These generator functions generate the actions for all agents for every step. Our approach uses two distinct generators, Action Distribution Sampling (based on RL), and Inter-Agent Negotiation (using LLMs). *Action Distribution Sampling:* Sampling from the policy's action distribution is the standard way to generate actions for MAPPO. The selected actions depend on the output of the model's network and are based on expected reward, gradually converging towards the best actions for a given state [27]. *Inter-Agent Negotiation:* While generating actions from the model's action distribution works well once the model has adequate experience, it can take significant time and data to reach such a standard. This leaves the model selecting poor-performing actions, especially at the beginning of training or, in the case of a dynamic environment, after the environment changes. By integrating each agent with an LLM instance, they can discuss the next actions to take. The inter-agent negotiation system is extended from [11] in which agents take turns discussing and critiquing solutions to a task.
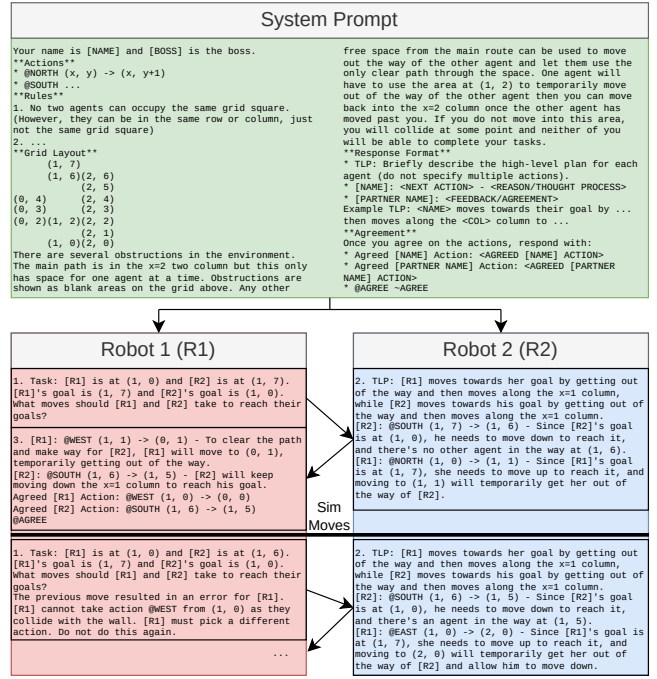


Fig. 2. A diagram of the inter-agent negotiation mechanism. Both agents share the same system prompt, then alternate in suggesting a Top-Level Plan (TLP) and moves for both agents. Critiques can be provided to improve moves until they are agreed by the pair. Moves are then simulated and the next moves are discussed.

### B. Inter-Robot Negotiation

To generate actions using language models, the robots discuss the next step to take. Due to the nature of LLMs, these discussions are conducted in natural language and hence are understandable by human observers. The language models act as planners instead of controllers. A lower layer is used to translate the high-level actions, e.g. @NORTH, into linear and angular velocities for the robot wheels. This also allows MARLIN to be applied to other discrete planning tasks by simply writing a translator from high-level to low-level actions, and redefining the LLMs' actions. In this work, the LLMs were hosted remotely and accessed over the Internet. At the beginning of a negotiation round, one of the robots is chosen to be the leader, this can either be deterministic (used in our experiments) or a random choice from all robots. This helps to mitigate agents yielding the decision to another agent, which assistant-style LLMs are prone to doing [28]. Then the robots alternate to discuss (in pairs) the best actions for both agents until they either come to a consensus or reach the message limit. Once an action has been agreed, the movements are stored and simulated, and the robots are given their updated positions. This process repeats until either all robots reach their goals, or they reach the maximum number of moves. Plans are cached for use in later training, or until they are replaced with higher-performing ones. In the system prompt, the robots are given the format in which responses should be written. This format can then be parsed to extract the agreed moves. If the models produce an output in a different format, they are prompted
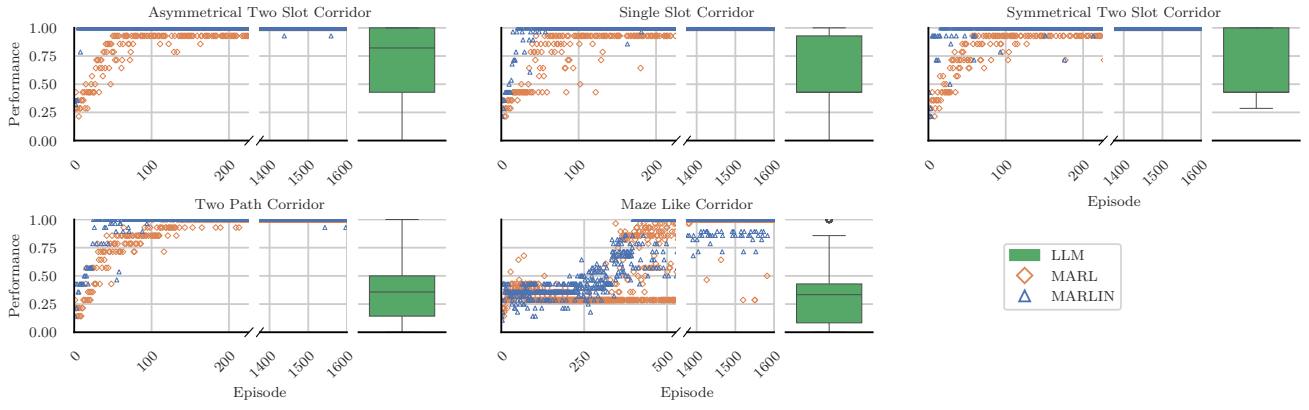
Fig. 3. Median performance of the MARLIN and MARL systems for different scenarios in simulation. The boxplot shows the distribution of performance scores for trials using the LLM-based negotiation.

to correct it. Operator preferences could also be included in the system prompt to further influence robot behaviour. A diagram of the negotiation mechanism can be seen in Fig. 2.

### C. Action Selection Process

---

**Algorithm 1** MARLIN Algorithm

---

**Require:** N is the set of agents
1: **while** episode $\leq$ episode$_{\text{max}}$ **do**
2:     $\overline{p} \leftarrow$ mean(pastPerformanceBuffer)
3:     **while** step $<$ step$_{\text{max}}$ **do**
4:         $P \leftarrow$ loadPlan(state)
5:         $p_{\text{plan}} \leftarrow$ perf($P$)
6:         **if** step $= 0$ **then**
7:             **if** episode $< m$ **then**
8:                 $G \leftarrow G_{\text{ADS}}$
9:             **else if** episode $< 2m$ **then**
10:                 $G \leftarrow G_{\text{IAN}}$
11:             **else**
12:                 $G \leftarrow \begin{cases} G_{\text{IAN}}, & \text{if } p_{\text{plan}} = 1 \\ G_{\text{IAN}}, & \text{if } \overline{p} < \overline{p_{\text{LLM}}} \\ G_{\text{ADS}}, & \text{otherwise} \end{cases}$
13:         **if** $P = \varnothing$ **and** $G = G_{\text{IAN}}$ **then**
14:             $P \leftarrow$ makePlan()
15:         **if** step $= \frac{\text{step}_{\text{max}}}{2}$ **and** rand() $\leq 0.1$ **then**
16:             $G \leftarrow$ toggleGenerator($G$)
17:         actions $\leftarrow$ genActionSingleStep($G$, step)
18:         trainStep(actions)      ▷ Update env. & model
19:         $p \leftarrow \frac{\sum_{i=1}^{|\mathbf{N}|} \max\left(0, 1 - \frac{d_i}{D_i}\right)}{|\mathbf{N}|}$    ▷ Evaluate performance
20:         pastPerformanceBuffer $\leftarrow p$

---

The algorithm for training using MARLIN is shown in Algorithm 1. An initialisation phase of $2m$ episodes establishes the initial performance of both generation methods. At the start of all proceeding episodes, it is dynamically determined which action generator, $G$, is to be used. $G_{\text{IAN}}$ and $G_{\text{ADS}}$ indicate the Inter-Agent Negotiation-based generator and the Action Distribution Sampling-based generator respectively;

$p_{\text{plan}}$ is the performance of the plan for the current state, $\overline{p}$ is the average performance of the past 5 episodes and $\overline{p_{\text{LLM}}}$ is the average performance for the LLM-only system in that state. For the performance metric, $d_i$ is the Manhattan distance from agent $i$'s end position to its goal position, and $D_i$ is the Manhattan distance from agent $i$'s starting point to its goal position; thus a successful trial has a performance of 1. This is similar but differs from the reward function of the MARL model, which includes penalties for collisions.

If $G_{\text{IAN}}$ is selected, the best plan for the current state is retrieved from memory. If a plan for the current state does not yet exist, or the retrieved plan has worse performance than the current rolling average, a new plan is created to ensure poor plans do not disrupt the training later in the process. If the regenerated plan has a higher performance than the cached version, the cache is updated with the new plan. To further enhance the effectiveness of training, a probabilistic element was implemented whereby there is a 10% chance that the generation technique of any given episode will be switched halfway through the episode. This helps to gather experience from elsewhere in the environment.

### IV. EVALUATION

The effectiveness of our approach is evaluated by testing each component and our combined hybrid system in each scenario. Episode length was limited to 50 moves for all tests and all common configurations were kept the same between the three systems. The Llama 3.1 8B Instruct model [9] was chosen as the LLM for all evaluation experiments, and the performance metric from Algorithm 1 was used to evaluate each trial. Models were trained for 1600 episodes.

*Simulation-Based Evaluation:* The first phase of experimentation was to test the performance of MARLIN against the MAPPO and LLM benchmarks in simulation. Fig. 3 shows that for all the scenarios, MARLIN outperforms both the benchmark MARL implementation and the system only using the inter-robot negotiation. Furthermore, in most scenarios the plans generated by the off-the-shelf LLMs still have difficulty navigating around the other robot. However, robot movement beyond the 50% performance mark

greatly increased the likelihood of task completion. When comparisons are made to the MAPPO benchmark, it is seen that peak performance is reached by MARLIN after fewer training episodes for most scenarios. This shows that plans generated through inter-robot discussion are effective in providing beneficial actions to help train the MARL policy. For all but the Maze-Like Corridor (see Fig. 1(e)), our method reached peak median performance earlier than the MARL system. When trying to solve the Maze-Like Corridor, MARL reached peak performance from episode 1000, but with some slight variation in the results. MARLIN plateaued slightly below the perfect score at around 0.95 but reached this score at episode 750, earlier than MARL.

TABLE I
AVERAGE PERFORMANCE THROUGHOUT TRAINING. **BOLD** INDICATES
SIGNIFICANCE ($p < 0.001$ OVER A 50 RUN PAIRED T-TEST)

| Scenario | 100 | | 850 | | 1600 | |
|---|---|---|---|---|---|---|
| | MARL | MARLIN | MARL | MARLIN | MARL | MARLIN |
| Asymmetrical Two Slot Corr. | 0.8543 | **1.0000** | 0.9814 | 0.9957 | 0.9843 | 0.9929 |
| Maze Like Corridor | 0.4300 | 0.4271 | 0.8529 | **0.9643** | 0.8214 | **1.0000** |
| Single Slot Corridor | 0.7157 | **0.9071** | 0.9329 | **0.9943** | 0.9700 | **1.0000** |
| Symmetrical Two Slot Corr. | 0.8671 | **0.9714** | 0.9929 | **1.0000** | 0.9857 | 0.9952 |
| Two Path Corr. | 0.7771 | **0.9271** | 0.9929 | **1.0000** | **1.0000** | 0.9667 |

These data show equal or higher performance by our approach, as compared to a standard MAPPO model, and an increased sample efficiency. These data highlight how, most of the time, our hybrid approach delivers better performance than using LLMs alone. Table I shows that MARLIN outperforms MARL in most cases. Only for the Maze Like Corridor (see Fig. 1(e)) did the MARL have better performance from the outset, but the performance of MARLIN was significantly higher later in the training period. This slower start may be due to the more complex environmental reasoning required by the LLM. As progress is made, the LLMs gain a better understanding of how to move around the environment, and the performance of their plans increases.

*Physical Robot Experiments:* To validate our findings from the simulations and to close the reality gap, further experiments were conducted on real hardware. For this, two TurtleBot3 robotic platforms were used, running ROS2 Humble Hawksbill. As the Maze-Like Corridor was the hardest scenario to solve in the simulation-based experiments, it was focused on for the physical testing (see Fig. 4(a))[1]. The performance of the system for both MARLIN and the MARL system was evaluated by testing the policy every 250 episodes, starting from episode 100. Fig. 5 shows the performance of our approach and the benchmark on physical hardware. This highlights that real-world performance follows the same trend as that seen in Fig. 3, with MARLIN reaching near-perfect performance earlier than the MARL

[1]Experimental videos can be found at https://sooratilab.com/z/marlin.php.



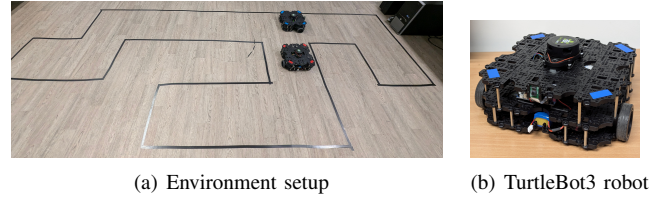(a) Environment setup      (b) TurtleBot3 robot

Fig. 4. The environment and robot platform used for the physical robot experiments in the Maze-Like Corridor. Results are reported in Fig. 5.
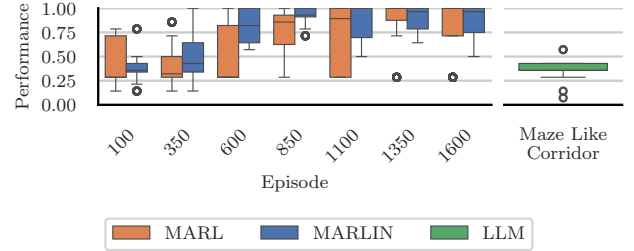


Fig. 5. Median performance of the system for the Maze-Like Corridor when executed on physical hardware. The boxplot shows the performance distribution when the LLM-only system was evaluated on physical hardware.

system, but with the MARL approach slightly outperforming our system after a considerable amount of training. These data reinforce our previous results that LLM-based, inter-agent negotiation can provide valuable information during MARL training to reduce the amount of training required.

*Scalability Experiments:* In order to ensure our inter-agent negotiation system scales to larger environments with more agents, we developed the environment shown in Fig. 6. MARLIN's negotiation system is designed for pairs of robots to discuss a task, therefore for larger systems, such a negotiation mechanism needs to be triggered only when robots must negotiate. For this proof-of-concept, one agent starts at each grey tile and moves randomly until either: (1) the agent reaches the green square and leaves the environment, or (2) the agent finds another agent in the same area and they cannot move past each other. If two agents are in
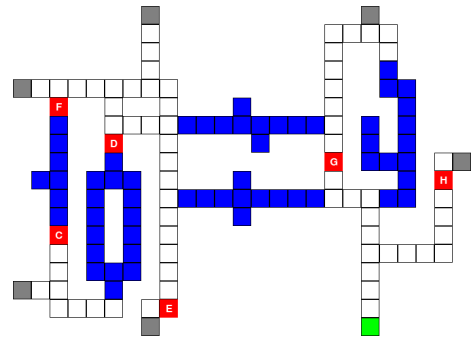


Fig. 6. A more complex scenario with a larger number of robots is used to test our method at a large scale. Agents negotiate in pairs to pass corridors with no alternative route (e.g., blue areas highlighting our examined corridors). The green cell is the exit and the grey cells indicate starting positions.

one of these areas, moves are generated using our inter-agent negotiation system. This highlights how our dyadic negotiation approach can be deployed to systems with more agents by detecting when negotiation pairs must be formed. This LLM-based negotiation system also enables its use in environments which may change during runtime. However, efficiently scaling LLM-based approaches is still an open challenge in multi-robot systems.

## V. Conclusions

MARLIN augments MARL training with inter-agent, natural language negotiation. Throughout MARL training, agents negotiate using the reasoning skills of off-the-shelf LLMs to generate actions which guide the MARL policy. This approach was evaluated in both simulation and the real-world and showed that, on average, MARLIN reaches peak performance in fewer training episodes than the MAPPO baseline. Therefore, systems trained using MARLIN can be deployed to real hardware after fewer episodes than typical MARL. This highlights how LLMs can utilise debate to act as a discrete planner for multi-robot systems. Whilst heuristic-based approaches could be useful for the tasks outlined in this work, the benefits of LLMs include more transparency and incorporation of human preferences, and adaptations to dynamic and more complex environments.

However, our system has a few limitations which could be addressed by future work. Firstly, only discrete actions and environments are used in the MARL model. Preliminary testing showed a decrease in LLM reasoning ability for continuous environments, but future work can alter our method to maintain performance in continuous tasks. This work used the Llama 3.1 8B Instruct model due to its availability and low-cost hosting, but extensions to this work could explore the deployment of local language models. Distilled small-scale models could run on the robots themselves to remove the requirement for an internet connection [29]. Another limitation of our system is the lack of environmental information being passed to the LLMs. It is possible that an extension in which robot sensor information is passed to the LLMs could provide better performance. Future work could also explore the performance implications of fine-tuning the models using human-written expert plans. Overall, our approach increased the speed of training with no loss in performance by utilising LLMs' inherent reasoning abilities.

## References

[1] S. V. Albrecht, F. Christianos, and L. Schäfer, *Multi-Agent Reinforcement Learning: Foundations and Modern Approaches*. MIT Press, 2024.

[2] Y. Chen, J. Arkin, Y. Zhang, N. Roy, and C. Fan, "Scalable multi-robot collaboration with large language models: Centralized or decentralized systems?" in *2024 IEEE International Conference on Robotics and Automation (ICRA)*. IEEE, 2024, pp. 4311–4317.

[3] Z. Mandi, S. Jain, and S. Song, "RoCo: Dialectic multi-robot collaboration with large language models," in *2024 IEEE International Conference on Robotics and Automation (ICRA)*, 2024, pp. 286–299.

[4] J. Schulman, F. Wolski, P. Dhariwal, A. Radford, and O. Klimov, "Proximal Policy Optimization Algorithms," Aug. 2017.

[5] C. Yu, A. Velu, E. Vinitsky, J. Gao, Y. Wang, A. Bayen, and YI. WU, "The surprising effectiveness of PPO in cooperative multi-agent games," in *Advances in Neural Information Processing Systems*, vol. 35, 2022, pp. 24 611–24 624.

[6] C. Schroeder de Witt, J. Foerster, G. Farquhar, P. Torr, W. Boehmer, and S. Whiteson, "Multi-Agent Common Knowledge Reinforcement Learning," in *Advances in Neural Information Processing Systems*, vol. 32. Curran Associates, Inc., 2019.

[7] M. Dorigo, G. Theraulaz, and V. Trianni, "Swarm Robotics: Past, Present, and Future [Point of View]," *Proceedings of the IEEE*, vol. 109, no. 7, pp. 1152–1165, July 2021.

[8] T. B. Brown, *et al.*, "Language Models are Few-Shot Learners," July 2020.

[9] A. Dubey, *et al.*, "The Llama 3 Herd of Models," Aug. 2024.

[10] W. Hunt, S. D. Ramchurn, and M. D. Soorati, "A Survey of Language-Based Communication in Robotics," June 2024.

[11] W. Hunt, T. Godfrey, and M. D. Soorati, "Conversational language models for human-in-the-loop multi-robot coordination," in *Proceedings of the 23rd International Conference on Autonomous Agents and Multiagent Systems*, 2024.

[12] A. Brohan, *et al.*, "RT-2: Vision-Language-Action Models Transfer Web Knowledge to Robotic Control," July 2023.

[13] D. Driess, *et al.*, "PaLM-E: An embodied multimodal language model," in *Proceedings of the 40th International Conference on Machine Learning*, ser. ICML'23. Honolulu, Hawaii, USA: JMLR.org, 2023.

[14] J. S. Park, J. O'Brien, C. J. Cai, M. R. Morris, P. Liang, and M. S. Bernstein, "Generative Agents: Interactive Simulacra of Human Behavior," in *Proceedings of the 36th Annual ACM Symposium on User Interface Software and Technology*, Oct. 2023.

[15] Y. Du, S. Li, A. Torralba, J. B. Tenenbaum, and I. Mordatch, "Improving Factuality and Reasoning in Language Models through Multiagent Debate," May 2023.

[16] C. Sun, S. Huang, and D. Pompili, "LLM-based Multi-Agent Reinforcement Learning: Current and Future Directions," May 2024.

[17] Z. Zhang, "Advancing sample efficiency and explainability in multi-agent reinforcement learning," in *Proceedings of the 23rd International Conference on Autonomous Agents and Multiagent Systems*, 2024, pp. 2791–2793.

[18] Z. Wang, M. Fang, T. Tomilin, F. Fang, and Y. Du, "Safe Multi-agent Reinforcement Learning with Natural Language Constraints," May 2024.

[19] S. Morad, A. Shankar, J. Blumenkamp, and A. Prorok, "Language-Conditioned Offline RL for Multi-Robot Navigation," July 2024.

[20] W. Zu, W. Song, R. Chen, Z. Guo, F. Sun, Z. Tian, W. Pan, and J. Wang, "Language and sketching: An LLM-driven interactive multimodal multitask robot navigation framework," in *2024 IEEE International Conference on Robotics and Automation (ICRA)*. IEEE, 2024, pp. 1019–1025.

[21] A. Zhang, A. Parashar, and D. Saha, "A Simple Framework for Intrinsic Reward-Shaping for RL using LLM Feedback," 2023.

[22] I. Tarakli, S. Vinanzi, and A. D. Nuovo, "Interactive Reinforcement Learning from Natural Language Feedback," in *2024 IEEE/RSJ International Conference on Intelligent Robots and Systems (IROS)*, Oct. 2024, pp. 11 478–11 484.

[23] Z. Liu, X. Yang, Z. Liu, Y. Xia, W. Jiang, Y. Zhang, L. Li, G. Fan, L. Song, and B. Jiang, "Knowing What Not to Do: Leverage Language Model Insights for Action Space Pruning in Multi-agent Reinforcement Learning," May 2024.

[24] Y. Zhou, S. Liu, Y. Qing, K. Chen, T. Zheng, Y. Huang, J. Song, and M. Song, "Is Centralized Training with Decentralized Execution Framework Centralized Enough for MARL?" May 2023.

[25] The Ray Team, "RLlib: Industry-Grade Reinforcement Learning — Ray 2.35.0," https://docs.ray.io/en/latest/rllib/index.html, 2024.

[26] M. Bettini, A. Shankar, and A. Prorok, "Heterogeneous multi-robot reinforcement learning," in *Proceedings of the 2023 International Conference on Autonomous Agents and Multiagent Systems*, 2023, pp. 1485–1494.

[27] R. S. Sutton and A. G. Barto, *Reinforcement Learning: An Introduction*. MIT press, 2018.

[28] G. Li, H. A. A. K. Hammoud, H. Itani, D. Khizbullin, and B. Ghanem, "CAMEL: Communicative Agents for "Mind" Exploration of Large Language Model Society," Nov. 2023.

[29] J. Li, *et al.*, "Large Language Model Inference Acceleration: A Comprehensive Hardware Perspective," Jan. 2025.