# LLM A*: Human in the Loop Large Language Models Enabled A* Search for Robotics
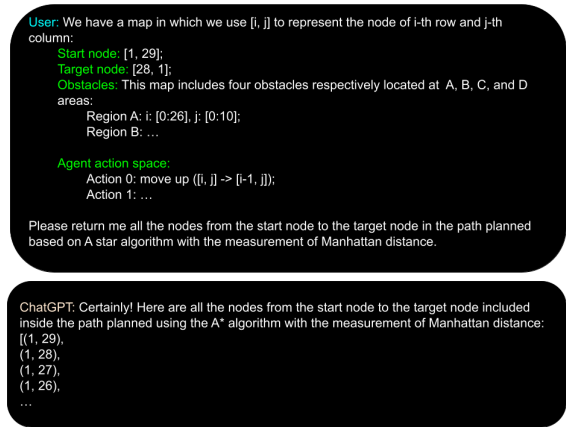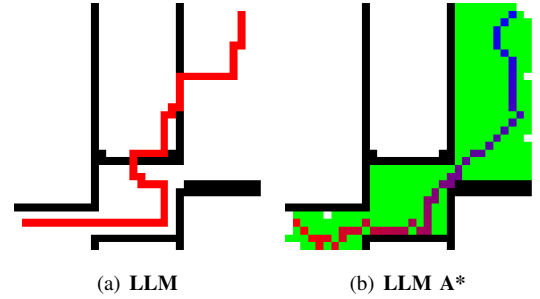
Hengjia Xiao[1], Peng Wang[2], Mingzhe Yu[3], Mattia Robbiani[2]

*Abstract*— This research explores how Large Language Models (LLMs) can assist in path planning for mobile embodied agents, such as robots, in an interactive, human-in-the-loop manner. We introduce LLM A*, a novel framework that leverages the commonsense reasoning of LLMs alongside the utility-optimal A* algorithm to enable few-shot, near-optimal path planning. By integrating LLMs within a reinforcement learning-inspired structure, LLM A* treats each interaction with the model as an alignment operation, akin to a reward function, allowing iterative refinement of the agent's policy at each step of the A* algorithm. Prompts serve two main functions: providing LLMs with essential information, such as environment details, costs, and heuristics, and relaying human feedback on intermediate planning results. This approach enhances transparency while enabling code-free path planning, fostering accessibility for users with limited coding expertise. Comparative analysis against A* and reinforcement learning (RL) shows that LLM A* achieves paths comparable to A* while significantly reducing search space and outperforming RL. Its interactive nature also makes it a promising tool for collaborative human-robot tasks. Code and supplemental materials are available on GitHub: https://github.com/speedhawk/LLM-A-.

## I. INTRODUCTION

Path planning, alongside mapping, localisation, and motion planning, has long been regarded as the cornerstone technique for enabling fully autonomous agents such as robots [1]–[3]. It typically involves taking the environment map, the initial and goal states as inputs, and aims to find the (near) optimal path in terms of action costs to move the agent from the initial state to the goal state. Various approaches to path planning have been proposed, which can be classified into 1) search-based algorithms such as the classical A* search which is guaranteed to find an optimal path when heuristics are admissible and consistent; 2) sampling-based planning such as the Rapidly-exploring Random Tree (RRT) [4] and Ant Colony Optimisation (ACO) [5], etc.; and 3) Data-driven planning such as Reinforcement Learning (RL) based search [6].

With the widespread application of neural networks, some hybrid algorithms such as neural A* [7] that combine traditional search algorithms with the multi-layer network cognitive model (MLP) have attracted attention from academia and industry. Although these works have greatly improved the applicability of traditional search algorithms in path planning tasks in complex environments, some of their significant

(a) **LLM**          (b) **LLM A***



(c) **Conversation between human and LLM**

Fig. 1: LLM-based path planning: (a) path planned by LLM directly; (b) path planned by the proposed LLM A*. The initial and goal states are at the upper right and the lower left corners, respectively. The white and black tiles represent free spaces and obstacles. The red tiles form the final paths and the green tiles are the total searched tiles. We can see the path planned by LLM goes through obstacles (results generated by GPT3.5-turbo), which is prohibitive in robotics.

shortcomings still face challenges and further research. For example,

Recent advancements in LLMs have facilitated a potential overhaul of data-driven techniques. LLMs exhibit capabilities akin to human-like text generation and task completion [8]. Their generated texts consistently demonstrate adept handling of commonsense information, often to a human-like level. Particularly noteworthy is their ability to engage in a chain-of-prompts interaction, wherein humans can guide LLMs to provide desired answers. This process resembles RL apart from that the rewards are assessed and granted by humans through interaction with LLMs. These developments have somewhat diverted attention from questioning the ex-

arXiv:2312.01797v3 [cs.RO] 15 May 2025

plainability of LLMs' underlying mechanisms, such as the 'transformer', a type of deep neural network. Instead, the focus has shifted towards the models' proficiency in promptly responding to queries posed to them.

This paper harnesses the advantages of A* being able to find an optimal path and LLMs being able to consider commonsense knowledge when interacting with humans (agents) and proposes the LLM A*, a human-in-the-loop solution for robot path planning.

Briefly, LLM A* can be summarised into three pivotal stages. The first stage is to set up the environment where a robot performs path planning, and communicates the information to LLMs. The second stage entails the definition of the initial state, the goal state, and the heuristics to be used. We have adopted a generic form of heuristic that considers obstacles, distances, and actions the robot can perform in the environment. It is worth mentioning that the heuristic can be tailored to specific setups based on the applications. Subsequently, LLMs identify feasible moves for path planning, evaluated against a cost function combining path cost and current-state heuristic, akin to A*. The best 'move' with the minimum cost will be selected and communicated to humans upon request, who may accept or decline it. The decision will be communicated back to LLMs to carry on completing or terminating the planning. One benefit of human guidance and assessments is to minimise the use of tokens, which remains a challenge to use LLMs such as ChatGPT [8].

LLM A* adheres to a pattern of '{state, action, rewards}' that is similar to RL, yet without necessitating explicit coding, and humans have full control of the planning process. To this end, we have compared it to the advanced RL algorithm PPO for performance evaluation. A quantitative comparison with A* is also conducted. The contributions of this paper include: 1). a first-of-its-kind LLM-based A* is proposed for robotic path planning and a set of metrics are defined to evaluate the performance; 2). the interactive and code-free nature of LLM A* makes it appealing to non-expert robotic users; 3). enhanced safety assurance in comparison to A* and RL as humans are involved in the planning and loop and have full control of the planning process.

The remainder of the paper is organised as follows. Some related works are introduced in Section II. Section III elaborates on the approach. Experiments, discussions, and an ablation study are provided in Section IV, and the paper is concluded in Section V.

## II. RELATED WORK

### A. LLM for Robotic Task Planning

The concept of commanding robots through human language or instructions has been a longstanding aspiration among robotcists [9], [10]. Before the notable success of LLMs, Anderson et al. [9] introduced the R2R navigation framework, where Transformers were used to achieve visually grounded natural language navigation by translating human instructions into robot action sequences. Each action facilitated the robot's movement from one viewpoint to another. However, details regarding how the path was planned

between viewpoints were not provided, and navigation errors of up to 10 meters were reported.

Trained on (text) datasets of immense size, LLM models come with commonsense knowledge, which is promising to help accelerate robotic task planning [11]. To the best of our knowledge, most research from literature follows the visual-linguistic grounding strategy and proposes to harness the commonsense knowledge of LLMs to build a better correspondence between human language/instructions and visual perceptions of the robots. Song et al. [11] have demonstrated that LLMs can help with few-shot planning, i.e., the commonsense knowledge of LLMs is used to generate hierarchical plans, with the high-level plan being a set of viable while commonsense aligned subgoals generated by LLMs. The high-level plan is next passed to a low-level planner that is independent of the instructions to plan paths between subgoals. Other works such as [12] still fall into visual-linguistic grounding but a collision check is proposed.

It is worth noting that the works aforementioned do not delve into how LLMs can be employed to plan paths between subgoals or viewpoints. Aghzal et al. [13] introduce path planning from natural language (PPNL), emphasising the utilisation of LLMs in spatial-temporal path planning. However, their method primarily relies on prompts, lacking a connection or comparison to traditional path planning methods like A* and RL.

### B. Reinforcement Learning for Robotic Task Planning

Reinforcement learning presents another data-driven learning-based framework for path planning. Within the RL framework, an agent endeavors to learn a policy that selects a sequence of actions leading it efficiently from one state to another, ultimately reaching the goal. This learning process is guided by a policy that rewards favorable actions while penalizing unfavorable ones.

The capacity of deep learning in data representation and processing has garnered significant attention from both the RL and robotic communities. Consequently, deep learning models have been integrated into the RL framework, giving rise to deep reinforcement learning (DRL) [14], to enable more sophisticated and effective policy learning, for robotic task planning, etc. [15].

Nevertheless, both RL and DRL still face the challenges of 1) low data efficiency that necessitates a substantial number of interactions with the environment to learn policies, as well as the need for large datasets for training deep learning models; and 2) the stability of convergence, i.e., RL/DRL models are not guaranteed to converge. As a result, recent DRL models have focused on striking a balance between optimal states and convergence stability. The PPO model with actor-critic structures is one of the most prominent solutions. PPO is a policy gradient algorithm designed to maximise the horizon returns, while the actor-critic structure enables separate training of states and policy. This arrangement enhances planning performance and convergence of the PPO [16].

The prompt-driven LLM-based path planning can be viewed as a variant of RL models, wherein rewards are administered by humans via prompts, and policies are implicitly learned by LLMs through interactions. Consequently, we undertake a comparison of path planning performance between LLM-based methods and RL-based methods in this paper.

## III. THE LLM A* APPROACH

In most real-world scenarios, path planning is defined as the process where the goal location is known, but the environment encountered by the agent during pathfinding is uncertain. In real-world A* path planning, both distance and environmental factors must be considered. Unlike many prior studies that divide known paths into stages, we propose that pathfinding should be achieved through staged inference. In this approach, the agent dynamically selects a sub-goal, advances toward it, and continues this process until reaching the final destination.

On the other hand, the uncertainty in the environment may lead the agent to visit points that do not directly contribute to the current pathfinding task. However, these seemingly irrelevant points could hold value for future exploration tasks. Therefore, the agent should retain these "key points" for future use. After each stage of A* path planning, these points are communicated to the large language model (LLM) to assist in determining the next sub-goal.

In the following section, we present our hierarchical model, which consists of two levels: a lower level that handles A* sub-goal path planning and self-adaptive auxiliary tasks, and a higher level where the LLM determines the sub-goal based on the agent's exploration history. In the final section, we explain the interaction between these two levels.

### A. A* Preliminary

In a segmented map, the cost function of A* algorithm for planning the $s$-th state is defined as:

$$f(s) = \min_{\zeta(s) \sim C} g(s) + \min_{\zeta(s) \sim C} h(s) \tag{1}$$

where $g(s)$ and $h(s)$ are respectively representing the cost from the current node to the start and the goal point. For the purpose of planning, we select the minimum cost value of both $g(s)$ and $h(s)$ based on the cost policy function $\zeta(s)$, which varies as the state $s$ changes, and the cost space $C$. For the segmented map with static and known environment, obviously, the cost policy function is to calculate the distance values as the unique planning information. However, if the environment is not public to the agent at the beginning of planning, it is certain that the agent needs to gradually enrich and improve its perception of the environment in the exploration in the future. Here, we further decomposed the cost policy function into two items:

$$\zeta(s) = d(s) + l(s) \tag{2}$$

where $d(s)$ represents the distance value from the start and goal and $l(s)$ the relative environment value that includes the

information except for distance. If we set $l(s)$ is learnable by the parameter $\theta$, the original cost function should be modified as:

$$f(s) = \min_{\zeta_\theta(s) \sim C} g(s) + \min_{\zeta_\theta(s) \sim C} h(s) \tag{3}$$

where

$$\zeta_\theta(s) = d(s) + l_\theta(s) \tag{4}$$

### B. Self-Adaption Auxiliary Task

Enlighten by [3], our strategy to achieve the self-adaptive learning of environment value is to introduce a pixel control auxiliary task (PCAT), in which there will be a reward value that does not depend on the current state but the pixel discrepancy between two different states. Commonly, the rewards with larger values are caused by more radical changes of the pixel colour in order that the agent can 'learn more from the unprecedented experience. This reward value can be gradually learned under an semi-A2C learning structure, in which for each step, the environment value $l_\theta(s)$ of A* is equal to the fixed $V(s)$ calculated by the value network but the policy network does not need to train because that it is just the heuristic function of A* itself. Therefore, we have

$$l_\theta(s) = V_\theta(s) \tag{5}$$

and the value network training obeys the loss function:

$$J(\theta) = \min_{sae} \frac{1}{2}(V(s) - V_\theta(s))^2 \tag{6}$$

where

$$V(s) := V(s) + \alpha(r(s, f(s)) + \gamma V(s') - V(s)) \tag{7}$$

### C. LLM-Based Initial Reward

It is important to note that the initial random reward $r(s, f(s))$ in deep Q-learning has not yet been fully defined. To address this issue, we introduce a large language model (LLM) as a higher-level determiner, forming a hierarchical structure in conjunction with a self-adaptive A* algorithm. At the beginning of each sub-stage of pathfinding using A*, the process from the previous stage is comprehensively evaluated and the results are fed into the LLM. This generates a data structure that informs the initial reward distribution for the current pathfinding stage.

Moreover, during A* path planning, certain jumps inevitably occur at specific points to minimize the cost function $g(s)$. In our experiments, we recorded these jump occurrences as additional feedback for the LLM. These jumps indicate that sufficient exploration has occurred in the current direction, even though the path along the jump point may not be followed further. Nevertheless, these points remain potentially valuable for future path planning tasks. By remembering these points, the LLM can more effectively determine feasible goals in subsequent stages.

## IV. EXPERIMENTS AND DISCUSSIONS

### A. Setup

To evaluate the performance of LLM A*, we conducted a series of experiments comparing it to A* and RL. We considered two variants of LLM A*:

1) Greedy LLM A*, which only considers the heuristic $h(s)$ as the cost function. To avoid ambiguity, we denote it as LLM Greedy hereafter.
2) LLM A*, which considers the combination of cumulative and heuristic costs $f(s)$ as the cost function.

The GPT3.5-turbo-16k LLM is used as it provides more tokens than GPT3.5-turbo or even GPT4. This is to ensure that we can get the planning results without disruption, not necessarily that the algorithm needs so many tokens. The GPT-3.5-turbo-16k LLM model allows for a total of 16,384 tokens, whereas GPT-3.5-turbo has a maximum of 4,096 tokens. In case there is a need to reduce the consumption of tokens to replicate the experiments, one can reset the request dictionary to null after each interaction. To enable full control/access to the planning process, one can split a single interaction process into two stages: 1) planning by LLM A*; and 2) outputting necessary results upon requests.

For evaluation purposes, we primarily utilise three different occupancy grid maps (which are popular in planning and navigation) denoted as Ailse, Canyon, and Double Door, respectively. Each map is of the same size $24 \times 24$ but with different obstacle distributions. Each map consists of both free spaces and obstacles, wherein the robot agent can only traverse through free spaces while avoiding collisions with obstacles. The agent can move in eight directions at most, provided it is safe to do so. All experiments are conducted at least three times using Python 3.8+ on Google Colab, and the average results are reported. The RL model employed in our experiments is based on the PPO configuration. We use this RL model alongside the standard A* algorithm for comparison purposes. Experiments on grid maps sized $16 \times 16$ and $32 \times 32$ were also conducted with results and analysis provided in the Ablation section.

### B. Evaluation Metrics

To evaluate the performance of the proposed approaches and their counterparts, we utilise three metrics: **1)** The Search Complexity: For both A* and LLM A*, the overall length of the open set (grids to visit) and the closed set (visited grids) maintained by the two algorithms is used as the search complexity indicator. For LLM Greedy and RL, search complexity is equivalent to the total number of grids accessed until the goal is reached. This metric is used to explore the potential of LLM in reducing search complexity. **2)** The Path Length: Measures the length of the final path found by each algorithm. **3)** The Maximum Deviation Times (MDT): The MDT serves as a measure of the smoothness of the final path and it is defined in Equation (8)

$$\text{MDT} = \sum_{t=1}^{T} \delta\Big(\langle v_d, v_t \rangle\Big), \tag{8}$$

where

$$\delta\Big(\langle v_d, v_t \rangle\Big) = \left\{ \begin{array}{ll} 1, & \pi/2 < \langle v_d, v_t \rangle \leqslant \pi \\ 0, & \text{else} \end{array} \right. . \tag{9}$$

In the definitions above, $v_d$ is a directional vector that connects the initial state and the goal state and it indicates a global direction the search should expand towards. $v_t$ is a vector originating from state $s_{t-1}$ and terminating at state $s_t$ and we use it to represent the direction of the agent's movement from state $s_{t-1}$ to $s_t$. We first calculate the angle $\langle v_d, v_t \rangle$ between $v_d$ and $v_t$. If the result satisfies Equation (10), it suggests that the agent is likely to move away from the goal. This can result in a longer path or a path with numerous back-and-forth movements, which undermines the smoothness of the path and is not preferred by agents such as robots.

$$\pi/2 < \langle v_d, v_t \rangle \leqslant \pi, \tag{10}$$

where $\langle v_d, v_t \rangle$ is the inner product of $v_d$ and $v_t$, and we use it to denote the angle between $v_d$ and $v_t$. Essentially, MDT is the maximum number of times Equation (10) is satisfied when the agent moves toward the goal along the planned path.

### C. LLM A* Training and Session Design

For experiments involving LLM, there are two stages: initialisation and interactive planning. During initialisation, essential information about the environment and the agent is prompted to be set up for planning. This includes 1) locations (coordinates) of the initial and goal states; 2) how obstacles are distributed in the environment; 3) the action space of the agent; 4) the Manhattan distance used for cumulative and heuristic cost calculation; and 5) the objective to plan a path between the initial and goal states. In addition, planning rules are communicated to LLM, such as 1) a viable path should avoid colliding with obstacles; 2) the path should ideally expand along the direction from heuristics or human guidance; 3) preference of actions that help accelerate the planning process.

In the second interactive planning stage, LLM can provide planning results based on the initial information and other prompts to assist humans in guiding or monitoring the planning process. This stage operates iteratively and interactively until a path is successfully planned. It's important to note that humans can request intermediate planning results at any stage, making the planning process transparent (white box) to humans and ensuring safety, among other benefits.

### D. RL Model Training

The PPO model serves as a data-driven benchmark alongside A* for comparison. In our setup, the PPO model adopts an actor-critic structure comprising two 3-layer deep neural networks for policy and value training, respectively. The model undergoes training for 2,000 episodes, with each episode having a maximum of 200 steps to ensure convergence. To prevent premature model fixation and achieve a balance between exploration and exploitation, we set the
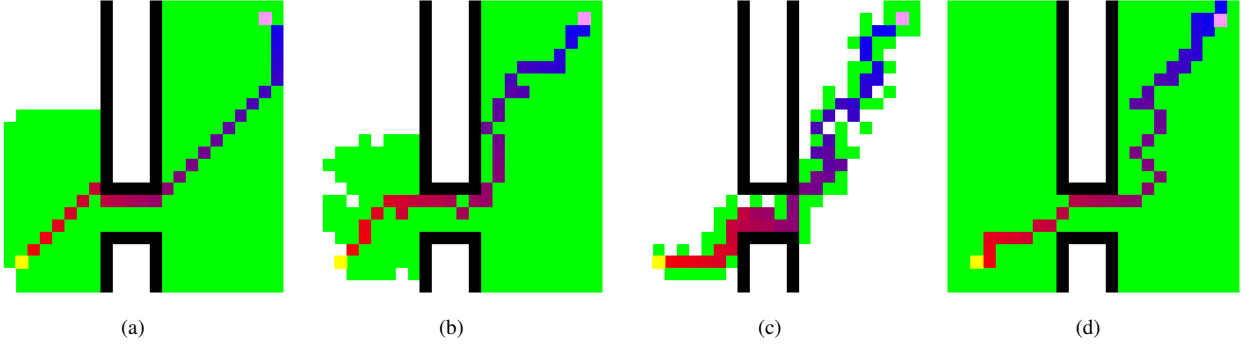
Fig. 2: Path planning results: (a), (b), (c), and (d) are results from A*, LLM A*, LLM Greedy, and PPO, respectively. The pink grids mark the initial states, while the yellow grids denote the goal states. White tiles depict free spaces, whereas black tiles represent obstacles. The green grids illustrate the search space. The final path is depicted by grids transitioning from blue to red, with color gradients aiding in visualising any back-and-forth movements within the paths.

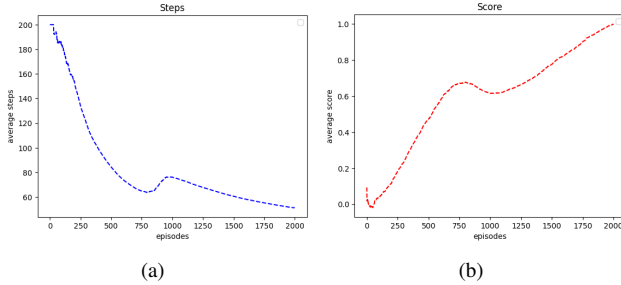learning rates of both the actor and critic to be smaller than 0.0005.



Fig. 3: Covergence of RL model training: (a) average steps per episode; (b) average scores achieved per episode.

To enhance the adaptability of the PPO model to diverse environments, we introduce a mechanism of randomising the initial state in each episode. This mechanism involves a progressive training approach, starting from states near the goal where the agent is close enough to reach it in one step. Subsequently, states are randomly selected from this subset for PPO training. As the model initially operates within a limited map scope, convergence is relatively easy. We then incrementally increase the map scope during subsequent training iterations until convergence is achieved for path planning from the initial state to the goal state. This process, known as the easy-to-difficult mechanism, enables the model to balance adaptability to varying environments while efficiently converging. Fig. 3 shows the average steps and scores against episodes for the PPO model training. We can see that the PPO model converges in 2,000 episodes and the results are eventually generated from this model.

### E. Main Results

Figure 2 illustrates the visual comparison of path planning results obtained from A*, LLM A*, LLM Greedy, and the PPO model. It is evident that the search complexity, indicated by the green tiles, of LLM A* is reduced compared to A* and

| Metrics | Environments | A* | LLM A* | LLM Greedy | PPO |
|---|---|---|---|---|---|
| Path Length | Aisle | **27** | **34** | 48.67 | 37.33 |
| | Canyon | **24** | **27.67** | 41.67 | 29 |
| | Double Door | **31** | **37.33** | 45.67 | 38 |
| MDT | Aisle | **1** | 4.67 | 12.33 | **3** |
| | Canyon | **0** | 2.33 | 7.67 | **0.33** |
| | Double Door | **0** | 3 | 7 | **2.33** |
| Complexity | Aisle | 372 | **352** | **100.67** | 471 |
| | Canyon | 145 | **133** | **74** | 156 |
| | Double Door | 210 | **197.67** | **103.67** | 214 |

TABLE I: Evaluation results of the algorithms against the defined metrics. Note each algorithm was executed three times and the average results are reported. The best two results are highlighted in the table.

PPO. However, it is notable that LLM Greedy outperforms LLM A*. This is attributed to the greedy nature of the heuristics, which consistently guides the agent towards the goal direction. However, the presence of obstacles may cause the agent to be redirected backwards, resulting in back-and-forth movements in the final path, making it less efficient. On the other hand, LLM A* considers both cumulative and heuristic costs, leading to an expanded search space but resulting in a smoother path with fewer back-and-forth movements.

Table I shows the quantitative results achieved by each algorithm against Path Length, MDT, and Search Complexity. Notably, LLM-based algorithms demonstrate superior performance compared to PPO in all environments and maintain an advantage over A* in most cases in terms of search complexity. In addition, LLM A* outperforms PPO and LLM Greedy in terms of path length, achieving comparable results with A*. However, A* exhibits superiority in terms of path smoothness, with LLM A* and PPO showing similar performance. Further studies have been carried out to further investigate and evaluate the performance in various environments, and the results are given in the Ablation section and the Supplemental Materials in our Github Repo: https://github.com/speedhawk/LLM-A-.
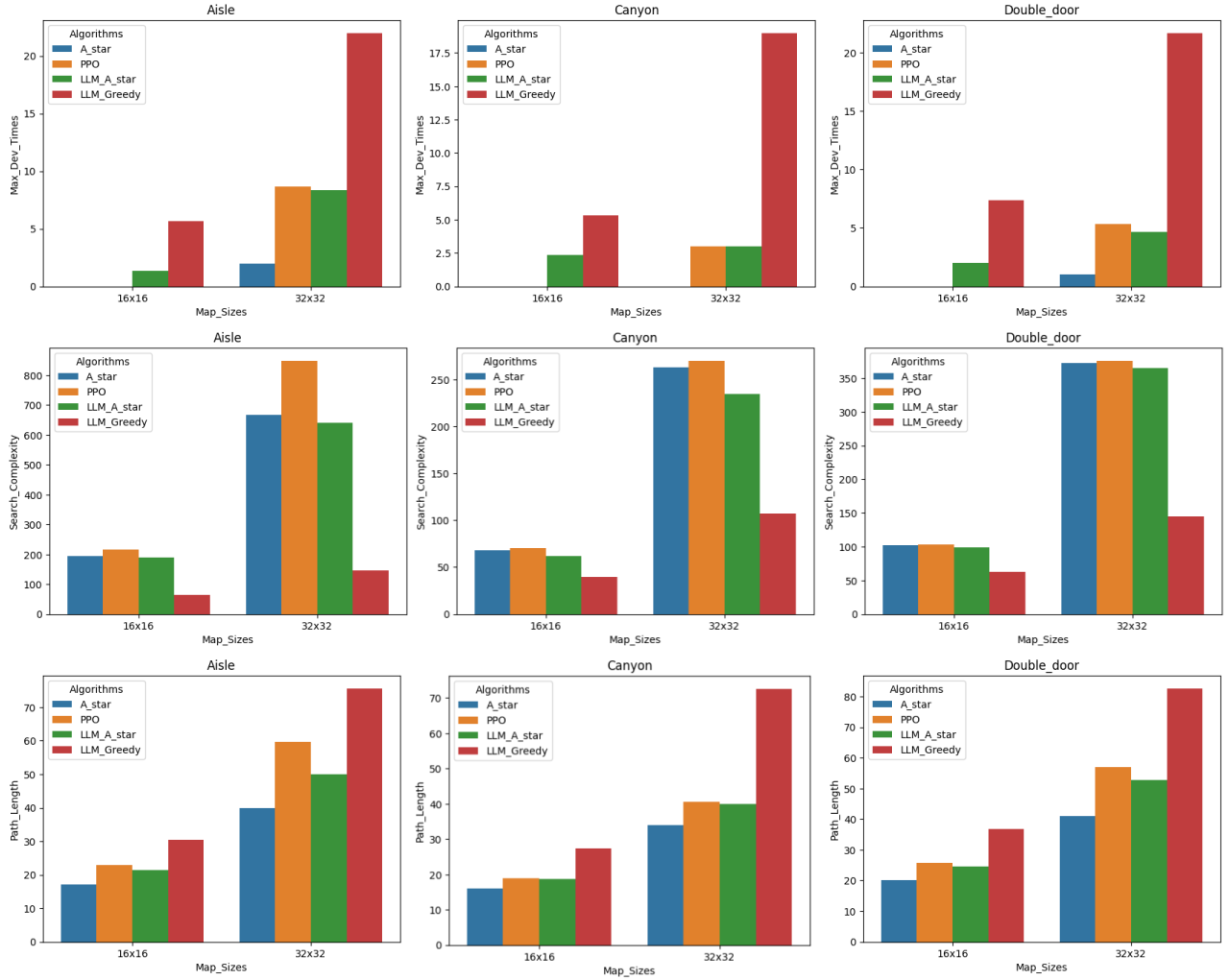
Fig. 4: Experimental results of A*, LLM A*, LLM Greedy, and PPO on the Aisle and Double Door environments with different sizes, i.e., $16 \times 16$ and $32 \times 32$.

*F. Ablation Study*

For a more comprehensive analysis of performance, additional experiments were conducted on grid maps sized $16 \times 16$ and $32 \times 32$. The results are depicted in Figure 4. It is evident that LLM A* generally outperforms both LLM Greedy and PPO across all metrics in each environment while achieving comparable results with A*.

Notably, there is an increasing trend in MDT for all algorithms as the size of the environment increases. However, we have observed a comparative advantage of LLM A* over PPO as the environment size increases, indicating that the agent following the path planned by LLM A* is less likely to experience back-and-forth movements within the final path compared to the path planned by RL. This underscores the potential of human knowledge and commonsense from LLMs to enhance path-planning outcomes in robotics.

*G. Discussion*

The integration of LLMs undoubtedly leads to increased interactions between humans and AI agents. Our research has shown that by incorporating LLMs into traditional path-planning algorithms like A*, we can achieve near-optimal path-planning results while still allowing humans to oversee the process. This represents a significant step towards making AI techniques transparent to humans and contributes to ensuring safety in human-robot interaction and collaboration.

It is important to highlight that the process of planning a path through interaction with LLMs bears a resemblance to using RL in path planning. Our findings demonstrate that leveraging commonsense knowledge from LLMs can greatly enhance path planning performance, particularly in terms of path length, path smoothness, and search complexity. This aspect will be further explored in our future work.

## V. CONCLUSIONS

This paper makes the first effort with the integration of Large Language Models (LLMs) with the classic A* algorithm, introducing LLM A* as a novel approach to human-in-the-loop interactive path planning for mobile embodied agents. By leveraging the inherent commonsense of

LLMs and the optimality of A*, LLM A* achieves few-shot near-optimal path planning compared to data-driven models such as PPO. Additionally, LLM A* offers the unique advantage of providing humans with complete access to the path-planning process, thereby enhancing safety when incorporating LLM A* in embodied agent path planning.

Despite the notable advantages of LLM A*, it is acknowledged that it shares similar inefficiencies (the LLM model training is data and computational resources demanding) with data-driven methods such as RL models. Future research will focus on enhancing the efficiency of LLM from both language modeling training/fine tuning and its integration with A*, aiming to strike a balance between interactive capacity and efficiency to facilitate the real-life deployment of such models.

## REFERENCES

[1] P. Wang, Q. Zhang, and Z. Chen, "A grey probability measure set based mobile robot position estimation algorithm," *International Journal of Control, Automation and Systems*, vol. 13, pp. 978–985, 2015.

[2] Q. Zhang, P. Wang, and Z. Chen, "An improved particle filter for mobile robot localization based on particle swarm optimization," *Expert Systems with Applications*, vol. 135, pp. 181–193, 2019.

[3] P. Wang, L. Mihaylova, P. Bonnifait, P. Xu, and J. Jiang, "Feature-refined box particle filtering for autonomous vehicle localisation with openstreetmap," *Engineering Applications of Artificial Intelligence*, vol. 105, p. 104445, 2021.

[4] D. González, J. Pérez, V. Milanés, and F. Nashashibi, "A review of motion planning techniques for automated vehicles," *IEEE Transactions on intelligent transportation systems*, vol. 17, no. 4, pp. 1135–1145, 2015.

[5] C. Carr and P. Wang, "Fast-spanning ant colony optimisation (fasaco) for mobile robot coverage path planning," *arXiv preprint arXiv:2205.15691*, 2022.

[6] C. Zhou, B. Huang, and P. Fränti, "A review of motion planning algorithms for intelligent robots," *Journal of Intelligent Manufacturing*, vol. 33, no. 2, pp. 387–424, 2022.

[7] R. Yonetani, T. Taniai, M. Barekatain, M. Nishimura, and A. Kanezaki, "Path planning using neural a* search," in *International Conference on Machine Learning*. PMLR, 2021, pp. 12 029–12 039.

[8] OpenAI, "ChatGPT: Engaging openai's conversational ai," https://openai.com/blog/chatgpt, 2023, [Accessed on November 20, 2023].

[9] P. Anderson, Q. Wu, D. Teney, J. Bruce, M. Johnson, N. Sünderhauf, I. Reid, S. Gould, and A. Van Den Hengel, "Vision-and-language navigation: Interpreting visually-grounded navigation instructions in real environments," in *Proceedings of the IEEE conference on computer vision and pattern recognition*, 2018, pp. 3674–3683.

[10] C. Carr, P. Wang, and S. Wang, "A human-friendly verbal communication platform for multi-robot systems: Design and principles," *arXiv preprint arXiv:2211.09519*, 2022.

[11] C. H. Song, J. Wu, C. Washington, B. M. Sadler, W.-L. Chao, and Y. Su, "Llm-planner: Few-shot grounded planning for embodied agents with large language models," in *Proceedings of the IEEE/CVF International Conference on Computer Vision*, 2023, pp. 2998–3009.

[12] A. Xie, Y. Lee, P. Abbeel, and S. James, "Language-conditioned path planning," *arXiv preprint arXiv:2308.16893*, 2023.

[13] M. Aghzal, E. Plaku, and Z. Yao, "Can large language models be good path planners? a benchmark and investigation on spatial-temporal reasoning," *arXiv preprint arXiv:2310.03249*, 2023.

[14] V. Mnih, K. Kavukcuoglu, D. Silver, A. Graves, I. Antonoglou, D. Wierstra, and M. Riedmiller, "Playing atari with deep reinforcement learning," *arXiv preprint arXiv:1312.5602*, 2013.

[15] C. Zhou, X. Lu, J. Dai, B. Huang, X. Liu, and P. Fränti, "Hybrid of representation learning and reinforcement learning for dynamic and complex robotic motion planning," *arXiv preprint arXiv:2309.03758*, 2023.

[16] M. Andrychowicz, A. Raichuk, P. Stańczyk, M. Orsini, S. Girgin, R. Marinier, L. Hussenot, M. Geist, O. Pietquin, M. Michalski, *et al.*, "What matters for on-policy deep actor-critic methods? a large-scale study," in *International conference on learning representations*, 2020.