
Mind-Map Agent: Enhancing Cooperative Task Planning through Communication Alignment with Large Language Models

Hobeom Jeon¹ Hyungmin Kim¹ Dohyung Kim^{1,2*}
Minsu Jang^{1,2} Jaehong Kim²

¹ University of Science and Technology, South Korea

² Electronics and Telecommunications Research Institute, South Korea
hobeomjeon@ust.ac.kr, dhkim008@etri.re.kr

Abstract

Embodied agents that collaborate with humans through natural language have become an active area of research, offering flexibility in cooperative planning and execution. Debate-based approaches often depend on repeated consensus procedures, which can increase dialogue frequency and risk over-communication. At the same time, LLMs are prone to hallucination during dialogue processing, sometimes causing confabulation and reducing consistency in long-term strategies. We introduce the Mind-Map Agent, an approach that guides reasoning with explicit cooperative strategies while maintaining structured long-term memory to disentangle dialogue, task state, and planning context. The generated Mind-Maps support coherent long-horizon planning, reduce redundant dialogue, and enhance interpretability in multi-agent interaction. Evaluations on Communicative Watch-and-Help and ThreeDWorld Multi-Agent Transport indicate that the Mind-Map Agent achieves more stable efficiency compared to classical planners and LLM agents across different model scales and environments. Our results suggest that Mind-Map reasoning enables cooperative agents to accomplish tasks with fewer conversations while sustaining effective collaboration.

1 Introduction

Recent advances in Large Language Models (LLMs) have enabled new paradigms in multi-agent embodied tasks, particularly in natural language-based reasoning, planning, and decision-making [1, 2]. Traditional multi-agent systems, while reliable under predefined communication protocols [3], often struggle in decentralized collaboration where agents must infer teammates’ intentions [4, 5] and dynamically coordinate without centralized control. The challenge becomes even more pronounced in human–robot cooperation, where free-form human communication frequently exceeds the capability of protocol-bound robot agents [6, 7]. To address these issues, recent research has explored LLM-powered agents that communicate through natural language, with the Cooperative Embodied Language Agent (CoELA) framework serving as a representative example [5].

CoELA advances multi-agent interaction by enabling free-form dialogue, but its decision-making largely focuses on immediate, single-step actions without maintaining structured long-term plans. In contrast, Cooperative Plan Optimization (CaPo) [8] introduces meta-planning to enhance coordination, but relies on rigid, rule-based communication protocols. In such protocols, robots communicate by following predefined discussion rules or debate-like turn-taking structures, which ensures precise

*Corresponding Author

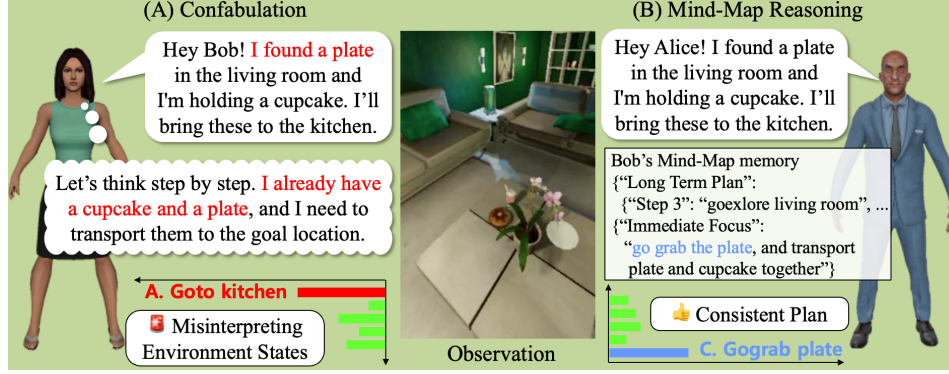


Figure 1: Comparison between confabulation and Mind-Map Reasoning in multi-agent communication. (A) Confabulation: The agent misinterprets dialogue as executed action, assuming possession of an object and leading to incomplete task execution. (B) Mind-Map Reasoning: By explicitly separating long-term plan and current state, the agent identifies missing objects and selects consistent actions, thereby mitigating confabulation and improving cooperation efficiency.

progress tracking in robot–robot settings [9]. However, this rigidity creates two practical limitations. First, repeated discussion sessions can introduce substantial communication overhead [10]. Second, when interacting with humans who use free-form and informal language, these strict rules often fail to support accurate and natural communication. Effective cooperation therefore requires agents that can infer goals from unconstrained dialogue, maintain consistent long-term strategies, and adapt plans dynamically without depending on strict communication rules.

Another critical barrier emerges in hallucination and confabulation [11, 12], where LLMs may misinterpret dialogue history as factual state, leading to faulty inferences and redundant actions. This risk becomes more likely with unconstrained Chain of Thought (CoT) prompting [13], where an agent may conflate spoken intentions with executed actions or assume environmental changes that have not actually occurred. For instance, as shown in Figure 1(A), an agent might treat a previous utterance as if the action had already been completed, disrupting task execution. In comparison, Figure 1(B) illustrates how a structured Mind-Map representation that separates dialogue, current state, and long-term plan can support more consistent decision-making, which in turn may reduce unnecessary retries and communication overhead [5, 14].

To address these limitations, we introduce the Mind-Map Agent, an LLM-powered collaborative agent that integrates long-term planning with structured reasoning. At the core of our approach is Mind-Map reasoning, which maintains explicit representations of teammates’ intentions, cooperative strategies, and communication plans. This structured representation helps agents remain consistent across interactions, mitigating confabulation and reducing redundant dialogue. By jointly capturing conversational cues and factual states, the approach is also applicable to human–robot collaboration, where agents must flexibly interpret free-form instructions without relying on rigid communication protocols.

- We introduce the Mind-Map Agent, which equips each agent with structured long-term reasoning by explicitly separating dialogue, task state, and cooperative strategies. This design enhances consistency in multi-step execution while reducing confabulation and redundant communication.
- We conduct extensive evaluation across diverse open-source LLMs, showing that Mind-Map reasoning consistently outperforms CoELA. These results demonstrate that structured reasoning provides robustness across model scales and architectures.
- Our analysis highlights both the benefits and limitations of Mind-Map reasoning. While it substantially mitigates confabulation, lightweight models remain more vulnerable in visual observation settings, underscoring ongoing challenges for current LLMs in generating reliable, grounded dialogue in collaborative tasks.

2 Related Work

2.1 Multi-Agent Cooperation

Research on multi-agent cooperation has developed a wide range of methods for coordination, communication, and decision-making in decentralized environments [3]. Early approaches such as MADDPG [15] introduced reinforcement learning techniques to improve strategic coordination among agents in cooperative–competitive settings [16]. Subsequent work explored cooperation in embodied household tasks [2], communication efficiency [8, 17], and dialogue grounding in environmental contexts [18]. Despite these advances, many methods still rely on opaque message-passing channels such as continuous vector embeddings [18], which limit interpretability and adaptability across heterogeneous agents. Moreover, protocol-driven communication and centralized planning, long established in robotic systems with networked control, scale poorly to open human-facing environments. For multi-agent systems to collaborate effectively with humans, agents must make autonomous decisions under partial observability [19, 4], coordinating tasks without globally shared knowledge or fixed communication protocols.

2.2 Cooperation with Natural Language Communication

Natural language has emerged as a promising alternative for enhancing cooperation, task planning, and information sharing in embodied multi-agent systems. Frameworks such as CoELA [5] integrate large language models (LLMs) to support flexible dialogue between agents and humans, enabling partners to share intentions and adapt cooperative behaviors. Built on the DEC-POMDP framework [20, 21] and echoing cognitive architectures such as Soar [22], CoELA organizes decision-making into perception, memory, communication, planning, and execution modules. While these systems benefit from the reasoning capabilities of LLMs, challenges arise when models trained on general dialogue corpora are used for situated collaboration. Agents may generate irrelevant or emotional messages (e.g., utterances with exaggerated punctuation or emoji-like expressions), hallucinate actions such as declaring that an object has been grasped before execution, or even confuse roles by redundantly performing tasks already delegated to teammates. Such forms of confabulation occur when generated text is mistaken for factual observation or when dialogue history is misinterpreted as action history, leading to failures in planning and coordination. These limitations highlight the need for methods that constrain utterances to task-relevant content and incorporate structured reasoning to ensure reliable cooperation.

2.3 Optimizing Cooperation and Communication Strategy

Beyond enabling natural language communication, another line of work has focused on optimizing the balance between communication frequency and cooperative efficiency [8, 10, 9]. For instance, CaPo [8] introduces meta-planning to generate shared high-level strategies and adjust execution adaptively, while REVECA [9] constrains dialogue through structured templates to support mutual awareness. While effective in controlled simulations, these frameworks assume rigid protocols such as always-on channels, explicit turn-taking, or repeated deliberation rounds. Such mechanisms provide transparency in robot–robot cooperation but introduce significant overhead [10] and reduce adaptability in human–robot interaction, where free-form dialogue is the norm. In practice, repeated meta-planning debates and verbose exchanges risk overwhelming both agents and human partners, especially when lightweight LLMs amplify “chattering” behavior that contaminates dialogue context and increases confabulation. These challenges highlight the need for mechanisms that support reliable, interpretable, and efficient communication without rigid protocol assumptions.

3 Mind Generation for Cooperative Planning

3.1 Mind-Map Agent: System Architecture and Operating Principle

Figure 2 presents the overall architecture of the Mind-Map Agent. The system is composed of modular components spanning from the internal agent modules to the external task environment. At its core, the agent operates with two distinct phases of large language model (LLM) reasoning: (1) Mind-Map Generation and (2) Plan Selection. Both phases make use of two complementary memory

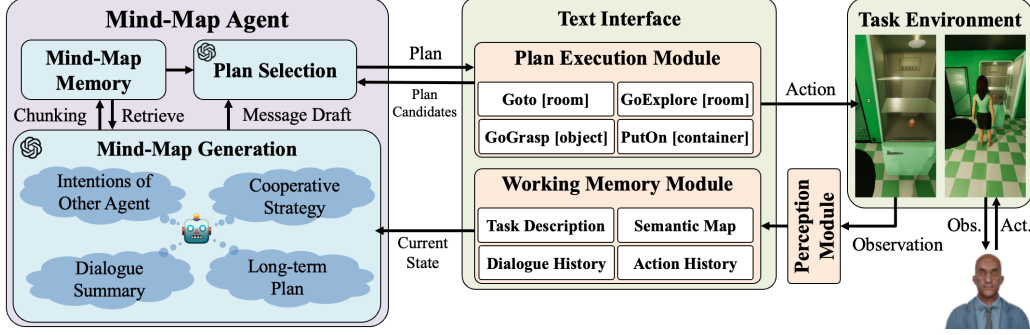


Figure 2: Overall architecture of the Mind-Map Agent. The system integrates two phases of LLM-based reasoning (Mind-Map Generation and Plan Selection) with complementary memory systems (Text Interface as working memory and Mind-Map Memory as structured long-term memory). The architecture supports perception-action loops while maintaining structured knowledge for long-horizon cooperative planning.

systems. The Text Interface functions as the agent’s working memory, while the Mind-Map Memory provides a structured long-term memory that is explicitly maintained by the agent.

The Mind-Map Generation phase constitutes the first reasoning step. In this step, the LLM takes the current state from the Text Interface and integrates it with the stored Mind-Map Memory. The output is an updated structured mind map that may include interpretations of intentions, cooperative strategies, dialogue summaries, and long-term plans. This process combines incoming observations with previously stored task knowledge. By refreshing the Mind-Map Memory over time, the agent can maintain a stable representation that persists across interactions. This design contributes to a clearer separation between short-term environmental descriptions and the structured long-term task context, which may help reduce confusion between immediate perception and accumulated memory.

The Plan Selection phase forms the second reasoning step. In this step, the LLM references both the updated Mind-Map Memory and the Text Interface. Based on this combined context, the model determines the next action or communication. Outputs are passed to the Plan Execution Module, which translates high-level plans into primitive operations such as navigation or object manipulation. Alternatively, the Plan Selection may generate a dialogue message to be transmitted to teammates. The reliance on both working memory and long-term structured memory is intended to keep each decision connected to the immediate environment while also aligned with the broader cooperative plan.

The Text Interface serves as the working memory of the agent and includes four elements: task description, dialogue history, semantic map, and action history. This interface reflects what the agent can directly access at runtime and interacts with the task environment through the Perception Module and the Plan Execution Module. The Perception Module converts raw environment signals into symbolic observations that are stored in the Text Interface. The Plan Execution Module translates symbolic action plans into concrete operations that affect the environment and generate new observations, thereby closing the perception-action loop. In contrast, the Mind-Map Memory maintains a structured representation that is chunked and retrieved explicitly through the LLM. This long-term memory captures higher-level constructs such as inferred intentions, cooperative strategies, and reasoning outcomes. The distinction between these two memory systems is important: the working memory is transient and environment-driven, whereas the long-term memory is agent-managed and schema-based.

This dual-phase architecture draws on insights from cognitive architectures and multi-agent cooperation research, which emphasize that effective collaboration depends on the ability to infer and represent the hidden goals of others, often described as Theory of Mind (ToM) [23, 24, 25]. Building on these principles, the Mind-Map Agent transforms implicit intention inference into explicit structured entries that can be transparently maintained and reasoned over. This design contrasts with prior frameworks such as CoELA [5], which pioneered natural language based coordination in multi-agent systems but primarily relied on step-by-step reactive reasoning without durable long-term structures. Inspired by recent advances that frame LLMs as both reasoning agents and generators of structured knowledge [26], as well as work on multi-persona collaboration for emergent reasoning [27], the

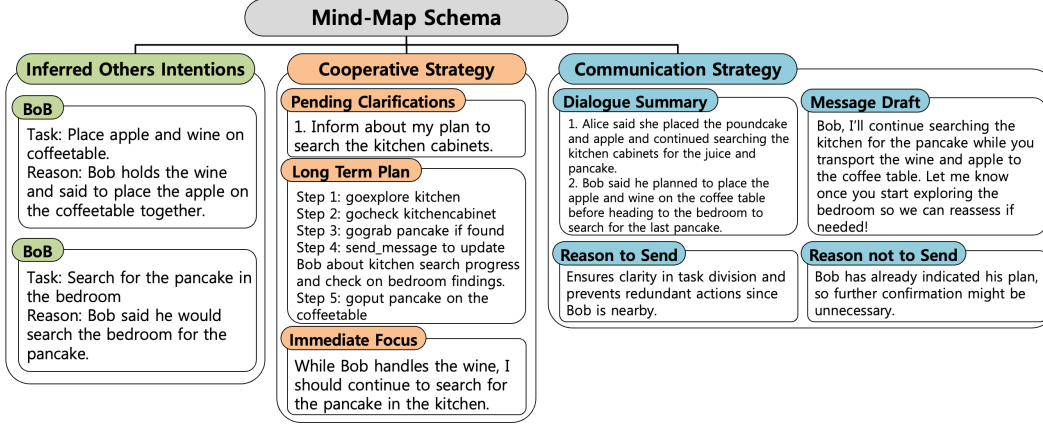


Figure 3: Mind-Map schema as maintained by the agent. The schema is divided into three main components: inferred others’ intentions, cooperative strategy, and communication strategy. Together, they provide a structured representation of team state and guide reasoning over multi-agent and human–robot collaboration.

Mind-Map Agent introduces structured memory to clearly separate dialogue history, environmental state, and cooperative strategy. This explicit representation reduces redundant or ambiguous utterances and enables reliable tracking of completed versus pending actions, thereby supporting more robust long-horizon cooperation in both multi-agent and human–robot contexts.

3.2 Mind-Map Schema and Key Components

At the core of our approach is the Mind-Map schema, a structured representation of the agent’s understanding of team state and task progress. During the Mind-Map Generation phase, the agent is guided by in-context exemplars expressed in a JSON format, ensuring that updates remain structured and consistent. The schema is dynamically updated as new observations and utterances are received. As illustrated in Figure 3, the schema consists of three main components: inferred others’ intentions, cooperative strategy, and communication strategy. Together, they enable the agent to track teammates’ goals, maintain alignment with the joint plan, and regulate communication to avoid redundant or irrelevant utterances.

Inferred Others’ Intentions This component guides the agent to interpret teammates’ utterances as explicit expressions of intent and transform them into structured entries within the mind map. In this way, natural language is not only treated as surface-level communication but also as a representation of hidden goals and action plans. By recording and organizing these inferred intentions, the agent externalizes a Theory-of-Mind perspective [23, 24], enabling it to track what others aim to achieve and why. This structured account of intent can be reused in subsequent reasoning steps, allowing the agent to anticipate teammates’ actions, align its own behavior with theirs, and reduce the need for repeated clarification. In this way, the process provides a systematic mechanism for capturing language-mediated intentions, supporting consistent interpretation of teammates’ goals and facilitating alignment in cooperative decision-making.

Cooperative Strategy This component organizes the agent’s reasoning about how to align its own role with the collective plan. It is structured into three layers: Pending Clarifications, which track unresolved issues requiring coordination with teammates; a Long-Term Plan, which outlines the global sequence of steps needed for successful task completion; and an Immediate Focus, which identifies the agent’s current local priority. These layers ensure that the agent’s decisions remain consistent across different time horizons. Pending Clarifications reduce ambiguity in teamwork by explicitly marking information gaps, while Long-Term Plans help prevent duplicated efforts through clear division of labor. Immediate Focus links ongoing execution to higher-level goals, reducing the risk of short-sighted or conflicting proposals.

Communication Strategy This component frames the agent’s process of generating utterances by embedding them in the dialogue context and requiring explicit justification before transmission. It is

divided into four sub-elements: Dialogue Summary, which condenses prior utterances into a concise history to prevent redundancy; Message Draft, which formulates candidate utterances grounded in the cooperative plan; and explicit reasoning fields that capture the Reason to Send and Reason not to Send. This reasoning process is integrated with the subsequent action selection phase, ensuring that the necessity of communication is evaluated before a message is produced. By discouraging unnecessary utterances, the Communication Strategy reduces coordination overhead and ensures that each message serves a clear collaborative function.

Building on these components, the Mind-Map Agent maintains a structured representation that guides foresight and long-horizon collaboration. Each new observation or teammate utterance updates the relevant fields of the schema, which in turn informs both action selection and communication decisions. By integrating these structured reasoning layers, the agent gains the ability to plan over longer horizons compared to single-step reactive approaches [5], while reducing coordination overhead that often arises in meta-plan deliberations [8]. In contrast to goal-alignment methods that emphasize convergence on a shared overall objective [25], the Mind-Map Agent complements this perspective by highlighting explicit and distributed articulation of sub-goals within the shared mission. To maintain coherence, updates are constrained to a JSON schema and guided by in-context exemplars, which reduces hallucinations caused by ambiguous or off-topic utterances. The schema thus functions as an externalized working memory that can be inspected to reveal the agent’s reasoning process, thereby supporting transparency and interpretability in multi-agent and human–robot cooperation.

4 Experiments

4.1 Implementation and Baselines

Benchmarks. We evaluate our approach on two cooperative embodied AI benchmarks from the CoELA framework. The first is Communicative Watch-and-Help (C-WAH) [28, 29], which focuses on household tasks that require two agents to coordinate through natural language communication. The second is ThreeDWorld Multi-Agent Transport (TDW-MAT) [30], where agents operate in a physically realistic environment to transport objects using containers. While C-WAH emphasizes communication and joint planning, TDW-MAT evaluates efficiency in collaborative object delivery.

Baselines. We compare the proposed Mind-Map Agent against both classical planning-based agents and LLM-based agents. The classical baselines include a Monte Carlo Tree Search hierarchical planner (MHP) and a rule-based hierarchical planner (RHP), both of which execute pre-defined routines without language use. The single-agent MHP setting serves as a reference for evaluating efficiency, while multi-agent MHP and RHP provide stronger non-LLM baselines. On the LLM side, we adopt the CoELA agent architecture and vary the LLMs: GPT-4o-mini, GPT-OSS-120B, Qwen-7B, and Llama3.3-70B. For each models, we evaluate both the CoELA variant and our Mind-Map variant to assess the contribution of structured reasoning. In addition, we include a Random Agent baseline which without LLM and selects randomly among plan candidates proposed by the Plan Execution Module. The Random Agent provides a lower bound and highlights the net contribution of LLM reasoning.

Metrics. For C-WAH, we measure cooperative efficiency using Average Steps (AS), defined as the number of steps required to complete all target tasks, and communication efficiency using Average Communication Steps (ACS), defined as the average number of utterances per episode. Lower values in both metrics indicate higher efficiency. To ensure robust evaluation, we repeat each experiment three times and report the mean and standard deviation. For TDW-MAT, we use Transport Rate (TR), defined as the fraction of target objects successfully delivered within 3000 timesteps, where higher values indicate better performance.

4.2 Experimental Results

4.2.1 Communicative Watch-and-Help

The results on C-WAH, summarized in Table 1, indicate that the Mind-Map Agent tends to outperform both classical planners and other LLM-based agents in terms of task completion efficiency and communication cost. Beyond the quantitative outcomes, several qualitative observations can be drawn.

Table 1: C-WAH performance across different agents. All values are reported as mean \pm standard deviation across three runs. Lower is better (\downarrow).

Model Type		Classic Agent		GPT-4o-mini		GPT-OSS-120B	
		Single	Multi	CoELA	MindMap	CoELA	MindMap
Symbolic Obs	AS \downarrow	111	75	65.4 \pm 23	57.0 \pm 19	74.3 \pm 21	53.2 \pm 18
	ACS \downarrow	-	-	12.1 \pm 5.0	1.5 \pm 1.5	4.2 \pm 2.0	4.3 \pm 3.0
Visual Obs	AS \downarrow	141	103	113.9 \pm 37	98.1 \pm 18	125.2 \pm 40	94.0 \pm 27
	ACS \downarrow	-	-	15.5 \pm 6.2	2.1 \pm 1.7	3.9 \pm 2.0	5.5 \pm 4.1

Model Type		Random Agent		Qwen2.5 7B		Llama3.3 70B	
		Single	Multi	CoELA	MindMap	CoELA	MindMap
Symbolic Obs	AS \downarrow	148 \pm 33	91.9 \pm 21	90.4 \pm 27	72.6 \pm 21	71.2 \pm 19	53.1 \pm 17
	ACS \downarrow	-	0	12.1 \pm 3.3	9.8 \pm 4.5	4.1 \pm 2.0	0.9 \pm 1.0
Visual Obs	AS \downarrow	186 \pm 34	128 \pm 29	174.1	134.3 \pm 10	128.6 \pm 39	90.6 \pm 25
	ACS \downarrow	-	0	14.9 \pm 21	10.2 \pm 5.1	5.1 \pm 2.6	1.3 \pm 1.5

First, smaller LLMs such as Qwen-7B often perform less effectively than even the non-reasoning Random Agent. This suggests that unstructured or hallucinated dialogue may in some cases interfere with coordination, resulting in less efficient task execution. For instance, we observed that Qwen2.5-7B frequently produced conversational turns containing emojis or multi-turn self-dialogues, reflecting tendencies from its pretraining data. Such behaviors, while natural in casual dialogue, were less suitable for cooperative task execution in this benchmark.

Second, the Mind-Map representation appears to mitigate these issues by separating dialogue history from the factual task state. This structured reasoning helps smaller models generate more consistent plans and, in some cases, allows them to approach or surpass the performance of much larger models without structured support. Explicit schema guidance focuses the agent on relevant task elements, leading to plans that are more coherent and aligned with the environment. It is also notable that, across most models, the CoELA baseline shows performance levels comparable to the Random Agent, indicating that language reasoning alone does not necessarily improve efficiency without structured memory. In part, this may be due to the CoELA prompt design being optimized for GPT-4, which other open models did not always follow strictly, resulting in mismatched or incomplete action selection.

Finally, Mind-Map Agents generate fewer utterances in most cases while maintaining effective cooperation. For example, with stronger LLMs such as Llama3.3-70B, the Mind-Map Agent achieves near-silent coordination once successful inference of the partner’s intentions has been established. In these cases, the agent appears to capture teammates’ goals with greater reliability and, based on environmental observations alone, can often construct long-horizon plans while monitoring the state of other agents. This enables effective collaboration with minimal reliance on explicit communication. The results suggest that structured reasoning not only facilitates early alignment of plans but also supports adaptation to evolving team states. By reducing the need for repeated exchanges and redundant clarifications, the Mind-Map Agent lowers overall communication overhead while sustaining cooperative efficiency in complex multi-agent settings.

4.2.2 ThreeDWorld Multi-Agent Transport

The advantages of the Mind-Map approach also extend to the TDW-MAT transport challenge (Table 2). In this setting, Mind-Map agents generally achieve higher transport success rates than the baselines by dividing labor more effectively and making use of tools such as containers to reduce redundant trips. The structured reasoning provided by the Mind-Map allows agents to exploit environment dynamics in a more systematic manner. While CoELA agents often engage in single-object deliveries, Mind-Map Agents utilize containers to transport multiple objects at once,

Table 2: TDW-MAT results (TR: transport rate)

TR \downarrow	Classic Agent		GPT-4o-mini	
	Single	Multi	CoELA	MindMap
Food	52	76	65	86
Stuff	49	74	76	79
Total	50	75	71	82

leading to faster progress and more successful deliveries within the timestep limit. Beyond task division, the Mind-Map approach also supports efficiency when acting alone: agents may temporarily defer picking up objects to first locate a container, or engage in dialogue with teammates to coordinate container use. For instance, when a teammate is already carrying two objects, an agent can provide information about the location of a container in another room to facilitate cooperation. These behaviors indicate that providing agents with an explicit schema of intentions and strategies contributes to more stable and coordinated teamwork, suggesting that the Mind-Map approach can serve as a general mechanism for enhancing decentralized multi-agent collaboration.

4.3 Ablation study on in-context example of the Mind-Map component

To examine which elements of the mind map contribute most to performance, we conduct an ablation study on C-WAH using GPT-4o-mini model. The agent is prompted to generate different subsets of the mind map, and performance is evaluated using Average Steps (AS) and Average Communication Steps (ACS). The tested settings include *Inf* (Inferred Others’ Intentions) only, *Coop* (Cooperative Strategy) only, *Comm* (Communication Strategy) only, and their combinations (Inf+Coop, Inf+Comm, Coop+Comm, and Inf+Coop+Comm). For this study, the key-matching process was deliberately omitted, allowing the agent to freely generate components without explicit structural constraints.

Table 3: Ablation study on different subsets of the Mind-Map component in C-WAH using GPT-4o-mini. Average Steps (AS), Average Communication Steps (ACS)

Mind-Map Example	AS ↓	ACS ↓
Inf	66.6	6.0
Coop	59.6	6.5
Comm	51.5	1.6
Inf + Coop	49.4	4.8
Inf + Comm	60.2	2.9
Coop + Comm	62.9	3.3
Inf + Coop + Comm	55.6	2.4

First, intention inference alone appears insufficient for effective collaboration. When relying only on inferred intentions (Inf), the agent exhibits higher step counts and substantial dialogue compared to the full Mind-Map agent. This suggests that predicting partner intentions without cooperative or communication planning does not provide enough structure for efficient execution.

Second, combining Cooperative Strategy and Communication Strategy (Coop+Comm) without explicit intention inference results in frequent overlapping behaviors. In this case, both agents often produced similar sub-plans independently, indicating that intention modeling is essential for task division and avoiding redundant effort.

Third, the Inf+Coop combination achieves the best numerical performance, with lower step counts than other subsets. However, in this setting, the agent did not generate the intended *Message Draft* within the mind map, instead defaulting to CoELA’s message generation. This observation suggests that refining the prompting process for Communication Strategy, perhaps with multi-phase message drafting, may further improve coordination when all components are included.

Finally, in the Comm-only setting, the agent displayed emergent flexibility by autonomously introducing additional elements beyond the predefined structure. These included tentative action plans, collaboration tactics, and shared goals, indicating that the agent adapts the reasoning structure when provided with minimal guidance. Such behavior highlights the potential value of allowing flexible mind map generation, where the model supplements missing components to support effective cooperation.

Overall, the ablation results suggest that diverse in-context examples of structured reasoning help agents organize their decision-making. Each component plays a distinct role: intention inference reduces redundancy, cooperative strategy supports task division, and communication strategy streamlines dialogue. When used together, they provide a balanced framework for adaptive and efficient multi-agent teamwork. Moreover, this result reveals that agents can correct and update their reasoning across steps, effectively managing their own long-term memory to support sustained cooperation.

5 Limitations and Discussion

5.1 Limitations

While the experimental results demonstrate clear advantages of the Mind-Map Agent, several limitations remain. First, our evaluation is restricted to two cooperative benchmarks (C-WAH and TDW-MAT), leaving open questions about the generality of the approach to other domains and real-world tasks. Second, smaller LLMs such as Qwen2.5-7B continue to exhibit unstable behavior, including hallucinated or noisy utterances, which the Mind-Map structure cannot fully mitigate. Third, it remains challenging to assess whether current LLMs possess a sufficiently robust grounding of commonsense knowledge. For instance, in certain cases the agent may fail to detect implausible object placements (e.g., a fork inside a refrigerator), leading to confusion and delayed task execution. Fourth, the effectiveness of Mind-Map reasoning is sensitive to prompt design and schema exemplars, raising concerns about reproducibility across different models and environments. Finally, while the agent demonstrates long-horizon planning in simulation, it remains uncertain whether these strategies will transfer reliably to open-ended settings involving human partners and noisy multimodal inputs.

5.2 Discussion

Despite these limitations, several insights emerge from our study. Mind-Map reasoning substantially mitigates confabulation by explicitly separating dialogue, factual state, and cooperative strategies, thereby improving consistency in multi-agent communication. Notably, emergent behaviors were observed in the Comm-only setting, where the agent autonomously expanded the reasoning structure to include missing elements such as collaboration tactics or shared goals. This suggests that schema-guided yet flexible reasoning allows agents to adaptively compensate when the provided structure is incomplete.

Another promising aspect lies in the transparency of the Mind-Map. By externalizing hidden goals and strategies into structured memory, agents produce interpretable representations that can potentially benefit human–AI collaboration. However, an important trade-off emerges: reducing explicit communication may increase efficiency but simultaneously reduce transparency for human teammates. Future work should therefore explore how to balance communication efficiency with interpretability, particularly in collaborative settings that require mutual trust and shared situational awareness.

6 Conclusion

In this work, we presented the Mind-Map Agent, a framework for structured reasoning in multi-agent cooperation. By disentangling dialogue, task state, and cooperative strategies, our approach mitigates confabulation, preserves long-term consistency, and improves interpretability. Extensive evaluation across multiple LLMs and environments showed that Mind-Map reasoning provides stable gains over both classical planners and unstructured LLM-based agents, highlighting the complementary value of intention inference, cooperative strategy, and communication strategy.

Beyond quantitative improvements, the findings suggest a broader implication: structured reasoning and memory can reconcile the trade-off between rigid protocol-driven planning and unconstrained dialogue-based coordination. In doing so, the Mind-Map Agent offers a generalizable mechanism for reliable multi-agent interaction, supporting both machine–machine and human–AI collaboration.

While open challenges remain in grounding perception, scaling to lightweight models, and adapting to open-ended human interaction, the proposed framework contributes toward the development of multi-agent systems that are not only more robust and efficient but also more transparent and better aligned with real-world cooperative demands.

Acknowledgment

This work was supported by the Electronics and Telecommunications Research Institute (ETRI) grant funded by the Ministry of Trade, Industry and Energy (MOTIE) in the year 2025 Robot Industrial Technology Project “Development of companion robot technologies capable of emotional

connection based on Human-Robot physical and cognitive interaction.” (No. 20018513, 20%), and by the Institute of Information & Communications Technology Planning & Evaluation (IITP) grants funded by the Korea government (MSIT) (RS-2024-00336738 Development of Complex Task Planning Technologies for Autonomous Agents, 20%, No. RS-2022-II220951 Development of Uncertainty-Aware Agents Learning by Asking Questions, 20%), by the National Research Council of Science & Technology (NST) grant by the Korea government (MSIT) (No. GTL25041-000, 20%) and by the Electronics and Telecommunications Research Institute (ETRI) (24ZB1200 Research of Human-centered Autonomous Intelligence System Original Technology, 20%).

References

- [1] Anthony Brohan, Yevgen Chebotar, Chelsea Finn, Karol Hausman, Alexander Herzog, Daniel Ho, Julian Ibarz, Alex Irpan, Eric Jang, Ryan Julian, et al. Do as i can, not as i say: Grounding language in robotic affordances. In *Conference on robot learning*, pages 287–318. PMLR, 2023.
- [2] Unnat Jain, Luca Weihs, Eric Kolve, Ali Farhadi, Svetlana Lazebnik, Aniruddha Kembhavi, and Alexander Schwing. A cordial sync: Going beyond marginal policies for multi-agent embodied tasks. In *ECCV*, 2020.
- [3] Peter Stone and Manuela Veloso. Multiagent systems: A survey from a machine learning perspective. *Autonomous Robots*, 8(3):345–383, 2000.
- [4] Matthew Chang, Gunjan Chhablani, Alexander Clegg, Mikael Dallaire Cote, Ruta Desai, Michal Hlavac, Vladimir Karashchuk, Jacob Krantz, Roozbeh Mottaghi, Priyam Parashar, et al. Partnr: A benchmark for planning and reasoning in embodied multi-agent tasks. In *Proceedings of the International Conference on Learning Representations (ICLR)*, 2025.
- [5] Hongxin Zhang, Weihua Du, Jiaming Shan, Qinhong Zhou, Yilun Du, Joshua B Tenenbaum, Tianmin Shu, and Chuang Gan. Building cooperative embodied agents modularly with large language models. In *Proceedings of the International Conference on Learning Representations (ICLR)*, 2024.
- [6] Yanming Wan, Jiayuan Mao, and Josh Tenenbaum. Handmethat: Human-robot communication in physical and social environments. *Advances in neural information processing systems*, 35:12014–12026, 2022.
- [7] Zhao Mandi, Shreeya Jain, and Shuran Song. Roco: Dialectic multi-robot collaboration with large language models. In *2024 IEEE International Conference on Robotics and Automation (ICRA)*, pages 286–299. IEEE, 2024.
- [8] Jie Liu, Pan Zhou, Yingjun Du, Ah-Hwee Tan, Cees GM Snoek, Jan-Jakob Sonke, and Efstratios Gavves. Capo: Cooperative plan optimization for efficient embodied multi-agent cooperation. In *Proceedings of the International Conference on Learning Representations (ICLR)*, 2025.
- [9] SeungWon Seo, SeongRae Noh, Junhyeok Lee, SooBin Lim, Won Hee Lee, and HyeongYeop Kang. Reveca: Adaptive planning and trajectory-based validation in cooperative language agents using information relevance and relative proximity. In *Proceedings of the AAAI Conference on Artificial Intelligence*, volume 39, pages 23295–23303, 2025.
- [10] Xinran Li, Chenjia Bai, Zijian Li, Jiakun Zheng, Ting Xiao, and Jun Zhang. Learn as individuals, evolve as a team: Multi-agent llms adaptation in embodied environments. *arXiv preprint arXiv:2506.07232*, 2025.
- [11] Andrew L Smith, Felix Greaves, and Trishan Panch. Hallucination or confabulation? neuroanatomy as metaphor in large language models. *PLOS Digital Health*, 2(11):e0000388, 2023.
- [12] Peiqi Sui, Eamon Duede, Sophie Wu, and Richard Jean So. Confabulation: The surprising value of large language model hallucinations. *arXiv preprint arXiv:2406.04175*, 2024.
- [13] Jason Wei, Xuezhi Wang, Dale Schuurmans, Maarten Bosma, Fei Xia, Ed Chi, Quoc V Le, Denny Zhou, et al. Chain-of-thought prompting elicits reasoning in large language models. *Advances in neural information processing systems*, 35:24824–24837, 2022.

- [14] Xudong Guo, Kaixuan Huang, Jiale Liu, Wenhui Fan, Natalia Vélez, Qingyun Wu, Huazheng Wang, Thomas L Griffiths, and Mengdi Wang. Embodied llm agents learn to cooperate in organized teams. *arXiv preprint arXiv:2403.12482*, 2024.
- [15] Ryan Lowe, Yi I Wu, Aviv Tamar, Jean Harb, OpenAI Pieter Abbeel, and Igor Mordatch. Multi-agent actor-critic for mixed cooperative-competitive environments. *Advances in neural information processing systems*, 30, 2017.
- [16] Tianmin Shu and Yuandong Tian. M³ rl: Mind-aware multi-agent management reinforcement learning. In *Proceedings of the International Conference on Learning Representations (ICLR)*, 2019.
- [17] Shivansh Patel, Saim Wani, Unnat Jain, Alexander G Schwing, Svetlana Lazebnik, Manolis Savva, and Angel X Chang. Interpretation of emergent communication in heterogeneous collaborative embodied agents. In *Proceedings of the IEEE/CVF International Conference on Computer Vision*, pages 15953–15963, 2021.
- [18] Unnat Jain, Luca Weihs, Eric Kolve, Mohammad Rastegari, Svetlana Lazebnik, Ali Farhadi, Alexander G Schwing, and Aniruddha Kembhavi. Two body problem: Collaborative visual task completion. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, 2019.
- [19] Matthijs TJ Spaan, Geoffrey J Gordon, and Nikos Vlassis. Decentralized planning under uncertainty for teams of communicating agents. In *Proceedings of the fifth international joint conference on Autonomous agents and multiagent systems (AAMAS)*, 2006.
- [20] Daniel S Bernstein, Robert Givan, Neil Immerman, and Shlomo Zilberstein. The complexity of decentralized control of markov decision processes. *Mathematics of operations research*, 27(4):819–840, 2002.
- [21] Claudia V Goldman and Shlomo Zilberstein. Optimizing information exchange in cooperative multi-agent systems. In *Proceedings of the second international joint conference on Autonomous agents and multiagent systems*, pages 137–144, 2003.
- [22] John E Laird. Introduction to soar. *arXiv preprint arXiv:2205.03854*, 2022.
- [23] Terence X Lim, Sidney Tio, and Desmond C Ong. Improving multi-agent cooperation using theory of mind. *arXiv preprint arXiv:2007.15703*, 2020.
- [24] Sarah A Wu, Rose E Wang, James A Evans, Joshua B Tenenbaum, David C Parkes, and Max Kleiman-Weiner. Too many cooks: Coordinating multi-agent collaboration through inverse planning. In *Proceedings of the annual meeting of the cognitive science society*, volume 42, 2020.
- [25] Lance Ying, Kunal Jha, Shivam Aarya, Joshua B Tenenbaum, Antonio Torralba, and Tianmin Shu. Goma: Proactive embodied cooperative communication via goal-oriented mental alignment. In *2024 IEEE/RSJ International Conference on Intelligent Robots and Systems (IROS)*, pages 7099–7106. IEEE, 2024.
- [26] Yao Xu, Shizhu He, Jiabei Chen, Zihao Wang, Yangqiu Song, Hanghang Tong, Guang Liu, Kang Liu, and Jun Zhao. Generate-on-graph: Treat llm as both agent and kg in incomplete knowledge graph question answering. *arXiv preprint arXiv:2404.14741*, 2024.
- [27] Zhenhailong Wang, Shaoguang Mao, Wenshan Wu, Tao Ge, Furu Wei, and Heng Ji. Unleashing the emergent cognitive synergy in large language models: A task-solving agent through multi-persona self-collaboration. *arXiv preprint arXiv:2307.05300*, 2023.
- [28] Xavier Puig, Tianmin Shu, Shuang Li, Zilin Wang, Yuan-Hong Liao, Joshua B Tenenbaum, Sanja Fidler, and Antonio Torralba. Watch-and-help: A challenge for social perception and human-ai collaboration. *arXiv preprint arXiv:2010.09890*, 2020.
- [29] Xavier Puig, Kevin Ra, Marko Boben, Jiaman Li, Tingwu Wang, Sanja Fidler, and Antonio Torralba. Virtualhome: Simulating household activities via programs. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pages 8494–8502, 2018.

- [30] Chuang Gan, Siyuan Zhou, Jeremy Schwartz, Seth Alter, Abhishek Bhandwalder, Dan Gutfreund, Daniel LK Yamins, James J DiCarlo, Josh McDermott, Antonio Torralba, et al. The threedworld transport challenge: A visually guided task-and-motion planning benchmark towards physically realistic embodied ai. In *2022 International conference on robotics and automation (ICRA)*, pages 8847–8854. IEEE, 2022.