

Graphormer-Guided Task Planning: Beyond Static Rules with LLM Safety Perception

Wanjing Huang¹, Tongjie Pan², Yalan Ye²

Abstract—Recent advancements in large language models (LLMs) have expanded their role in robotic task planning. However, while LLMs have been explored for generating feasible task sequences, their ability to ensure safe task execution remains underdeveloped. Existing methods struggle with structured risk perception, making them inadequate for safety-critical applications where low-latency hazard adaptation is required. To address this limitation, we propose a Graphormer-enhanced risk-aware task planning framework that combines LLM-based decision-making with structured safety modeling. Our approach constructs a dynamic spatio-semantic safety graph, capturing spatial and contextual risk factors to enable online hazard detection and adaptive task refinement. Unlike existing methods that rely on predefined safety constraints, our framework introduces a context-aware risk perception module that continuously refines safety predictions based on real-time task execution. This enables a more flexible and scalable approach to robotic planning, allowing for adaptive safety compliance beyond static rules. To validate our framework, we conduct experiments in the AI2-THOR environment. The experiments results validates improvements in risk detection accuracy, rising safety notice, and task adaptability of our framework in continuous environments compared to static rule-based and LLM-only baselines. Our project is available at <https://github.com/hwj20/GGTP>

I. INTRODUCTION

A household robot retrieving vegetables from the refrigerator may entirely fail to notice a child approaching a hot stove. Existing robotic task planning frameworks, including rule-based systems and large language model (LLM)-driven planners, focus exclusively on predefined task sequences while lacking low-latency environmental risk perception [1]. These systems evaluate task execution based on internal constraints, ensuring rule compliance but ignoring external hazards. For instance, a rule-based model may enforce that a knife is only used on a cutting board but fails to recognize a child reaching for it as a risk. applications. While LLMs enable flexible, generalizable task reasoning [2], [3], they lack structured risk perception, making them unsuitable for safety-critical Reinforcement learning (RL) approaches [4] can optimize decision-making in controlled settings but require extensive training and fail to generalize in low-latency scenarios [2]. We summarize three key limitations in current approaches to LLM-driven robotic task planning [2], [5]:

- Lack of structured risk perception: Most LLM-based planners operate purely on semantic reasoning, failing

to incorporate spatially structured risk modeling. As a result, they struggle to detect high-risk interactions, such as a child approaching a hazardous object.

- Static rule-based safety mechanisms: Traditional rule-based safety approaches require manually predefined constraints, making them rigid and difficult to scale.
- Limited task adaptability: Existing frameworks typically generate task sequences in a feedforward manner, with limited capacity to modify execution plans in response to safety emergencies.

To address these challenges, we propose a Graphormer-enhanced risk-aware task planning framework that integrates structured safety perception with LLM-driven decision making. We compare with the static rule-based method in Fig. 2. Our work is inspired by the paradigm shift introduced by graph neural networks (GNNs) in autonomous driving, where structured reasoning has proven instrumental in risk prediction and trajectory optimization [6]–[9]. GNN architectures have been extensively applied to model spatial-temporal interactions, facilitating predictive control by encoding agent-object dependencies within a graph-based framework. These approaches [8], [9] employ hierarchical attention mechanisms to selectively prioritize critical relational structures, thereby refining motion forecasting and enhancing decision robustness in unstructured scenarios.

Building upon the above principles, our approach adaptively constructs a spatio-semantic safety graph, encoding both spatial and contextual risk attributes to elevate real-time hazard perception. The Graphormer [10] model is designed to emphasize high-risk interactions through an attention-weighted graph representation, enabling continuous risk assessment and adaptive task refinement. During execution, the system iteratively updates its risk evaluation, ensuring that, upon the emergence of a hazardous condition, a task re-planning mechanism is invoked. Using LLM-based semantic reasoning in conjunction with structured risk modeling, our framework synthesizes online task adaptation with proactive safety measures, significantly improving robustness in open-world robotic applications.

In summary, the main contributions of this study can be outlined as follows:

- We introduce a spatio-semantic safety graph that models environmental risk factors based on Graphormer, providing a structured representation of hazards beyond conventional rule-based methods.
- We develop an adaptive LLM decision framework that enables adaptive task modifications when detecting

¹Wanjing Huang is with the Department of Computer Science, University of California, Santa Barbara.

²Yalan Ye and Tongjie Pan are with the School of Computer Science and Engineering, University of Electronic Science and Technology of China, Chengdu, Sichuan, China. Corresponding email: yalanye@uestc.edu.cn

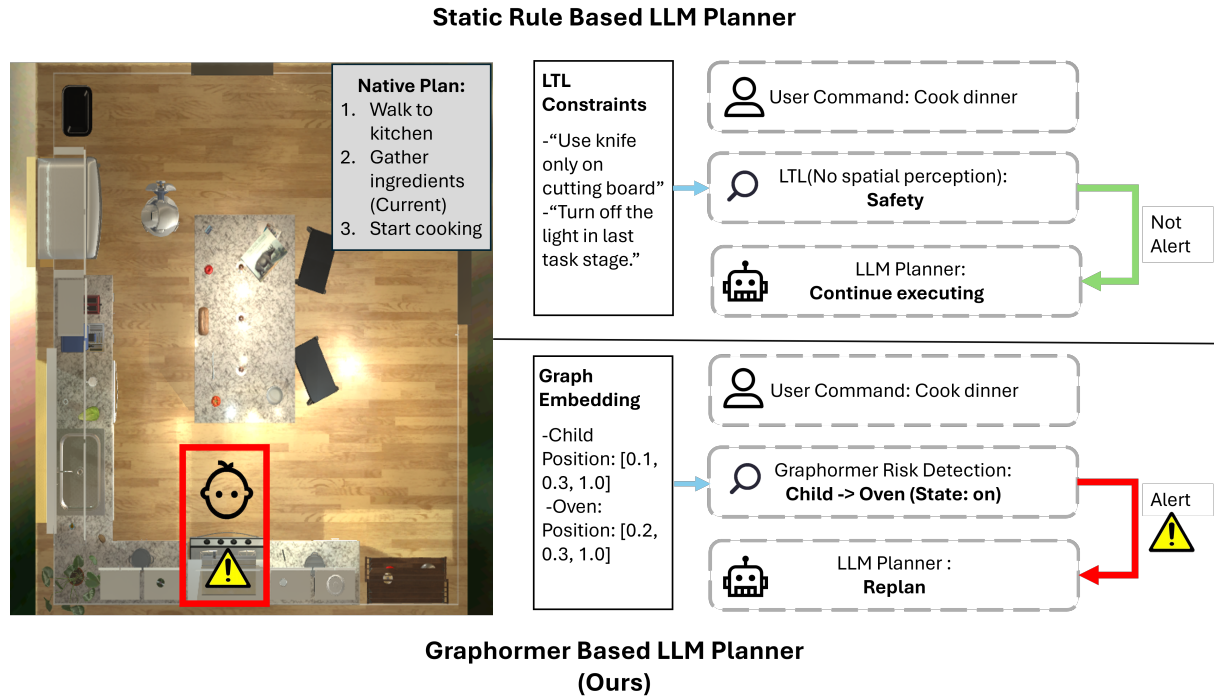


Fig. 1. Comparison between the static rule-based method and our approach. LTL-based safety enforcement relies on predefined constraints, which inherently lack spatial perception and fail to capture unenumerated risks. As a result, when a robot executes a task such as retrieving ingredients from the refrigerator, it is unable to recognize a child approaching a high-temperature oven. Consequently, LTL transmits an incorrect “safe” signal to the LLM planner, leading to a failure in proactive hazard avoidance. In contrast, our method models spatial-semantic relationships, enabling real-time risk assessment and adaptive task modification.

risks, with flexible safety compliance for continuous environments.

- We validate our framework through extensive experiments in AI2-THOR [11], demonstrating substantial improvements in risk detection and adaptive task execution compared to static rule-based and LLM-only baselines.

II. RELATED WORK

A. LLMs in Robotic Task Planning

Recently rich literature has tried to employ LLMs and to accomplish comprehensive planning in complex long-horizon robotics tasks [18]. Studies [2], [5] have shown that LLMs can generate high-level action sequences with world knowledge and semantic reasoning. LLM-driven agents have been applied to domestic robotics [15], multi-agent collaboration [3], [18], and interactive task execution [19].

Despite their remarkable capabilities in task reasoning and generalization [20], LLMs exhibit several fundamental challenges when applied to robotic planning and control. First, LLM-generated task sequences often suffer from inconsistencies and unpredictability, as minor variations in input prompts can lead to significantly different outputs [5]. This lack of determinism raises critical safety concerns, particularly in real-world robotic deployments where failure to adhere to strict operational constraints can result in hazardous consequences.

Second, existing LLM-based planning frameworks rely heavily on prompt engineering, yet there is no standardized methodology for designing prompts that ensure robust and

contextually aware decision-making [5], [21]. Without a structured framework for encoding safety constraints, robots may misinterpret task objectives, leading to unsafe execution behaviors.

B. Rule-Based Safety Constraints in Robotics

Linear Temporal Logic (LTL) has found utility in enforcing formal safety constraints in robotic planning and control [17], [22], [23]. Due to its expressivity and mathematically rigorous semantics, LTL provides a structured approach to specifying task constraints over time, ensuring compliance with pre-defined operational safety rules [24], [25]. LTL-based safety enforcement has been extensively explored in runtime verification and monitoring, where predefined temporal logic specifications are continuously checked against system execution traces to prevent unsafe behaviors [26], [27].

In industrial robotics, LTL has been adopted for programmable logic controllers (PLCs) to ensure compliance with safety standards such as ISO 61508 [28]. LTL has also been applied in human-robot interaction (HRI), where safety constraints are viewed as language-specified conditions that must be satisfied for seamless collaboration between humans and robots [29], [30].

However, while LTL provides strong formal guarantees, it suffers from poor adaptability to continuous environments. LTL-based methods rely on predefined logical constraints [1], which require exhaustive enumeration of all possible safety conditions. This rigidity makes LTL-based approaches

TABLE I
COMPARISON OF RELATED WORK AND OUR APPROACH

Method	LLM Planning	Safety Model	Adaptivity	Real-time Risk Perception
AiGem [6]	✗	Graph-based	✓	✓
SMART-LLM [3]	✓	✗	✗	✗
Plug in the Safety Chip [1]	✓	Rule-based	✗	✗
SafePlanner [12]	✓	Pretrained Safety Prediction Model	✓	✗
RoCo [13]	✓	✗	✓	✗
Co-NavGPT [14]	✓	✗	✓	✗
TidyBot [15]	✓	✗	✗	✗
AutoRL [16]	✓	Human Guardrails	✓	✗
Cross-Layer Sequence Supervision [17]	✓	LTL Constraints	✗	✗
Ours	✓	Graphormer-enhanced	✓	✓

impractical for open-world robotic applications, where environmental hazards can emerge unpredictably.

C. Safety Perception in Robotics

Recent studies have explored integrating vision-language models (VLMs) with LLM-based robotic planning to enhance environmental awareness [12], [16]. While such methods improve perceptual understanding, they primarily focus on task feasibility rather than structured safety reasoning. Furthermore, these approaches rely on human-in-the-loop supervision to ensure compliance, rather than enabling autonomous risk adaptation [16].

In autonomous systems, graph-based models have been adopted for structured spatial data processing and risk modeling. GNNs have been used for trajectory prediction [6], pedestrian intent estimation [31], and interaction-aware motion planning [32] in autonomous driving. Recent studies have also explored graph-based representations for robotic navigation and spatial reasoning [33], [34], enabling agents to model complex object-agent interactions. However, the integration of graph-based risk perception with LLM-driven task planning remains underexplored.

Our work builds upon these advancements by integrating LLM-based semantic reasoning with Graphormer-enhanced risk-aware task adaptation. Unlike prior approaches that either rely solely on LLMs for planning or employ static rule-based safety constraints, our framework constructs a spatio-semantic safety graph that continuously updates during task execution, ensuring immediate hazard detection and proactive task modifications.

To highlight the limitations of existing LLM-driven task planning and rule-based safety mechanisms, we summarize key differences in Table I.

III. METHODOLOGY

A. Problem Formulation

Robotic agents operating in open-world environments must balance task efficiency with operational safety [35]. We propose a framework that integrates a spatio-semantic safety graph with an LLM-driven task planner to achieve risk-aware task execution. Our method constructs a instantaneous risk

representation of the environment, detects high-risk interactions, and triggers adaptive replanning when necessary.

Given an initial task plan $T = \{a_1, a_2, \dots, a_n\}$, where each action a_i corresponds to a specific interaction with an object or agent, our goal is to construct an adaptive spatio-semantic graph $G = (V, E)$, where V represents environmental entities (e.g., objects, humans, robots), and E encodes risk-aware and spatial relationships between them. The agent must update T in response to high-risk edges $e_{ij} \in E$, ensuring safe task execution.

B. Graph-Based Risk Representation

We employ LLM-based semantic reasoning to annotate interactions between entities, assigning risk levels and explainable hazard assessments. Specifically, given a scene with entities $\{v_1, v_2, \dots, v_m\}$, an LLM generates structured risk annotations:

$$r(v_i, v_j) = \text{LLM}(v_i, v_j) \quad (1)$$

$$r \in \{\text{low:0.25, medium:0.50, high:1.00}\}$$

where $r(v_i, v_j)$ defines the inferred risk level between entities v_i and v_j . Example risk annotations include:

```

1 {
2   "type1": "Baby",
3   "type2": "Kettle",
4   "danger_level": "high",
5   "risk_type": [
6     "thermal",
7     "physical",
8     "water"
9   ],
10  "llm_reason": "A kettle, when in use,
11  can become extremely hot and poses a risk
12  of burns or scalds to a baby. Additionally,
13  if a kettle is tipped over, it can lead to
14  significant injury. The risk of water
15  spillage further increases the danger
16  as it can lead to slips or electrical
17  hazards if near power outlets."
18 }

```

We construct a synthetic risk-aware dataset by embedding these risk relationships into AI2-THOR environments. We introduce randomized human-agent interactions (e.g., chil-

dren, adults, pets) to ensure diverse safety-critical scenarios. Each sample is labeled with a danger score, computed as:

$$S(v_i, v_j) = r(v_i, v_j) \times \text{SP}(v_i, v_j) \quad (2)$$

where the SP refers to spatial proximity, which accounts for spatial scene configurations, and in experiment part we chose:

$$\text{SP} = \begin{cases} \frac{1}{\text{distance}}, & \text{if distance} \leq \text{DT} \\ \frac{1}{\text{DT}}, & \text{otherwise} \end{cases} \quad (3)$$

where DT refers to the threshold to prevent SP from diverging.

C. Graphormer Pretraining for Risk Detection

We train a Graphormer [10] model to predict high-risk interactions based on the structured safety graph. Given an environment represented as a graph $\mathcal{G} = (\mathcal{V}, \mathcal{E})$ with node features $X \in \mathbb{R}^{|\mathcal{V}| \times d}$ and adjacency matrix $A \in \mathbb{R}^{|\mathcal{V}| \times |\mathcal{V}|}$, Graphormer computes risk-weighted attention scores via:

$$H = \text{Graphormer}(X, A, E) \quad (4)$$

where $H \in \mathbb{R}^{|\mathcal{V}| \times d}$ represents the learned risk embeddings, and $E \in \mathbb{R}^{|\mathcal{E}| \times d}$ encodes edge features. The model applies multi-head attention over graph structures:

$$Z = \text{MultiHeadAttn}(Q, K, V) = \sum_{i=1}^h \alpha_i W_i V \quad (5)$$

where queries, keys, and values are computed as:

$$Q = W_Q H, \quad K = W_K H, \quad V = W_V H \quad (6)$$

The attention weights α_i incorporate both structural and semantic information:

$$\alpha_{ij} = \frac{\exp\left(\frac{(QW_Q)(KW_K)^T}{\sqrt{d}} + b_{ij}\right)}{\sum_k \exp\left(\frac{(QW_Q)(KW_K)^T}{\sqrt{d}} + b_{ik}\right)} \quad (7)$$

where b_{ij} represents edge encodings extracted from E . To address data imbalance—where most edges are non-hazardous—we employ a focal loss [36]:

$$\mathcal{L} = - \sum_{(i,j) \in \mathcal{E}} w_{ij} (1 - p_{ij})^\gamma \log p_{ij} \quad (8)$$

where p_{ij} is the predicted probability of an edge being hazardous, w_{ij} is a re-weighting factor prioritizing risky edges, and γ adjusts the focus on hard-to-classify edges.

D. Risk-Aware Task Replanning

At runtime, Graphormer continuously evaluates the environment, identifying high-risk edges e_{ij} in the safety graph. When a high-risk interaction is detected, the LLM is queried to generate a revised task plan:

$$T' = \text{LLM}(T, G) \quad \text{subject to safety constraints.} \quad (9)$$

The updated task plan incorporates safety-aware modifications, ensuring secure execution. An example of updated plan is as follows:

Task: Prepare a meal

Description: Cook the meal.

LLM Initial Plan:

0. Walk to kitchen
1. Gather ingredients
2. Start cooking

Graphormer Risk Detection: Upon entering the kitchen, the system updates the environment state and detects a high-risk interaction:

- High-risk edge detected:
Baby \rightarrow Knife (Risk level: High)
- Reason: A child is standing close to the knife, posing a potential injury hazard.

Replanning Triggered: The LLM generates a revised task sequence to mitigate the identified risk.

LLM Updated Plan:

0. Walk to kitchen
1. Ensure child is in a safe location
2. Secure knife in a designated area
3. Gather ingredients
4. Start cooking
5. DONE

IV. EXPERIMENTS

A. Risk Detection Experiments

To evaluate our approach, we conduct experiments in AI2-THOR using a dataset of 120 diverse household environments, covering four common domestic settings: kitchens, living rooms, bedrooms, and bathrooms. To simulate real-world safety-critical scenarios, we manually introduce human agent nodes near hazardous objects (e.g., a child near a stove or a knife) in half of the scenes. The dataset is split into 90 training samples, 15 validation samples, and 15 test samples. We evaluate our method against two baseline approaches to analyze its performance in in-time, risk-aware task execution.

1) *Evaluation Metrics and Data Imbalance:* Our dataset exhibits significant class imbalance, with hazardous edges constituting only 1% of the total edges. A naive classifier that predicts all edges as "safe" would trivially achieve 99% accuracy while failing to identify critical safety risks.

2) *Comparison with Alternative Methods:* The results confirms that:

- **Random Guess Baseline:** Performs no better than chance, with near-zero precision.
- **Rule-Based System:** Requires manual rule specification, which is inherently limited. No existing rule set comprehensively covers all indoor safety hazards, and manually defining exhaustive safety constraints is infeasible. In addition, static rule enforcement lacks adaptability, making it ineffective in continuous environments where unforeseen risks arise.

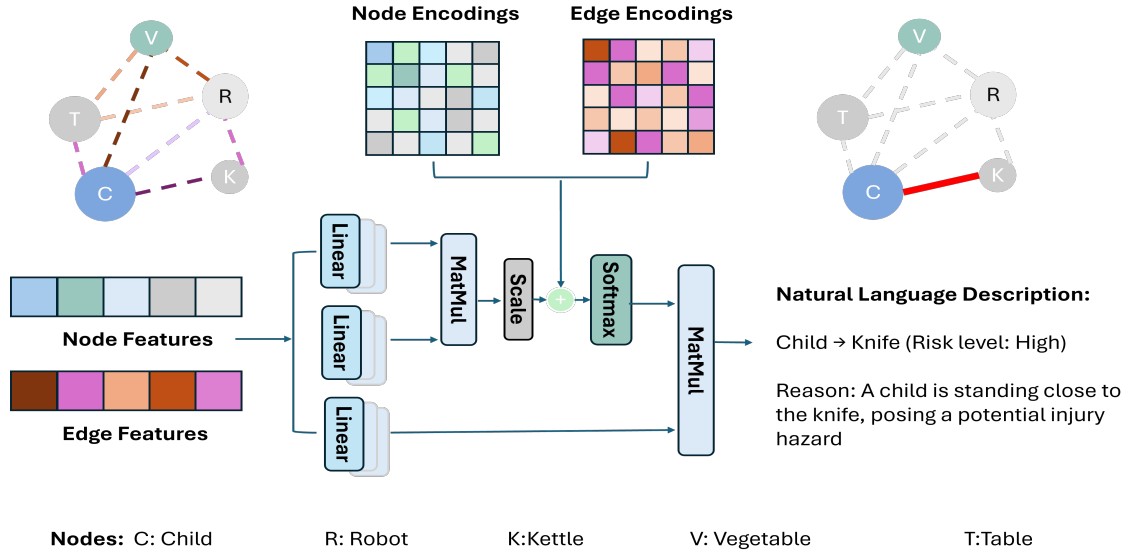


Fig. 2. Graphormer Overview. Our framework integrates Graphormer-based risk modeling with LLM-driven task planning to enable instantaneous safety adaptation. The system constructs a context-aware spatio-semantic safety graph from environmental observations, where high-risk interactions are identified using attention-weighted edge representations. The Graphormer will auto translate the dangerous edges into natural language.

- **Our Model:** Achieves the best trade-off between hazard detection and task efficiency, with online risk perception to assess and adapt to environmental threats.

Figure 3 illustrates the trade-off between Recall and Precision across different models. To ensure a meaningful evaluation, we prioritize Recall as our primary metric, ensuring that hazardous interactions are detected. Specifically, we select a decision threshold=0.21 that yields:

- **Recall:** 91.39% (covering 90% of hazardous edges)
- **Precision:** 29.27% (filtering out 70% of false positives)

This threshold ensures that safety-critical hazards are not overlooked while maintaining practical alert filtering.

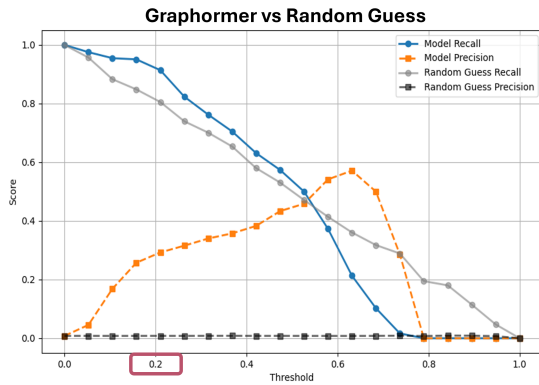


Fig. 3. Precision-Recall analysis of our method. Due to the severe class imbalance in our dataset, random guessing yields near-zero precision. In contrast, by setting a decision threshold of 0.21, our method achieves a precision of 30% while maintaining a recall above 90%. This means that in a dataset with 10,000 edges, where only 100 are hazardous, our model identifies 300 edges as potentially dangerous, successfully capturing 90 out of the 100 true hazardous edges. This balance ensures that critical risks are detected while minimizing false alarms.

B. LLM-Guided Risk-Aware Task Planning Experiments

To evaluate our proposed framework, we define 5 household activities categorized in three levels of complexities, including two spatial reorganization tasks (simple), two object manipulation tasks (intermediate), and one cooking tasks (complex). Each task requires no high-risk actions to be executed without proper intervention. We conducted real task planning experiments for these methods:

- **LLM-Only (No Risk Awareness):** The LLM receives the entire scene description and generates a task sequence without explicit risk perception.
- **Prompted LLM (Risk-Aware Prompting):** The LLM is provided with a general instruction to consider safety risks but only receives the information of all objects without guidelines from environment.
- **LLM + Static Rules (LTL):** The LLM is given a general instruction to consider safety risks and receives signals from the LTL module. In our experiments, we manually define rules to detect all hazards; however, this approach is not feasible in real-world settings.
- **LLM + Graphormer (Ours):** Our model integrates risk-aware graph-based reasoning with LLM-based planning.

1) *Evaluation Metrics:* For additional analysis, we test our framework in AI2-THOR first 20 kitchen scenes with evaluation metrics:

- **Task Success Rate (TSR):** The percentage of successfully completed tasks.
- **Safety Notice Rate (SNR):** The percentage of scenarios where a high-risk situation is noticed by agents.
- **Risk Handling Success (RHS):** The percentage of cases where the system correctly identified and mitigated a potential hazard.

TABLE II
TASK PLANNING PERFORMANCE ACROSS DIFFERENT TASKS FOR THREE CATEGORIES OF COMPLEXITIES

Model	Simple			Intermediate			Complex		
	TSR (%)	SNR (%)	RHS (%)	TSR (%)	SNR (%)	RHS (%)	TSR (%)	SNR (%)	RHS (%)
LLM-Only	100	0	0	97.5	0	0	85.0	0	10
LLM-Safe-Prompting	100	0	0	97.5	0	0	85.0	0	0
LLM + LTL	95.0	100	90	90.0	100	80	60.0	100	40
LLM + Graphormer (Ours)	100	100	100	97.5	100	95.0	92.5	100	95.0



Fig. 4. Stages of a complete cooking task in AI2-THOR (FloorPlan2) from the perspective of the executing agent. The task includes picking up ingredients and placing them into a pan. The first stage, "HandleSafetyIssue" for the "Baby", is not natively supported in AI2-THOR; however, we implement a script-based check to determine whether the model actively resolves safety issues.

To further validate generalization capabilities, we introduce modified flexible human-agent interactions within AI2-THOR, simulating a child is near a dangerous object (e.g. knife).

2) *Results and Analysis*: As shown in Table II, our approach verifies superior risk detection and task modification capabilities compared to baselines.

LLM-based planners fail to recognize dynamic risk relationships, even when prompted for safety. They may identify hazardous objects (e.g., knives) but lack structured perception to detect evolving threats (e.g., a child approaching).

LTL-based methods achieve high accuracy in controlled settings due to manually predefined constraints, but real-world hazard enumeration is infeasible. Their static nature causes them to misclassify unforeseen risks as "safe," misleading the LLM.

Our method surpasses baselines in unstructured risk scenarios by leveraging a risk-aware safety graph, enabling proactive hazard detection and adaptive task planning for safer execution in dynamic environments.

Fig. 4 illustrates task execution in AI2-THOR, highlighting our system's ability to adaptively modify task sequences based on emerging risks.

V. CONCLUSION AND FUTURE WORK

In this work, we introduced a graphormer-enhanced risk-aware task planning framework that integrates structured safety perception with LLM-driven decision-making for robotic agents in real-time environments. Through extensive

evaluations in AI2-THOR, we validated the effectiveness of our method in three key aspects:

- **Risk Prediction Accuracy**: Our Graphormer model achieves high precision and recall in hazard identification.
- **Risk-Aware Task Planning**: Our LLM + Graphormer framework successfully adapts task sequences based on environmental risks, reducing unsafe actions while maintaining high task efficiency.
- **Generalization in Real-time Environments**: Our method consistently outperforms baselines in the rate of safety notice and handling safety issues when tested in AI2-THOR settings.

While our framework demonstrates strong performance in risk-aware robotic task planning, several directions remain open for future exploration:

- Expanding the AI2-THOR environment with a broader set of interactive, safety-related objects to enhance risk-aware interactions. Specifically, we aim to release a systematically curated dataset to benchmark models on hazardous factor perception.
- Extending our current static spatial risk perception to a temporal framework, enabling not just the recognition of hazards but also their prediction over time.

By addressing these challenges, we aim to further advance the intersection of LLM-driven task planning and graph-based safety modeling, paving the way for safer, more intelligent autonomous systems in unstructured environments.

APPENDIX

TABLE III
COMPARISON OF TASK EXECUTION TIMES

Stage	Graphormer (seconds)	LTL (seconds)
Retrieve Object Information	0.1552	0.1257
Build Environment Graph	0.0185	-
Receive Safety Notice	1.1926	0.0002 (1 rule)
Generate Task Sequence	1.4581	0.9931
Parse Task Sequence	0.0000	0.0000

Table III presents a comparison of task execution times for a single task between our method, Graphormer (loading model from disk), and LTL. It is important to emphasize that in the LTL setup, we used only one rule, whereas real-world applications require thousands of safety rules. As a result, the

majority of time consumption is concentrated on waiting for the LLM to generate the task sequence, rather than on risk detection.

REFERENCES

- [1] Z. Yang, S. S. Raman, A. Shah, and S. Tellex, "Plug in the safety chip: Enforcing constraints for llm-driven robot agents," in *2024 IEEE International Conference on Robotics and Automation (ICRA)*, pp. 14435–14442, IEEE, 2024.
- [2] Y. Kim, D. Kim, J. Choi, J. Park, N. Oh, and D. Park, "A survey on integration of large language models with intelligent robots," *Intelligent Service Robotics*, vol. 17, no. 5, pp. 1091–1107, 2024.
- [3] S. Kannan, V. Venkatesh, and B. C. S.-L. Min, "Smart multi-agent robot task planning using large language models," *arXiv preprint arXiv:2309.10062*, 2023.
- [4] L. Ouyang, J. Wu, X. Jiang, D. Almeida, C. Wainwright, P. Mishkin, C. Zhang, S. Agarwal, K. Slama, A. Ray, *et al.*, "Training language models to follow instructions with human feedback," *Advances in neural information processing systems*, vol. 35, pp. 27730–27744, 2022.
- [5] H. Jeong, H. Lee, C. Kim, and S. Shin, "A survey of robot intelligence with large language models," *Applied Sciences*, vol. 14, no. 19, p. 8868, 2024.
- [6] J. Samiuddin, B. Boulet, and D. Wu, "Trajectory prediction for autonomous driving using agent-interaction graph embedding," *arXiv preprint arXiv:2410.23298*, 2024.
- [7] J. Gao, C. Sun, H. Zhao, Y. Shen, D. Anguelov, C. Li, and C. Schmid, "Vectormet: Encoding hd maps and agent dynamics from vectorized representation," in *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*, pp. 11525–11533, 2020.
- [8] X. Kong, W. Xing, X. Wei, P. Bao, J. Zhang, and W. Lu, "Stgat: Spatial-temporal graph attention networks for traffic flow forecasting," *IEEE Access*, vol. 8, pp. 134363–134372, 2020.
- [9] N. Nayakanti, R. Al-Rfou, A. Zhou, K. Goel, K. S. Refaat, and B. Sapp, "Wayformer: Motion forecasting via simple & efficient attention networks," in *2023 IEEE International Conference on Robotics and Automation (ICRA)*, pp. 2980–2987, IEEE, 2023.
- [10] C. Ying, T. Cai, S. Luo, S. Zheng, G. Ke, D. He, Y. Shen, and T.-Y. Liu, "Do transformers really perform badly for graph representation?," *Advances in neural information processing systems*, vol. 34, pp. 28877–28888, 2021.
- [11] E. Kolve, R. Mottaghi, W. Han, E. VanderBilt, L. Weihs, A. Herrasti, D. Gordon, Y. Zhu, A. Gupta, and A. Farhadi, "AI2-THOR: An Interactive 3D Environment for Visual AI," *arXiv*, 2017.
- [12] S. Li, Z. Ma, F. Liu, J. Lu, Q. Xiao, K. Sun, L. Cui, X. Yang, P. Liu, and X. Wang, "Safe planner: Empowering safety awareness in large pre-trained models for robot task planning," *arXiv preprint arXiv:2411.06920*, 2024.
- [13] 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)*, pp. 286–299, IEEE, 2024.
- [14] B. Yu, H. Kasaei, and M. Cao, "Co-navgpt: Multi-robot cooperative visual semantic navigation using large language models," *arXiv preprint arXiv:2310.07937*, 2023.
- [15] J. Wu, R. Antonova, A. Kan, M. Lepert, A. Zeng, S. Song, J. Bohg, S. Rusinkiewicz, and T. Funkhouser, "Tidybot: Personalized robot assistance with large language models," *Autonomous Robots*, vol. 47, no. 8, pp. 1087–1102, 2023.
- [16] M. Ahn, D. Dwibedi, C. Finn, M. G. Arenas, K. Gopalakrishnan, K. Hausman, B. Ichter, A. Irpan, N. Joshi, R. Julian, *et al.*, "Autort: Embodied foundation models for large scale orchestration of robotic agents," *arXiv preprint arXiv:2401.12963*, 2024.
- [17] Z. Wang, Q. Liu, J. Qin, and M. Li, "Ensuring safety in llm-driven robotics: A cross-layer sequence supervision mechanism," in *2024 IEEE/RSJ International Conference on Intelligent Robots and Systems (IROS)*, pp. 9620–9627, IEEE, 2024.
- [18] M. Dalal, T. Chiruvolu, D. Chaplot, and R. Salakhutdinov, "Plan-seq-learn: Language model guided rl for solving long horizon robotics tasks," *arXiv preprint arXiv:2405.01534*, 2024.
- [19] Y. Du, K. Konyushkova, M. Denil, A. Raju, J. Landon, F. Hill, N. de Freitas, and S. Cabi, "Vision-language models as success detectors," *arXiv preprint arXiv:2303.07280*, 2023.
- [20] J. Achiam, S. Adler, S. Agarwal, L. Ahmad, I. Akkaya, F. L. Aleman, D. Almeida, J. Altenschmidt, S. Altman, S. Anadkat, *et al.*, "Gpt-4 technical report," *arXiv preprint arXiv:2303.08774*, 2023.
- [21] Y.-J. Wang, B. Zhang, J. Chen, and K. Sreenath, "Prompt a robot to walk with large language models," *arXiv preprint arXiv:2309.09969*, 2023.
- [22] H. Barringer, A. Goldberg, K. Havelund, and K. Sen, "Rule-based runtime verification," in *Verification, Model Checking, and Abstract Interpretation: 5th International Conference, VMCAI 2004 Venice, Italy, January 11-13, 2004 Proceedings 5*, pp. 44–57, Springer, 2004.
- [23] Y. Wu, Z. Xiong, Y. Hu, S. S. Iyengar, N. Jiang, A. Bera, L. Tan, and S. Jagannathan, "Selp: Generating safe and efficient task plans for robot agents with large language models," *arXiv preprint arXiv:2409.19471*, 2024.
- [24] D. S. Grigorev, A. K. Kovalev, and A. I. Panov, "Common sense plan verification with large language models," in *International Conference on Hybrid Artificial Intelligence Systems*, pp. 224–236, Springer, 2024.
- [25] T. Reinbacher, K. Y. Rozier, and J. Schumann, "Temporal-logic based runtime observer pairs for system health management of real-time systems," in *Tools and Algorithms for the Construction and Analysis of Systems: 20th International Conference, TACAS 2014, Held as Part of the European Joint Conferences on Theory and Practice of Software, ETAPS 2014, Grenoble, France, April 5-13, 2014. Proceedings 20*, pp. 357–372, Springer, 2014.
- [26] P. Thati and G. Roşu, "Monitoring algorithms for metric temporal logic specifications," *Electronic Notes in Theoretical Computer Science*, vol. 113, pp. 145–162, 2005.
- [27] J. Schumann, P. Moosbrugger, and K. Y. Rozier, "R2u2: monitoring and diagnosis of security threats for unmanned aerial systems," in *Runtime Verification: 6th International Conference, RV 2015, Vienna, Austria, September 22-25, 2015. Proceedings*, pp. 233–249, Springer, 2015.
- [28] D. J. Smith and K. G. Simpson, *The safety critical systems handbook: a straightforward guide to functional safety: IEC 61508 (2010 Edition), IEC 61511 (2015 edition) and related guidance*. Butterworth-Heinemann, 2020.
- [29] B. J. Choi, J. Park, and C. H. Park, "Formal verification for human-robot interaction in medical environments," in *Companion of the 2021 ACM/IEEE International Conference on Human-Robot Interaction*, pp. 181–185, 2021.
- [30] S. Srinivas, R. Kermani, K. Kim, Y. Kobayashi, and G. Fainekos, "A graphical language for ltl motion and mission planning," in *2013 IEEE International Conference on Robotics and Biomimetics (ROBIO)*, pp. 704–709, IEEE, 2013.
- [31] H. Zhou, D. Ren, H. Xia, M. Fan, X. Yang, and H. Huang, "Ast-gnn: An attention-based spatio-temporal graph neural network for interaction-aware pedestrian trajectory prediction," *Neurocomputing*, vol. 445, pp. 298–308, 2021.
- [32] J. Li, Z. Su, and Y. Qiu, "Dynamic motion planning model for multirobot using graph neural network and historical information," *Advanced Intelligent Systems*, vol. 5, no. 8, p. 2300036, 2023.
- [33] F. Zhang, C. Xuan, H.-K. Lam, and S.-H. Tsai, "Navigation control of mobile robots using graph neural network and reinforcement learning with fuzzy reward," in *2022 International Conference on Fuzzy Theory and Its Applications (iFUZZY)*, pp. 01–06, IEEE, 2022.
- [34] A. Imran, G. Beltrame, and D. St-Onge, "Gnn-based decentralized perception in multirobot systems for predicting worker actions," *arXiv preprint arXiv:2501.04193*, 2025.
- [35] J. S. Albus and J. S. Albus, *A reference model architecture for intelligent systems design*. Citeseer, 1994.
- [36] T.-Y. Lin, P. Goyal, R. Girshick, K. He, and P. Dollár, "Focal loss for dense object detection," in *Proceedings of the IEEE international conference on computer vision*, pp. 2980–2988, 2017.