

TeamWeaver: Hybrid LLM and Optimization-based Planning with Transparent Constraints in Heterogeneous Multi-Robot Teams

Abstract—Multi-robot task allocation in heterogeneous teams remains challenging due to dynamic environments and semantic ambiguity in natural language instructions. We propose *TeamWeaver*, a hybrid framework that integrates large language models with mixed-integer quadratic programming for robust, interpretable, and real-time task allocation. Language models parse unstructured instructions into structured task graphs, while the optimizer ensures feasibility under robot capabilities, resource constraints, and spatial-temporal dependencies. A closed-loop mechanism enables online capability learning and replanning under perturbations. We evaluate the framework on three datasets—PARTNR-val-mini, Dynamic-TeamWeaver (83 tasks, 4 scenarios), and real-world deployments (5 tasks, 40+ trials)—and show that it consistently outperforms baselines, achieving up to 10.6% higher success rate, 6.2% better completion, and 18.9% fewer planning steps on PARTNR-val-mini, 49.3% relative improvement under ambiguous conditions, and above 90% task success in real-world tests with up to 18.6% gain under communication failure. Transparency metrics remain above 0.66 with fewer than 1.2 replanning cycles per task. These results validate *TeamWeaver* as a scalable, interpretable, and deployable solution for heterogeneous multi-robot coordination in complex, real-world settings. The supplementary video, dataset and code are available online¹.

I. INTRODUCTION

A. Background and Motivation

Multi-robot task allocation (MRTA) is a fundamental problem in enabling effective collaboration among robot teams [1]–[4]. It is particularly critical in dynamic and unstructured scenarios, where tasks may need to be reassigned due to environmental changes or fluctuating resources [5]. Recent years have witnessed a rapid diversification of robotic platforms, driven by advances in morphology design, embodied intelligence, and specialized operating capabilities [6]–[8]. These developments have led to the widespread formation of heterogeneous multi-robot systems (HMRS), in which robots with complementary sensing, mobility, and manipulation skills operate together [9]–[11]. Such heterogeneity greatly enhances the potential of robot teams in real-world environments, enabling them to tackle more diverse, specialized, and resource-intensive tasks than homogeneous systems. However, it also poses significant challenges for MRTA. Differences in robot capabilities, resource constraints, and feasible action spaces complicate the design of allocation strategies that remain efficient, adaptive, and reliable under dynamic and uncertain conditions [3], [11], [12]. As recent studies highlight, robust task allocation in HMRS continues to be a central yet unresolved challenge in multi-robot coordination [13].

Recent multi-robot applications increasingly operate in open environments, where tasks emerge from high-level human

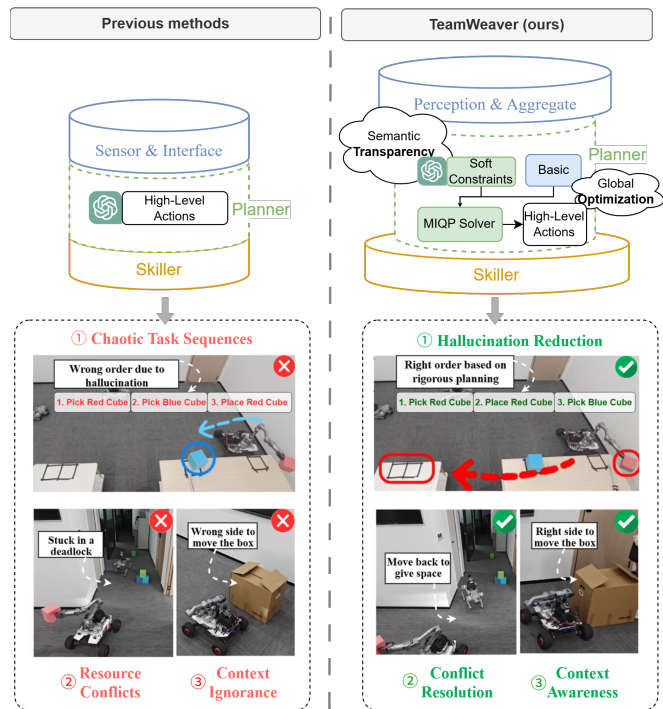


Fig. 1: Architectural comparison between traditional LLM Multi-Robot Planner (e.g., PARTNR [14]) and TeamWeaver. Left (PARTNR): The planner directly generates high-level actions from sensor inputs using LLM-based reasoning. Right (TeamWeaver): The planner employs a hybrid approach where LLMs provide semantic transparency and soft constraints, which are then processed by an MIQP solver along with hard constraints to produce globally optimized high-level actions. The key innovation is the separation of semantic reasoning (LLM) from constraint optimization (MIQP), enabling both interpretability and mathematical guarantees.

instructions or evolving conditions rather than from predefined templates. This shift requires MRTA to interpret open-ended goals and convert them into structured, executable tasks. Advances in large language models (LLMs) offer new opportunities to this problem by enabling the transformation of high-level instructions into representations suitable for downstream allocation [15]–[18]. However, when applied directly to reasoning, LLMs may introduce hallucinated task details or dependencies and may not consistently account for capability or resource constraints, leading to allocations that are linguistically coherent but not always guaranteed to be feasible or rigorous for execution in multi-robot systems [19]–[21].

¹<https://anonymous.4open.science/r/Dynamic-TeamWeaver-2766>

B. Proposed Approach

In this article, we propose a hybrid planning framework that integrates the semantic interpretability of LLMs with the constraint satisfaction guarantees of mathematical optimization. Our goal is not to replace LLM reasoning, but to address a specific mismatch that arises when natural-language task descriptions must be executed by heterogeneous robot teams: LLMs can extract high-level structures from open-ended instructions, but their outputs may contain implicit assumptions, missing dependencies, or softened constraints that do not fully correspond to what the robots can execute. This discrepancy becomes especially pronounced in multi-robot settings, where feasibility hinges on capability matching, resource limits, spatial-temporal consistency, and coordinated task sequencing.

To resolve this gap, we treat the language model as a semantic perception layer that converts unstructured instructions into intermediate world graphs explicitly capturing task primitives, their ordering constraints, and capability requirements. These representations are then grounded through a mathematically rigorous planner, instantiated in our system as a mixed-integer quadratic programming (MIQP) module. MIQP reformulates the LM-produced structure into a constraint-aware allocation problem, ensuring that the final plan satisfies resource feasibility, temporal consistency, and global coordination across heterogeneous robots. We term this integrated algorithm *LMs-MIQP* and build the *TeamWeaver* multi-robot system on top of it. Through this two-level design, the framework maintains the flexibility needed to interpret open-ended instructions while guaranteeing the execution-level correctness required for reliable multi-robot operation.

We evaluate the *LMs-MIQP* algorithm and the full *TeamWeaver* system in both simulated and real-world multi-robot experiments, spanning dynamic home-service tasks and warehouse-like environments. Experimental results show that our system consistently outperforms baseline methods, including pure LM-based and optimization-only approaches. Across five representative real-world tasks, *TeamWeaver* achieves task success rates ranging from $90.5 \pm 4.1\%$ to $98.3 \pm 5.1\%$, with improvements of up to 18.6% over PARTNR [14] in communication failure scenarios. Notably, *TeamWeaver* reduces replanning counts by 46–65% compared with PARTNR and achieves significantly higher Semantic Coherence via Temporal Logic (S_{TL}) values ranging from 0.66 ± 0.05 to 0.78 ± 0.04 , compared to PARTNR’s 0.10 ± 0.04 to 0.18 ± 0.05 , representing approximately 4–6 times improvement. As illustrated in Fig. 1, *TeamWeaver* reduces hallucinations, resolves conflicts, and enhances context awareness in real-world deployment compared with PARTNR.

C. Contributions

Our main contributions are threefold:

- **Design of the LMs-MIQP algorithm:** We develop LMs-MIQP, a hybrid planning algorithm that combines the semantic interpretation capabilities of large language models (LLMs) with the constraint satisfaction guarantees of mixed-integer quadratic programming (MIQP).

The algorithm grounds natural-language instructions into feasible, globally consistent multi-robot task allocations. In dynamic simulated home-service and warehouse-like environments, LMs-MIQP improves task success rate by 11.0%, percent complete by 8.0%, and reduces simulation steps by 18.9% compared with PARTNR [14].

- **Implementation of the TeamWeaver multi-robot system:** We integrate the proposed LMs-MIQP algorithm into *TeamWeaver*, a deployable heterogeneous robotic system capable of executing instruction-driven tasks under dynamic and partially observable conditions. The system features a unified semantic optimization pipeline, real-time coordination mechanisms, and interfaces for cross-platform deployment. In real-world experiments, *TeamWeaver* outperforms PARTNR across three representative scenarios, achieving improvements in task success rate (up to 4.9%), percent complete (up to 4.6%), and factual accuracy (up to 0.08), demonstrating robust coordination and reliable execution in everyday environments.
- **Comprehensive evaluation framework with transparency and interpretability metrics:** We adapt semantic entropy methods—originally developed for assessing uncertainty in LLM reasoning—to multi-robot system interpretability evaluation. This represents the first application of semantic-level uncertainty quantification to transparency assessment in heterogeneous multi-robot coordination. We establish a hierarchical evaluation framework that systematically assesses multi-robot systems across task execution, allocation efficiency, and interpretability dimensions. The framework introduces three transparency metrics: Semantic Coherence via Temporal Logic (S_{TL}), Semantic Consistency via Temporal Execution (S_{TE}), and Factual Accuracy (FA). S_{TL} and S_{TE} quantify the alignment between LLM-generated actions and optimization-based decisions, while FA measures hallucination reduction. These metrics, combined with traditional performance indicators (Success Rate, Percent Complete, Task Stability), provide a comprehensive assessment tool for evaluating hybrid LLM-optimization systems in dynamic multi-robot collaboration scenarios.

II. RELATED WORKS

A. Practices and Challenges of Heterogeneous Multi-Robot Systems

Heterogeneous multi-robot systems have found widespread applications across diverse domains, including warehouse automation [10], [22], search and rescue operations [9], [23], [24], domestic service [5], [25], [26], collaborative manufacturing [11], [27], [28], and physical human-robot collaborative transportation and manipulation [29]. These systems leverage complementary capabilities of heterogeneous robots to accomplish complex multi-robot collaborative tasks that would be infeasible for individual robots or homogeneous teams. As illustrated in Fig. 2, such collaborative scenarios typically involve multiple interconnected agents (mobile platforms, manipulation arms, and humanoid companions) navigating

1 and interacting within complex environments, requiring co-
2 ordinated task allocation and execution across heterogeneous
3 capabilities. The fundamental problem in these applications is
4 MRTA [1], [2], [4], which must account for diverse robot capa-
5 bilities, varying resource constraints, complex spatiotemporal
6 relationships, and dynamic environmental conditions.

7 Real-world multi-robot applications operate in dynamic and
8 uncertain environments, where multiple challenges emerge [2],
9 [4]: task requirements and environmental conditions change
10 unpredictably, requiring frequent replanning [9], [30], [31];
11 robot states fluctuate due to battery depletion, mechanical fail-
12 ures, or performance degradation, with tight coupling causing
13 cascading effects across the system [9], [30]; and open-world
14 scenarios involve natural language instructions with implicit
15 requirements difficult to formalize using traditional constraint-
16 based approaches [3]. These challenges create a fundamental
17 tension between rigorous constraint satisfaction and flexi-
18 bility. Traditional optimization-based approaches, formulated
19 as Mixed-Integer Programming problems [1], [9], excel at
20 constraint satisfaction when constraints are known and static,
21 but struggle with open-world instructions containing implicit
22 requirements [3]. This limitation has motivated alternative
23 paradigms, particularly LLMs for semantic understanding,
24 though LLM-based approaches introduce challenges related
25 to reliability and constraint verification [15], leading to hybrid
26 methods that combine the strengths of both paradigms.

27 *B. LLM-Driven Planning and Task Allocation for Multi Robot* 28 *Systems*

29 Recent surveys have shown that LLMs can support robot
30 task planning by translating high-level instructions into ex-
31 ecutable actions [32], [33]. Representative methods can be
32 categorized into dialog-driven approaches and modular frame-
33 works. Dialog-driven approaches employ interactive conver-
34 sations for task understanding: RoCo [13] equips each
35 robot with an LLM proxy for decentralized decision-making
36 through interactive dialogs; EMOS [30] adopts a “centralized
37 discussion–distributed execution” framework; SMART-LLM
38 [34] employs conversational feedback for task decomposition.
39 Modular frameworks integrate LLMs within structured sys-
40 tem architectures: CoELA [35] leverages LLMs as semantic
41 controllers coordinating perception, planning, and execution
42 modules; PARTNR [14], built on Habitat-Sim, validates col-
43 laborative systems in complex scenarios and contributes large-
44 scale datasets.

45 However, LLMs may fail to perform mathematical con-
46 straint verification on their output plans, leading to “hallucina-
47 tory allocations” [15], [36] where seemingly reasonable plans
48 are actually infeasible. These methods suffer from limitations
49 including poor adaptability to dynamic environments, high
50 latency, lack of large-scale validation, and hallucination issues
51 [15], [36]. The limitations of pure LLM-based approaches
52 have motivated hybrid methods that combine LLM semantic
53 understanding with formal optimization, positioning LLMs as
54 semantic perception modules that generate soft constraints
55 rather than directly making allocation decisions.

C. Hybrid Multi-Robot Planning on LLMs with Optimization *Methods*

3 The combination of LLM semantic understanding and op-
4 timizer rigor has emerged as a promising solution. Recent
5 studies have combined LLMs with formal planning methods to
6 improve task allocation efficiency in HMRS [37], [38]. LiP-
7 LLM [39] generates skill lists and task dependency graphs
8 using LLMs and then applies linear programming (LP) for
9 efficient allocation. LaMMA-P [40] combines LLMs with
10 the Planning Domain Definition Language (PDDL) for long-
11 horizon task planning and allocation. These methods mitigate
12 the black-box nature of LLMs and improve execution flexibil-
13 ity.

14 However, existing hybrid methods reveal several limitations.
15 First, regarding optimization expressiveness, many methods
16 employ LP as the optimization backend [39], but LP’s ex-
17 pressiveness may be limited in handling integer decisions,
18 mutual exclusion constraints, and nonlinear spatiotemporal re-
19 lationships inherent in heterogeneous robot allocation. Mixed-
20 Integer Programming formulations offer greater expressive-
21 ness, but their adoption in hybrid LLM-optimizer frameworks
22 remains limited. Second, many hybrid schemes adopt an open-
23 loop design, using the LLM as a one-time pre-processor
24 [39], [40], limiting the system’s ability to re-reason and
25 make closed-loop adjustments based on robot state feedback.
26 Closed-loop architectures have been explored [30], [35], but
27 their integration with LLM-based semantic reasoning in multi-
28 robot task allocation remains an open challenge. Third, re-
29 garding interpretability, the decision-making process is often
30 split between the LLM and the optimizer, with limited in-
31 tegration. In dynamic and unstructured environments, these
32 methods struggle to accurately model multi-agent constraints
33 and provide sufficient interpretability [4], [13], [36], limiting
34 end-to-end transparency from human intent to final action.

35 In contrast, our LMs-MIQP framework achieves tighter
36 integration between semantic reasoning and mathematical op-
37 timization. Unlike methods using LP, we employ MIQP to
38 handle integer decisions, mutual exclusion constraints, and
39 nonlinear relationships. Unlike open-loop designs, our frame-
40 work implements a closed-loop architecture where the LLM
41 continuously interprets high-level instructions and environ-
42 mental feedback, enabling dynamic adaptation. Unlike designs
43 that split decision-making, our framework provides end-to-end
44 transparency by embedding LLM assessments as interpretable
45 soft constraints in the optimization objective, enabling both
46 open-world adaptability and rigorous constraint satisfaction.

III. PROBLEM FORMULATION

A. Notation

49 To facilitate readability and ensure consistency, we establish
50 a unified notation system for the multi-robot task allocation
51 problem and the LMs-MIQP algorithm. Throughout this pa-
52 per, we consistently use $i \in \{1, \dots, N\}$ to index agents
53 and $j \in \{1, \dots, M\}$ to index tasks, where N denotes the
54 number of heterogeneous agents and M denotes the number
55 of decomposed sub-tasks.

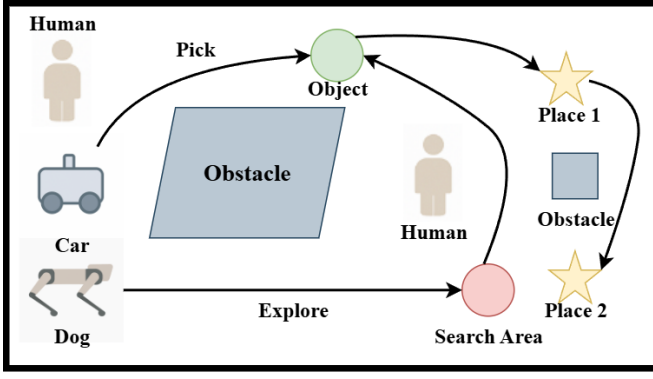


Fig. 2: Representative heterogeneous multi-robot collaborative scenario in Habitat-Sim, illustrating coordinated task execution among agents with complementary capabilities (mobile platforms, manipulation arms, and humanoid companions) in complex indoor environments.

Algorithm 1 Graph Fusion Operator

Require: Local graphs $\{\mathcal{G}_t^i = (\mathbf{V}_t^i, \mathbf{E}_t^i)\}_{i=1}^N$
Ensure: Global graph $\mathcal{G}_t^{\text{global}} = (\mathbf{V}_t^{\text{global}}, \mathbf{E}_t^{\text{global}})$

```

1: Initialize:  $\mathbf{V}_t^{\text{global}} \leftarrow \emptyset, \mathbf{E}_t^{\text{global}} \leftarrow \emptyset$ 
2: for each agent  $i \in \{1, \dots, N\}$  do
3:   for each node  $v \in \mathbf{V}_t^i$  do
4:     if  $v \notin \mathbf{V}_t^{\text{global}}$  then
5:       Add  $v$  to  $\mathbf{V}_t^{\text{global}}$ 
6:     else
7:       Update properties of existing node in  $\mathbf{V}_t^{\text{global}}$ 
8:     end if
9:   end for
10:  for each edge  $(u, v) \in \mathbf{E}_t^i$  do
11:    if  $(u, v) \notin \mathbf{E}_t^{\text{global}}$  then
12:      Add  $(u, v)$  to  $\mathbf{E}_t^{\text{global}}$ 
13:    end if
14:  end for
15: end for
16: Remove obsolete edges for dynamic entities
17: return  $\mathcal{G}_t^{\text{global}}$ 

```

dependencies) [42]. These local graphs capture both task-relevant information and dynamic environmental context from each agent’s perspective.

A fusion operator $\mathbf{Fusion}(\cdot)$ aggregates these local semantic graphs across all agents into a unified global world model:

$$\mathcal{G}_t^{\text{global}} = \mathbf{Fusion}(\{(\mathbf{V}_t^i, \mathbf{E}_t^i)\}_{i=1}^N) \quad (1)$$

The fusion operator implements a graph merge operation that constructs $\mathcal{G}_t^{\text{global}} = (\mathbf{V}_t^{\text{global}}, \mathbf{E}_t^{\text{global}})$ from local graphs $\{\mathcal{G}_t^i = (\mathbf{V}_t^i, \mathbf{E}_t^i)\}_{i=1}^N$. Algorithm 1 outlines the computational procedure, where v and u denote graph nodes (temporary variables within the algorithm):

This merge process ensures that the global graph $\mathcal{G}_t^{\text{global}}$ maintains a consistent, unified representation while preserving the temporal dynamics captured by individual agents’ partial observations. The fusion is performed incrementally as each agent’s local graph is updated, enabling real-time maintenance of the shared world model. This global world graph $\mathcal{G}_t^{\text{global}}$ serves as the core representation for multi-robot shared environmental understanding [14] and the backbone for downstream modules, including task allocation, high-level reasoning, and multi-agent coordination. It enables agents to maintain consistency across local decision-making processes while adapting to dynamic and partially observable environments.

C. MRTA Problem Definition

Building on the system representation and shared semantic graph defined in the previous section, we now formulate the multi-robot task allocation (MRTA) problem. Given a high-level natural language instruction \mathcal{I} , a heterogeneous agent team of N agents (denoted by indices $i \in \{1, \dots, N\}$), the goal is to assign these N agents to a set of M decomposed

1 The notation is organized into several categories: (1) *System*
2 *Parameters* define the fundamental dimensions and indices
3 of the problem; (2) *Inputs and System State* represent the
4 external inputs (natural language instructions, agent capabil-
5 ities) and the internal system representation (world graph);
6 (3) *Optimization Variables* are the decision variables in the
7 MIQP formulation; (4) *LLM and Cost Functions* capture the
8 semantic assessment and cost computation mechanisms; (5)
9 *Parameters and Constants* include weighting coefficients and
10 threshold values that tune the optimization behavior; and (6)
11 *Other Notation* covers auxiliary sets and derived quantities.

12 Table I provides a comprehensive summary of all key sym-
13 bols. When a symbol appears with subscripts or superscripts
14 (e.g., π_j for task j , \mathcal{C}_i for agent i), the indices follow the
15 convention established above. Matrix and vector dimensions
16 are explicitly stated where first introduced, and all optimization
17 variables are defined with their feasible domains.

B. System Model

19 We consider a team of N heterogeneous agents operating
20 in a shared workspace. Each agent $i \in \{1, \dots, N\}$ is charac-
21 terized by a capability profile $\mathcal{C}_i \in \mathbb{R}^D$ that defines the agent’s
22 proficiency across D capability dimensions (e.g., mobility, ma-
23 nipulation, control, liquid handling, power management) [41].
24 Each component $\mathcal{C}_i[k]$ quantifies agent i ’s capability level in
25 dimension k , where values range from 0 (incapable) to higher
26 positive values (highly proficient). These capability profiles
27 capture the diverse skills and resources available across the
28 heterogeneous team, enabling complementary task execution
29 through coordinated allocation. The capability profile \mathcal{C}_i is
30 used to determine task feasibility: task j can be assigned to
31 agent i only if $\mathcal{C}_i \geq \mathbf{r}_{\text{req}}(j)$, where $\mathbf{r}_{\text{req}}(j) \in \mathbb{R}^D$ specifies the
32 required capability vector for task j .

33 To support decision-making under uncertainty and partial
34 observability, each agent maintains a structured semantic belief
35 graph $\mathcal{G}_t^i = (\mathbf{V}_t^i, \mathbf{E}_t^i)$ at time t , where \mathbf{V}_t^i represents perceived
36 entities (objects, locations, other agents) and \mathbf{E}_t^i encodes
37 their pairwise relationships (spatial, temporal, or functional

TABLE I: Notation summary.

Symbol	Description	Symbol	Description
<i>System Parameters</i>			
N	Number of heterogeneous agents	i	Agent index, $i \in \{1, \dots, N\}$
M	Number of decomposed sub-tasks	j	Task index, $j \in \{1, \dots, M\}$
D	Dimension of the capability space	k	Capability dimension index, $k \in \{1, \dots, D\}$
t	Time step		
<i>Inputs and System State</i>			
\mathcal{I}	Natural language instruction for task goals	\mathcal{C}_i	Capability profile of agent i
$\mathcal{G}_t^i = (\mathbf{V}_t^i, \mathbf{E}_t^i)$	Local semantic belief graph of agent i at time t	\mathbf{V}_t^i	Set of perceived entities in agent i 's local graph
\mathbf{E}_t^i	Set of pairwise relationships in agent i 's local graph	$\mathcal{G}_t^{\text{global}}$	Global world graph encoding environmental context
$\mathbf{V}_t^{\text{global}}$	Set of entities in the global world graph	$\mathbf{E}_t^{\text{global}}$	Set of relationships in the global world graph
Fusion (\cdot)	Fusion operator aggregating local graphs	$\mathbf{r}_{\text{req}}(j)$	Required capability vector for task j
<i>Optimization Variables</i>			
$\mathbf{A} \in \{0, 1\}^{M \times N}$	Binary task assignment matrix	$\phi \in \{0, 1\}^M$	Task relaxation vector
$\pi_j \in \mathbb{R}^N$	Planning policy for task j	$\delta \in \mathbb{R}^M$	Load penalty vector
$\mathbf{A}_p \in \{0, 1\}^{M \times N}$	Previous allocation matrix		
<i>LLM and Cost Functions</i>			
$J(\mathbf{A}, \mathcal{G}_t^{\text{global}})$	Global cost function for MRTA	\mathcal{L}	LLM's contextual capability assessment
$\mathbf{c}(\mathbf{A}, \mathcal{L})$	Execution cost function	$\hat{\mathbf{c}}_j(\mathbf{A}, \mathcal{L})$	Estimated cost for task j
$\mathbf{q}_j \in \mathbb{R}^N$	Expected reward vector for task j	$\mathbf{Y}_j^* \in \mathbb{R}^N$	Policy-expected outcome vector for task j
$\mathbf{r}_j \in \mathbb{R}^N$	Allocation reward vector for task j	\mathbf{A}_j	j -th row of \mathbf{A}
<i>Parameters and Constants</i>			
$\mathbf{w}_j \in \mathbb{R}^M$	Weight vector for relaxation penalties	$\lambda_D, \lambda_T \geq 0$	Weighting coefficients for margin and transition cost
$\delta_{\text{max}} > 0$	Maximum allowed slack for load penalties	$\epsilon_j \geq 0$	Estimation bias for task j
$\mathbf{T}_{\text{trans}} \in \mathbb{R}^{M \times M}$	Transition cost matrix	$\mathbf{H}_{\text{perf}} \in \mathbb{R}^{M \times D}$	Performance history matrix
\mathcal{S}_i	Observed performance for agent i	$\hat{\mathcal{S}}_i$	Expected performance for agent i
$\gamma_1, \gamma_2 \geq 0$	Weighting coefficients for capability updates	$w_k \geq 0$	Weighting coefficient for capability dimension k
$\epsilon > 0$	Update threshold for capability adaptation		
<i>Other Notation</i>			
Ω	Feasible set (MRTA + MIQP constraints)	\mathcal{Y}_{gen}	Generated set from LLM outputs
\mathcal{Y}	Valid generated content, $\mathcal{Y} = \mathcal{Y}_{\text{gen}} \cap \Omega$		

1 sub-tasks (denoted by indices $j \in \{1, \dots, M\}$). These sub-
2 tasks are grounded from \mathcal{I} through parsing and graph-level
3 reasoning over the global world graph $\mathcal{G}_t^{\text{global}}$. Task allocation
4 decisions are conditioned on each agent's capability profile \mathcal{C}_i
5 and the global task context encoded in $\mathcal{G}_t^{\text{global}}$.

6 Formally, the MRTA problem is cast as a combinatorial opti-
7 mization over a binary assignment matrix $\mathbf{A} \in \{0, 1\}^{M \times N}$,
8 where $A_{ji} = 1$ indicates that task j is assigned to agent
9 i (we use j to index tasks and i to index agents through-
10 out). The goal is to find an optimal assignment $\mathbf{A}^* =$
11 $\arg \min_{\mathbf{A}} J(\mathbf{A}, \mathcal{G}_t^{\text{global}})$, minimizing a global cost function J
12 that depends on the assignment matrix and the world graph,
13 subject to the following constraints:

- 14 • **Agent capacity:** each agent is assigned to at most one
15 task, $\sum_{j=1}^M A_{ji} \leq 1, \quad \forall i \in \{1, \dots, N\}$;
- 16 • **Task exclusivity:** each task is assigned to at most one
17 agent, $\sum_{i=1}^N A_{ji} \leq 1, \quad \forall j \in \{1, \dots, M\}$;
- 18 • **Capability feasibility:** task j is assigned to agent i only
19 if the agent's capability \mathcal{C}_i satisfies the task requirement
20 $\mathbf{r}_{\text{req}}(j) \in \mathbb{R}^D$, where D is the dimension of the capability

space, i.e., $\mathcal{C}_i \geq \mathbf{r}_{\text{req}}(j) \cdot A_{ji}, \quad \forall i \in \{1, \dots, N\}, j \in$
 $\{1, \dots, M\}$;

- **System constraints:** including kinodynamic limits, safety
zones, and collision avoidance [9], [43], [44].

Together, the above formulation defines the MRTA problem
as a constrained optimization over discrete task-agent assign-
ments. The input includes: (1) a natural language instruction
 \mathcal{I} that specifies the high-level task goals, (2) the capability
profiles \mathcal{C}_i for each of the N agents that define what each
agent can do, and (3) the global world graph $\mathcal{G}_t^{\text{global}}$ that
captures the current environmental context. The output is a
binary assignment matrix $\mathbf{A} \in \{0, 1\}^{M \times N}$ satisfying both
logical constraints (agent capacity, task exclusivity, capabil-
ity feasibility) and physical constraints (kinodynamic limits,
safety zones, collision avoidance). This formulation provides
a unified interface between high-level task semantics derived
from \mathcal{I} and low-level execution feasibility determined by
agent capabilities \mathcal{C}_i and environmental constraints encoded in
 $\mathcal{G}_t^{\text{global}}$. In the following section, we present the LMs-MIQP
algorithm to solve this problem efficiently in dynamic, real-

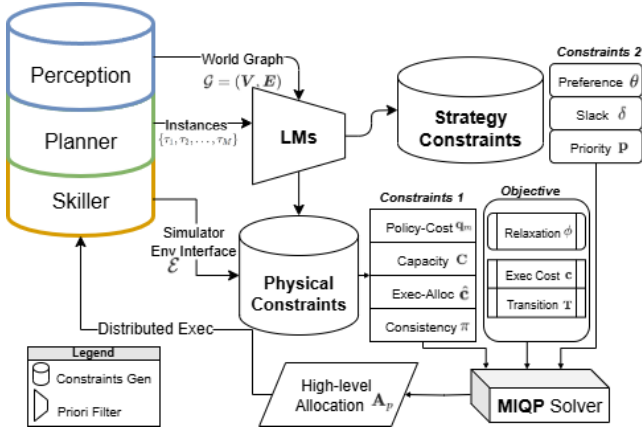


Fig. 3: The algorithm framework of LMs-MIQP. Language models generate strategy constraints, which are combined with physical constraints and solved by the MIQP optimizer. Notations are defined in Table I.

load penalty vector; (4) $\mathbf{A}_p \in \{0, 1\}^{M \times N}$ is the previous allocation matrix; and (5) \mathcal{L} is the LLM's contextual capability assessment derived from $\mathcal{G}_t^{\text{global}}$ and agent capability profiles C_i for $i \in \{1, \dots, N\}$.

The MIQP optimization determines the optimal allocation matrix \mathbf{A} , relaxation vector ϕ , and planning policies $\pi = [\pi_1, \pi_2, \dots, \pi_M]$ by solving:

$$\min_{\mathbf{A}, \phi, \pi} \underbrace{\mathbf{c}(\mathbf{A}, \mathcal{L})}_{\text{Execution Cost}} + \underbrace{\mathbf{w}_J^T \phi}_{\text{Relaxation}} + \underbrace{\lambda_D \|\delta\|_2^2}_{\text{Margin Penalty}} + \underbrace{\lambda_T \|\mathbf{T}_{\text{trans}}(\mathbf{A} - \mathbf{A}_p)\|_1}_{\text{Transition Cost}} \quad (2)$$

subject to:

$$\hat{c}_j(\mathbf{A}, \mathcal{L}) \leq \mathbf{q}_j^T \pi_j + \epsilon_j, \quad \forall j \in \{1, \dots, M\} \quad (a)$$

$$\delta_j \geq -\phi_j \delta_{\max}, \quad \forall j \in \{1, \dots, M\} \quad (b)$$

$$\pi_j^T \mathbf{Y}_j^* - \mathbf{r}_j^T \mathbf{A}_j = \delta_j, \quad \forall j \in \{1, \dots, M\} \quad (c)$$

$$\pi_j^T \mathbf{1} = 1, \quad \forall j \in \{1, \dots, M\} \quad (d)$$

The objective function in Eq. (2) consists of four terms: execution cost $\mathbf{c}(\mathbf{A}, \mathcal{L})$ modulated by LLM assessment \mathcal{L} , relaxation penalty $\mathbf{w}_J^T \phi$ with weights $\mathbf{w}_J \in \mathbb{R}^M$ informed by LLM interpretation of task criticality from \mathcal{I} , margin penalty $\lambda_D \|\delta\|_2^2$ for load balancing, and transition cost $\lambda_T \|\mathbf{T}_{\text{trans}}(\mathbf{A} - \mathbf{A}_p)\|_1$ that penalizes changes from previous allocation \mathbf{A}_p . The coefficients $\lambda_D, \lambda_T \geq 0$ control the relative importance of load balancing versus allocation stability.

Constraints (a)–(d) in Eq. (a)–(d) regulate the optimization for all tasks $j \in \{1, \dots, M\}$: constraint (a) bounds estimated cost $\hat{c}_j(\mathbf{A}, \mathcal{L})$ with expected reward $\mathbf{q}_j^T \pi_j$ plus estimation bias $\epsilon_j \geq 0$; constraint (b) allows load penalties $\delta_j \in \delta$ to be negative (up to $-\delta_{\max}$) when relaxation variable $\phi_j = 1$; constraint (c) equates the difference between policy-expected outcomes \mathbf{Y}_j^* and allocation rewards $\mathbf{r}_j^T \mathbf{A}_j$ to δ_j ; and constraint (d) enforces normalization so π_j forms a valid probability distribution over the N agents. These constraints extend the basic MRTA constraints from Section III by incorporating LLM-derived semantic assessments, adaptive task relaxation, and probabilistic planning policies. Together with the basic MRTA constraints (agent capacity, task exclusivity, capability feasibility), constraints (a)–(d) define the feasible set Ω for the optimization.

The execution cost $\mathbf{c}(\mathbf{A}, \mathcal{L})$ aggregates task-specific costs for all agent-task pairs (i, j) where $A_{ji} = 1$. For each pair, the cost component c_{ji} combines negative task utility (computed from task-specific functions based on agent state and $\mathcal{G}_t^{\text{global}}$) and an aptitude penalty term inversely proportional to the LLM compatibility score \mathcal{L}_{ji} from \mathcal{L} . The total cost is $\mathbf{c}(\mathbf{A}, \mathcal{L}) = \sum_{i=1}^N \sum_{j=1}^M c_{ji} \cdot A_{ji}$, and the estimated cost for task j is $\hat{c}_j(\mathbf{A}, \mathcal{L}) = \sum_{i=1}^N c_{ji} \cdot A_{ji}$. The LLM generates compatibility scores \mathcal{L}_{ji} by analyzing \mathcal{I} and $\mathcal{G}_t^{\text{global}}$ to evaluate how well agent capabilities C_i match task requirements $\mathbf{r}_{\text{req}}(j)$ given the environmental context encoded in $\mathcal{G}_t^{\text{global}}$. These scores are derived from semantic analysis of task descriptions in \mathcal{I} , object locations and relationships in $\mathcal{G}_t^{\text{global}}$, and agent capability profiles C_i for $i \in \{1, \dots, N\}$.

1 time settings.

2 IV. METHODOLOGY

3 Building on the MRTA problem formulation in Section III, this section presents our methodology that combines LLM semantic reasoning with MIQP optimization. As illustrated in Fig. 3, LLMs process natural language instructions \mathcal{I} and the global world graph $\mathcal{G}_t^{\text{global}}$ to generate contextual assessments, encoded as soft constraints in MIQP. The MIQP solver enforces hard constraints (agent capacity, task exclusivity, capability feasibility) while optimizing LLM-informed preferences. This methodology consists of two components: (1) the LMs-MIQP algorithm (Section IV-A) that formulates task allocation as constrained optimization; and (2) a closed-loop adaptation mechanism (Section IV-B) for dynamic capability learning.

15 A. LMs-MIQP: Constrained Multi-Robot Task Allocation Algorithm

17 To solve the MRTA problem formulated in Section III, we develop the LMs-MIQP algorithm that formulates task allocation as a Mixed Integer Quadratic Programming (MIQP) problem [45]. The algorithm optimizes over the binary assignment matrix $\mathbf{A} \in \{0, 1\}^{M \times N}$ that maps M tasks to N agents. LLMs interpret natural language instructions \mathcal{I} and the global world graph $\mathcal{G}_t^{\text{global}}$ to generate contextual capability assessments \mathcal{L} , which are encoded as soft constraints and cost terms in the MIQP objective. The MIQP solver enforces hard constraints from Section III while optimizing LLM-informed preferences. This approach mitigates LLM hallucination by constraining optimization to the feasible set Ω .

29 The LMs-MIQP algorithm extends the basic MRTA formulation by introducing additional optimization variables: (1) $\phi \in \{0, 1\}^M$ is the task relaxation vector; (2) $\pi_j \in \mathbb{R}^N$ is the planning policy for task j ; (3) $\delta \in \mathbb{R}^M$ is the

¹For symbol definitions (cylinder, rectangle, arrow, etc.), see the anonymous reference at https://anonymous.4open.science/t/uml-notation-reference-5F74/UML_notation_reference.pdf.

The MIQP problem is solved using standard solvers that handle the mixed-integer quadratic structure [46]. The algorithm mitigates LLM hallucination by constraining optimization to the feasible set Ω , ensuring LLM outputs from \mathcal{I} are validated against hard constraints involving \mathcal{C}_i and $\mathcal{G}_t^{\text{global}}$.

B. Closed-Loop Adaptation for Dynamic Capability Learning

Building upon the MIQP foundation, TeamWeaver incorporates a closed-loop adaptation mechanism that enables continuous learning and parameter adjustment in response to dynamic environments and evolving agent capabilities. The mechanism operates through continuous monitoring of system performance and systematic updates to agent capability profiles \mathcal{C}_i based on execution feedback.

The adaptation process, formalized in Algorithm 2, begins when performance deviations exceed a predefined threshold. The system updates agent capability profiles \mathcal{C}_i by incorporating both immediate performance feedback and historical execution patterns. For agent i , the capability update is computed as:

$$\Delta \mathcal{C}_i = \gamma_1 \frac{|\mathcal{S}_i - \hat{\mathcal{S}}_i|}{|\hat{\mathcal{S}}_i|} + \gamma_2 \sum_{j=1}^M \sum_{k=1}^D \mathbf{H}_{\text{perf}}[j, k] \cdot \mathcal{C}_i[k] \cdot w_k \quad (3)$$

where \mathcal{S}_i and $\hat{\mathcal{S}}_i$ denote the observed and expected performance for agent i , $\mathbf{H}_{\text{perf}} \in \mathbb{R}^{M \times D}$ represents the performance history matrix encoding task execution quality, $w_k \geq 0$ are weighting coefficients for capability dimension $k \in \{1, \dots, D\}$, and $\gamma_1, \gamma_2 \geq 0$ control the relative importance of immediate feedback versus historical patterns. This formulation ensures that capability updates reflect both immediate deviations and historical performance trends, enabling the system to distinguish between temporary fluctuations and systematic capability changes.

The updated capability profiles \mathcal{C}_i subsequently inform task prioritization and specialization assessment. For each agent i , the specialization score for task j is computed as the alignment between the updated capability profile \mathcal{C}_i and the task requirement vector $\mathbf{r}_{\text{req}}(j)$, accounting for the performance history encoded in \mathbf{H}_{perf} . Task priorities are dynamically adjusted through a constrained optimization that maximizes alignment between agent specializations and task assignments, ensuring that high-priority tasks are allocated to agents with the best matching capabilities. This prioritization mechanism operates in conjunction with the MIQP allocation process: when capability updates exceed a threshold, the system triggers a replanning cycle that recomputes the assignment matrix \mathbf{A} using the updated \mathcal{C}_i values, creating a feedback loop between execution performance and allocation decisions.

The adaptation mechanism operates asynchronously and event-driven, monitoring performance continuously but triggering updates only when deviations exceed predefined thresholds. This design balances responsiveness to environmental changes with computational efficiency, avoiding unnecessary replanning cycles while ensuring the system adapts promptly to significant capability mismatches or environmental disturbances. The threshold parameter ϵ in Algorithm 2 controls the sensitivity of the adaptation mechanism: smaller values lead to

more frequent updates and higher responsiveness, while larger values reduce computational overhead and filter out transient fluctuations.

This integrated feedback loop creates a self-improving system that maintains robust task allocation performance by continuously learning from execution experience and adapting to evolving operational conditions. The mathematical rigor of the MIQP foundation combined with the adaptability of the closed-loop mechanism provides a principled approach to handling dynamic multi-robot collaboration scenarios. The two-level design (static optimization via MIQP and dynamic adaptation via feedback learning) enables TeamWeaver to maintain optimal performance in both well-modeled scenarios (where MIQP provides near-optimal solutions) and uncertain, evolving environments (where adaptation compensates for model inaccuracies and changing conditions).

Together, the LMs-MIQP algorithm and the closed-loop adaptation mechanism form a unified methodology for dynamic multi-robot task allocation. The LMs-MIQP component (Section IV-A) provides the core optimization framework that combines LLM semantic reasoning with mathematical constraint satisfaction, ensuring feasible and interpretable task assignments. The closed-loop adaptation component (Section IV-B) enables the system to learn from execution experience and adapt to changing conditions, updating agent capability profiles \mathcal{C}_i and task priorities based on performance feedback. This two-level design addresses both the static optimization challenge (through MIQP) and the dynamic adaptation challenge (through feedback learning), creating a robust system that maintains optimal performance while responding to environmental changes and evolving agent capabilities. The integration of these components is realized through the system architecture presented in Section V, which provides the implementation framework for deploying TeamWeaver in real-world scenarios.

V. SYSTEM ARCHITECTURE

To realize the LMs-MIQP algorithm and closed-loop adaptation mechanism presented in Section IV, TeamWeaver adopts a **Perception-Plan-Skill** architecture that enables closed-loop multi-robot execution. As illustrated in Fig. 4, this three-layer design decouples semantic understanding, optimization-based planning, and decentralized execution while maintaining tight integration through the global world graph $\mathcal{G}_t^{\text{global}}$ and feedback mechanisms.

A. Perception Layer

The perception layer serves as the foundation for multi-robot environmental understanding by fusing heterogeneous multimodal data from distributed sensors into a unified world graph $\mathcal{G}_t^{\text{global}}$. As shown in Fig. 4, the layer processes high-level commands (task goals) from human operators, such as “Rearrange toy cars and clean the kitchen table,” along with raw sensor data from multiple sources.

The layer operates through three main stages: *aggregation*, *fusion*, and *representation*. In the aggregation stage, each sensor module processes its raw observations and performs

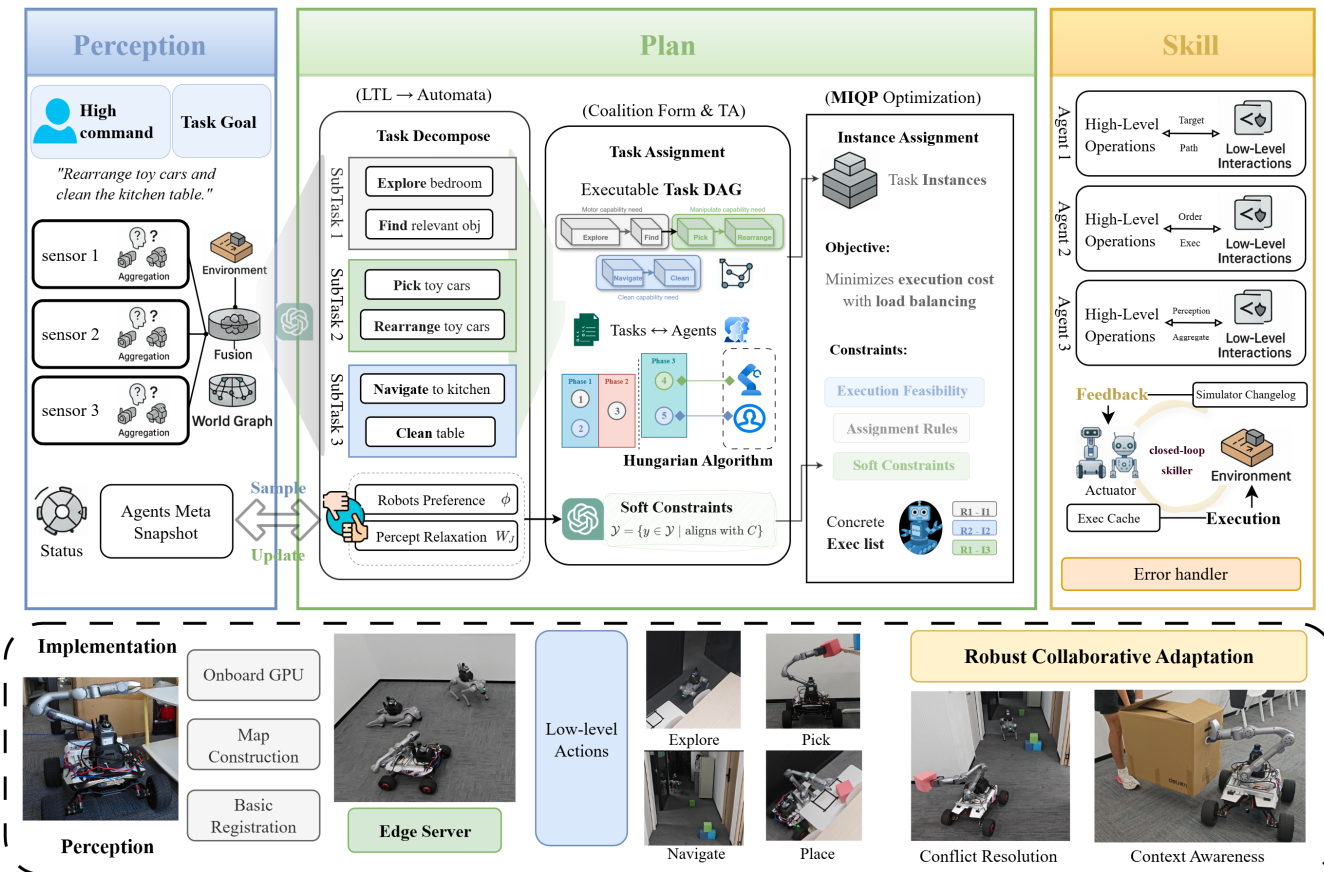


Fig. 4: Our TeamWeaver’s **Perception–Plan–Skill** system architecture for multi-robot closed-loop task execution. The framework consists of a **perception** layer that fuses heterogeneous sensing and task goals into a unified world graph, a **planner** layer that decomposes tasks and allocates them via MIQP optimization with LLM-derived soft constraints, and a **skill** layer that distributes primitives to decentralized agents with integrated state and performance feedback. A real-world **implementation** is demonstrated with onboard perception, edge-based planning, and collaborative execution in real-world deployment.

1 local aggregation to extract semantic entities and relationships.
 2 These aggregated observations are then fused using the fusion
 3 operator $\text{Fusion}(\cdot)$ (Algorithm 1) to construct the global
 4 world graph $\mathcal{G}_t^{\text{global}}$, which encodes the environment state
 5 including object locations, spatial relationships, and dynamic
 6 entities. Simultaneously, the layer maintains an *agents meta*
 7 *snapshot* that tracks the current status of all robots in the
 8 system, including their positions, capabilities, and operational
 9 states.

10 In real-world deployment, the *aggregation* stage runs *on-*
 11 *board* each robot, performing real-time local sensor processing
 12 to extract semantic entities and relationships. This distributed
 13 aggregation enables low-latency perception updates and re-
 14 duces communication bandwidth requirements. The *fusion*
 15 stage executes on the *edge workstation*, where locally aggre-
 16 gated observations from all robots are merged to construct the
 17 unified world graph $\mathcal{G}_t^{\text{global}}$. This hybrid design separates time-
 18 critical local processing from global coordination: onboard
 19 aggregation maintains real-time responsiveness, while cen-
 20 tralized fusion ensures consistent global state representation
 21 across the multi-robot team. The output of the perception layer
 22 consists of *sample* data that includes: (1) updates to the world
 23 graph $\mathcal{G}_t^{\text{global}}$ reflecting the latest environmental observations,

and (2) the agents meta snapshot providing current robot status
 information. This data is continuously fed to the plan layer,
 enabling real-time task planning based on the most up-to-date
 environmental context.

B. Plan Layer

The plan layer serves as the computational core that imple-
 ments the LMs-MIQP algorithm (Section IV-A) and closed-
 loop adaptation (Section IV-B). As illustrated in Fig. 4, the
 plan layer operates through three sequential stages: task de-
 composition, task assignment, and instance assignment.

a) *Task Decomposition (LTL \rightarrow Automata)*: The plan
 layer receives natural language instructions \mathcal{I} (task goals)
 and the world graph $\mathcal{G}_t^{\text{global}}$ from the perception layer. Linear
 temporal logic (LTL) specifications are compiled into automata
 representations, enabling systematic breakdown of high-level
 goals into executable primitives. For example, the task “Re-
 arrange toy cars and clean the kitchen table” is decomposed
 into subtasks such as “Explore bedroom,” “Find relevant obj,”
 “Pick toy cars,” “Rearrange toy cars,” “Navigate to kitchen,”
 and “Clean table.” These primitives are then grouped into
 task phases to support coalition formation and task alignment,
 forming the basis for multi-robot coordination.

b) *Task Assignment (Coalition Form & TA)*: The decomposed primitives are arranged into a directed acyclic graph (DAG) called the *executable task DAG*, whose edges encode precedence relationships and shared-resource constraints. The tasks \leftrightarrow agents module employs a Hungarian algorithm to facilitate initial task-agent matching, providing a preliminary assignment that considers basic compatibility between tasks and agent capabilities. This initial matching serves as a starting point for the subsequent optimization-based allocation.

c) *Instance Assignment (MIQP Optimization)*: The core optimization problem is solved using the LMs-MIQP algorithm (Section IV-A), which formulates task allocation as a mixed-integer quadratic program. The optimization minimizes the objective function in Eq. (2), subject to the constraints in Eq. (a)–(d), along with the basic MRTA constraints (agent capacity, task exclusivity, capability feasibility) from Section III. Soft constraints are derived from LLM assessments: $\mathcal{Y} = \mathcal{Y}_{\text{gen}} \cap \Omega$, where \mathcal{Y}_{gen} represents LLM-generated preferences and Ω is the feasible set, enabling hallucination mitigation through constraint intersection.

The plan layer receives *robots preference* information and *percept relaxation weights* \mathbf{w}_J from the LLM, which inform the soft constraints in the optimization. The output is a *concrete execution list* specifying task assignments for each robot (e.g., R1-T1, R2-T2, R1-T3), which is passed to the skill layer for execution.

C. Skill Layer

The skill layer grounds the optimized plans from the plan layer by instantiating high-level primitives into decentralized low-level actions executed by individual agents. As shown in Fig. 4, the layer consists of multiple agent modules (Agent 1, Agent 2, Agent 3), each responsible for executing tasks assigned by the plan layer.

Each agent module operates through a two-level abstraction: *high-level operations* (Target, Order, Perception) that receive commands from the plan layer, and *low-level interactions*

(Path, Exec, Perception) that interface with the physical environment through actuators. The *closed-loop skiller* coordinates the execution process, translating high-level primitives such as “navigate,” “pick,” or “clean” into sequences of low-level control commands. During execution, the layer maintains an *execution cache* that stores execution data, including action histories, performance metrics, and environmental feedback.

In real-world deployment, the skill layer runs *onboard* each robot, enabling decentralized execution and real-time control. The *closed-loop skiller* translates high-level commands into robot-specific low-level control sequences locally, adapting to each robot’s hardware capabilities and kinematic constraints. The *execution cache* maintains action histories and performance metrics onboard, providing low-latency access for local error handling and rapid adaptation to perturbations. The layer also includes an *error handler* that manages execution failures, adapting to perturbations and environmental changes to ensure robustness in dynamic collaborative environments. Execution status and performance feedback are transmitted to the edge workstation through the *simulator changelog* (in simulation) or *execution status messages* (in real-world deployment), which are published as ROS topics. The edge workstation receives this feedback and enables the closed-loop adaptation mechanism (Algorithm 2) in the plan layer to monitor performance, update agent capability profiles \mathcal{C}_i , and trigger replanning when significant deviations are detected. This distributed execution architecture ensures responsive local control while maintaining global coordination through centralized monitoring.

D. Real-World Deployment

TeamWeaver is deployed in a hybrid onboard–edge architecture that mirrors the structure shown in Fig. 4. In the real-robot setup, Zhongling and Go2-A are each equipped with an external onboard computer running Ubuntu 20.04 and ROS, interfaced with the robot base, arm controllers, and a Mid-360 LiDAR. These machines run SLAM-based mapping, localization, and navigation stacks, so Zhongling+Z1 and Go2-A provide stronger *explore* and *navigate* capabilities and can maintain richer local world models. In contrast, G1 and Go2-B rely on the vendor-provided onboard Ubuntu 20.04 systems and ROS navigation modules shipped with the Unitree platforms. On these robots, we enable speech recognition and lightweight waypoint navigation, which supports basic exploration and human interaction but with weaker long-horizon autonomy. This intentionally asymmetric configuration exposes non-uniform sensing and control capabilities across agents and allows TeamWeaver to be evaluated on whether it schedules robots according to their strengths.

All high-level reasoning runs on an edge workstation with an Intel Core i9 CPU, an NVIDIA RTX 4060 GPU, and 64 GB RAM. The LMs-MIQP planner and world-graph construction modules execute on this edge node, which connects to each robot’s Ubuntu host via *ssh* instruction and ROS. Robots publish perception results, pose estimates, and task-status messages as ROS topics, while the edge node sends back task allocations, navigation goals, and manipulation subgoals. In

Algorithm 2 Closed-Loop Capability Adaptation

Require: Capability profiles $\{\mathcal{C}_i\}_{i=1}^N$, history \mathbf{H}_{perf} , threshold $\epsilon > 0$

Ensure: Updated profiles $\{\mathcal{C}_i^{(t+1)}\}_{i=1}^N$

- 1: Compute allocation \mathbf{A} using MIQP
 - 2: Execute tasks and collect performance data
 - 3: **for** each agent $i \in \{1, \dots, N\}$ **do**
 - 4: Observe \mathcal{S}_i and compare with $\hat{\mathcal{S}}_i$
 - 5: **if** $|\mathcal{S}_i - \hat{\mathcal{S}}_i| \geq \epsilon$ **then**
 - 6: **Update** \mathcal{C}_i **using Eq. (3)**
 - 7: Recompute task priorities
 - 8: **end if**
 - 9: **end for**
 - 10: **if** significant capability changes detected **then**
 - 11: **Replan** \mathbf{A} **using updated** $\{\mathcal{C}_i\}_{i=1}^N$
 - 12: **end if**
 - 13: **return** $\{\mathcal{C}_i^{(t+1)}\}_{i=1}^N$
-

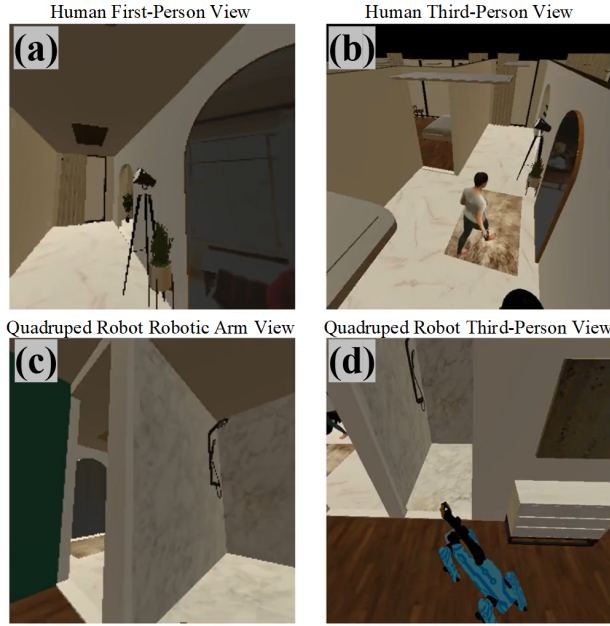


Fig. 5: Illustrative example of the Habitat-Sim environment used in our simulation study. The figure shows four representative viewpoints: (a) human first-person view, (b) human third-person view, (c) quadruped-mounted robotic arm view, and (d) quadruped robot third-person view.

1 this way, perception and skills remain physically distributed
 2 on the robots, whereas task planning and optimization are
 3 centralized but lightweight enough for closed-loop operation.

4 Communication between the robots and the edge worksta-
 5 tion is provided by a shared high-speed Wi-Fi network in
 6 the indoor environment. ROS is used as the common mid-
 7 dleware layer to handle message passing and synchronization.
 8 Although the wireless network introduces realistic latency and
 9 occasional jitter, the bidirectional bandwidth is sufficient for
 10 TeamWeaver to perform frequent replanning and semantic
 11 updates. By combining stronger compute and SLAM pipelines
 12 on Zhongling and Go2-A with lighter onboard stacks on G1
 13 and Go2-B over the same network, the real-world deployment
 14 directly tests whether the proposed architecture can maintain
 15 coherent multi-robot coordination under heterogeneous sens-
 16 ing, control, and communication conditions.

VI. EVALUATION METRICS AND PROTOCOLS

18 We adopt a hierarchical set of quantitative metrics to
 19 evaluate task execution, allocation stability, and system trans-
 20 parency, defined as follows:

21 *a) Task Execution Metrics: Success Rate (SR)* measures
 22 the ratio of successfully completed tasks to the total number
 23 of assigned tasks: $SR = \frac{N_{\text{success}}}{N_{\text{total}}}$, where N_{success} and N_{total}
 24 denote the number of successfully completed and total tasks,
 25 respectively. *Percent Complete (PC)* quantifies the average
 26 completion ratio of all sub-tasks within an episode: $PC =$
 27 $\frac{1}{N_{\text{task}}} \sum_{i=1}^{N_{\text{task}}} \frac{n_i^{\text{done}}}{n_i^{\text{total}}}$, where n_i^{done} and n_i^{total} represent the number

of completed and total steps for sub-task i . *Simulation Steps*
 (SS) is the total number of simulation timesteps required to
 complete an episode, reflecting overall execution efficiency.
Planning Cycles (PCy) counts the number of global re-
 planning events triggered due to environmental changes or
 violated constraints, thus reflecting system adaptability.

7 *b) Transparency and Interpretability Metrics:* To further
 8 evaluate the planner’s planning fidelity, we refer to prior LLM
 9 evaluation approaches [47] that assess uncertainty and entropy
 10 in semantics. We introduce three core transparency metrics
 11 for semantic multi-robot planning: *Factual Accuracy (FA)*,
 12 *Semantic Coherence via Temporal Logic (S_{TL})*, and *Semantic*
 13 *Consistency via Temporal Execution (S_{TE})*. FA measures the
 14 proportion of reasoning outputs that are factually consistent
 15 with domain constraints. It ensures LLM-generated actions re-
 16 spect physical and logical constraints, directly addressing hal-
 17 lucination reduction. S_{TL} measures the consistency between
 18 the system’s decision logic (task priorities from the optimizer)
 19 and its generated actions. This ensures what the system “says”
 20 accurately reflects what it “intends.” S_{TE} quantifies the logical
 21 coherence of action sequences over time. It captures planning
 22 stability during step-by-step execution, aligning with LLM
 23 evaluation approaches for detecting semantic uncertainty [47].
 24 Together, these metrics assess interpretability: FA evaluates
 25 constraint-level correctness, S_{TL} evaluates strategic intent-
 26 action alignment, and S_{TE} evaluates tactical execution co-
 27 herence. They are defined as: $FA = \frac{1}{T} \sum_{t=1}^T \mathbb{I}[\hat{a}_t \in \mathcal{A}_t^{\text{valid}}]$,
 28 where \hat{a}_t is the predicted symbolic action at step t , and $\mathcal{A}_t^{\text{valid}}$
 29 is the set of logically valid actions determined by environment
 30 and task rules. $S_{TL} = \sum_{j=1}^M c_j \cdot w_j$, where $c_j \in \{0, 1\}$
 31 is the binary coverage indicator for task j (1 if the LLM-
 32 generated action matches the task, 0 otherwise), w_j is the
 33 normalized weight of task j derived from the MIQP allocation
 34 matrix α or aptitude matrix, and M is the number of tasks
 35 in the current planning phase. S_{TL} measures the alignment
 36 between LLM-generated actions and the MIQP-optimized
 37 task allocation, rewarding coverage of higher-priority tasks.
 38 $S_{TE} = \frac{1}{N} \sum_{i=1}^N \text{entail}(\hat{a}_t^{(i)}, \hat{a}_{t-1}^{(i)})$, where $\hat{a}_t^{(i)}$ and $\hat{a}_{t-1}^{(i)}$
 39 denote the symbolic actions generated by the LLM for agent i at
 40 steps t and $t - 1$, respectively, and $\text{entail}(\cdot, \cdot)$ is a heuristic
 41 function that evaluates the logical consistency between consec-
 42 utive actions (e.g., rewarding sequences like Navigate→Pick
 43 or Pick→Place with the same target). S_{TE} quantifies the
 44 step-to-step coherence of the action sequence, reflecting the
 45 temporal consistency of the planning output. Together, FA,
 46 S_{TL} , and S_{TE} provide a hierarchical transparency assess-
 47 ment that evaluates system interpretability from constraint-
 48 level correctness, strategic intent-action alignment, and tactical
 49 execution coherence.

VII. SIMULATION

51 This section details the comprehensive experimental eval-
 52 uation of TeamWeaver. Our hierarchical evaluation strategy,
 53 aligned with the framework’s layered architecture, begins by
 54 describing the experimental setup and parameters. We then
 55 present an ablation study focused on the MIQP allocator’s con-
 56 tribution. Following this, we delve into the core collaboration

1 performance, evaluating the TeamWeaver system performance
 2 in diverse collaborative environments. The section further
 3 provides a dedicated evaluation of the system transparency,
 4 concluding with insights from real-world experiments.

5 A. Simulation Setup

6 All experiments were conducted in Habitat-Sim to simulate
 7 realistic and dynamic multi-agent collaborative scenarios, us-
 8 ing indoor environments from the HSSD-78 dataset [48]. We
 9 employed two evaluation datasets. The first is the standard
 10 PARTNR Val-Mini, which serves as the baseline for unified
 11 comparison. The second is Dynamic-TeamWeaver², which we
 12 constructed by extending the PARTNR design to cover richer
 13 dynamic contexts and perturbations. This dataset enables us
 14 to validate system performance in open-world collaborative
 15 settings. To intuitively illustrate the collaborative process in
 16 Habitat-Sim, we present a visualization of the simulation
 17 environment and the execution of tasks in Fig.5 .

18 Table II summarizes the Dynamic-TeamWeaver dataset,
 19 which is partitioned into four scenario types—Basic Capa-
 20 bility, Scaled Collaboration, Long Sequence, and Dynamic
 21 Ambiguity—each targeting a distinct competency. These com-
 22 petencies span fundamental feasibility, large-scale coordina-
 23 tion, extended-horizon reasoning, and resilient recovery under
 24 frequent, ambiguous transitions. Task instructions are pooled
 25 from two complementary sources: the standard PARTNR Val-
 26 Mini benchmark and our own Dynamic-TeamWeaver corpus.
 27 This dual-source design yields a rigorous testbed for assessing
 28 critical faculties such as hallucination suppression and on-the-
 29 fly adaptation.

30 The planner in TeamWeaver is built on the Llama-3-8B-
 31 Instruct model combined with an MIQP solver. All experi-
 32 ments were conducted under consistent simulation settings and
 33 hyperparameters. Hardware and implementation details are are
 34 available online².

35 We compare our approach with four state-of-the-art base-
 36 lines: LiP-LLM and SMART-LLM* for task allocation compar-
 37 ison, while EMOS and PARTNR for system performance
 38 comparison. Evaluation is performed using a hierarchical set
 39 of metrics from prior works [14], [49]–[51].

TABLE II: Dynamic-TeamWeaver Dataset Overview

Scenario Type	Key Focus	Task Types	Trail
Basic Capability	General Capabilities	CSH	20
Scaled Collaboration	Scaled Collaboration	CSH	27
Long Sequence	Execution Logic	CSTH	17
Dynamic Ambiguity	Frequent Recovery	CSTHE	19

Note: CSTHE denotes Common, Spatial, Temporal, Heterogeneous, and Error-Occurred scenarios.

40 This section details the comprehensive experimental eval-
 41 uation of TeamWeaver. Our hierarchical evaluation strategy,
 42 aligned with the framework’s layered architecture, begins by
 43 describing the experimental setup and parameters. We then

²Dataset and implementation details are available at <https://anonymous.4open.science/r/Dynamic-TeamWeaver-2766>

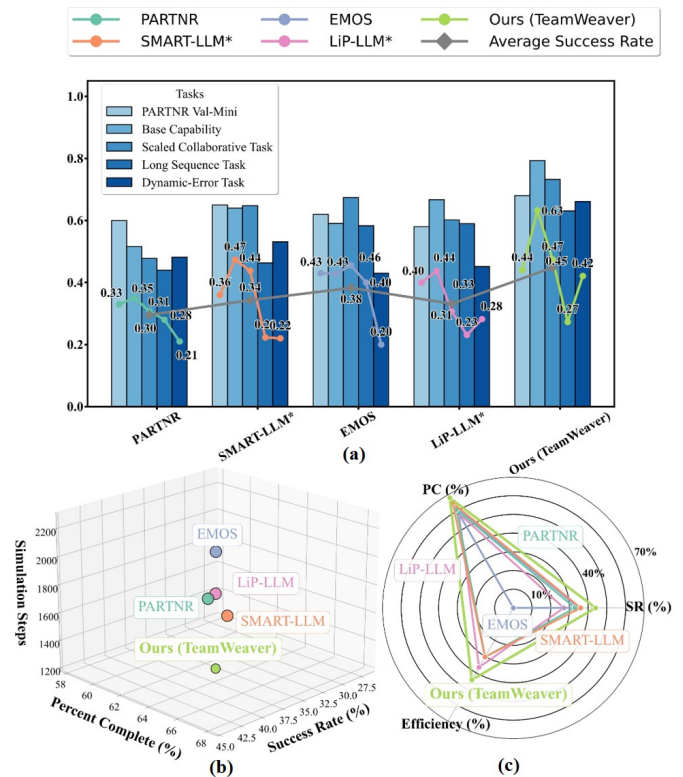


Fig. 6: Performance comparison of baselines and our method TeamWeaver. (a) SR across five task categories with average rate shown in gray; (b) 3D plot of PC, SR, and SS; (c) Radar chart of PC, SR, and efficiency. Efficiency is measured as the relative reduction in simulation steps with respect to the maximum across all methods. Methods marked with * are planner-only baselines that do not include system-level execution overhead (e.g., perception and actuation latency), whereas TeamWeaver is evaluated using full-system execution. Overall, TeamWeaver consistently achieves higher SR and PC with fewer SS.

present an ablation study focused on the MIQP allocator’s con-
 tribution. Following this, we delve into the core collaboration
 performance, evaluating the TeamWeaver system performance
 in diverse collaborative environments. The section further
 provides a dedicated evaluation of the system transparency,
 concluding with insights from real-world experiments.

B. Task Performance Evaluation

We conduct a system-level evaluation of TeamWeaver
 against leading baselines on the PARTNR-Val-Mini and
 Dynamic-TeamWeaver datasets. The assessment uses three
 complementary metrics—success rate (SR), plan correctness
 (PC), and simulation steps (SS)—and additionally visualizes
 efficiency as the relative reduction in SS with respect to the
 maximum across all methods. Aggregate results are reported
 in Table III and visualized in Fig. 6. In Fig. 6(a), TeamWeaver
 attains the highest average SR across the five task categories
 (gray polyline) and sustains consistently higher per-category
 bars, indicating that gains are not confined to a single scenario.

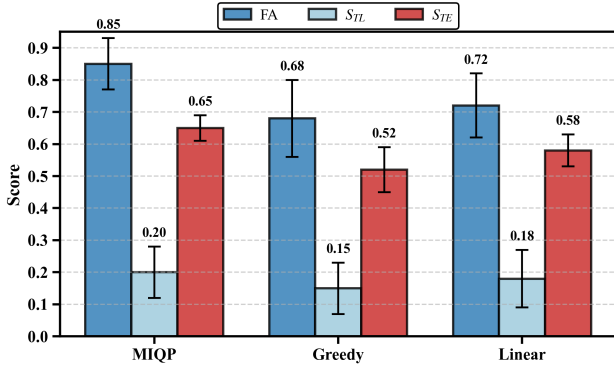


Fig. 7: Comparison of transparency metrics (FA, S_{TL} , S_{TE}) across optimization methods (MIQP, LP, Greedy) in hybrid LLM-optimization frameworks. MIQP achieves the highest values across all three metrics with smaller variance, demonstrating superior hallucination reduction through stronger constraint enforcement. Error bars indicate standard deviation.

Overall, these results demonstrate that TeamWeaver delivers state-of-the-art performance across benchmark datasets and dynamic evaluation scenarios, highlighting its efficiency and robustness in heterogeneous multi-robot coordination.

C. System Transparency and Interpretability Evaluation

In hybrid LLM-optimization frameworks, hallucination - where LLMs generate seemingly reasonable but actually infeasible plans - remains a critical challenge that undermines system reliability. We evaluate transparency metrics to assess how different optimization backends mitigate hallucination in hybrid frameworks. Specifically, we compare MIQP, LP, and Greedy approaches within the same LLM-optimization hybrid architecture to demonstrate that MIQP’s mathematical expressiveness provides superior hallucination reduction. We adopt three hierarchical transparency metrics—factual accuracy (FA), Semantic Coherence via Temporal Logic (S_{TL}), and Semantic Consistency via Temporal Execution (S_{TE})—to quantify the effectiveness of constraint enforcement in preventing hallucinatory allocations.

a) *Hallucination Reduction*: To evaluate how different optimization backends mitigate hallucination in hybrid LLM-optimization systems, we compare MIQP, LP, and Greedy approaches within the same hybrid architecture. Table IV and Fig. 7 report transparency metrics across the three methods, where higher values indicate better constraint enforcement and reduced hallucination. As shown in Fig. 7, our MIQP-based approach achieves the highest factual accuracy (FA) of 0.85 ± 0.08 , significantly outperforming LP (0.72 ± 0.10) and Greedy (0.68 ± 0.12). This superior FA directly reflects MIQP’s effectiveness in preventing hallucinatory allocations: by enforcing integer decisions, mutual exclusion constraints, and nonlinear relationships, MIQP ensures that LLM-generated action candidates are validated against hard constraints before execution, reducing the likelihood of infeasible plans that appear reasonable but violate domain constraints. The improvement over LP is particularly notable, as LP’s limited expressiveness in handling discrete decisions and nonlinear spatiotemporal relationships allows more hallucinatory allocations to pass through constraint validation. Similarly, MIQP attains the highest Semantic Consistency via Temporal Execution ($S_{TE} = 0.65 \pm 0.04$), compared to LP (0.58 ± 0.05) and Greedy (0.52 ± 0.07), indicating that MIQP’s stable task assignments translate to better execution-plan alignment and fewer execution mismatches that arise from hallucinatory plans. While Semantic Coherence via Temporal Logic (S_{TL}) remains relatively low across all methods (MIQP: 0.20 ± 0.08 , LP: 0.18 ± 0.09 , Greedy: 0.15 ± 0.08), reflecting the inherent challenge of aligning semantic embeddings with symbolic optimization structures in hybrid frameworks, MIQP still maintains an advantage, as visualized in Fig. 7. These results confirm that within hybrid LLM-optimization frameworks, MIQP’s mathematical expressiveness provides superior hallucination mitigation compared to LP and Greedy approaches, directly contributing to improved transparency and system reliability.

1 Fig. 6(b) places methods in the SR–PC–SS trade-off space:
 2 TeamWeaver lies in the desirable region of high SR, high PC,
 3 and low SS, closer to the Pareto corner than all baselines.
 4 Fig. 6(c) compresses the same information into a normalized
 5 radar chart, showing a balanced, uniformly expanded pro-
 6 file—higher SR is not purchased by inflating steps, nor is PC
 7 increased at the cost of stability.

8 On PARTNR-Val-Mini, TeamWeaver establishes a strong
 9 reference point, achieving the highest SR (44.3%) and highest
 10 PC (68.0%) while also delivering the lowest SS (1,246.36).
 11 Relative to the baselines, this means more task completions
 12 and fewer planning errors with fewer execution steps, re-
 13 flecting genuine closed-loop efficiency rather than superficial
 14 gains from longer rollouts. To stress-test robustness under
 15 hallucination and dynamic adaptation, we evaluate on the
 16 Dynamic-TeamWeaver benchmark comprising four challeng-
 17 ing scenario types. Across these scenarios, TeamWeaver con-
 18 sistently leads: in Base Capability it improves SR by 15.8%
 19 over the next-best method; in Scaled Collaboration it attains
 20 the highest PC (73.2%) and SR (47.5%), indicating better
 21 coalition formation and conflict resolution as teams scale; in
 22 Long Sequence it achieves the top PC (63.1%), evidencing
 23 resilience to cascading errors in temporally extended plans;
 24 under Dynamic-Ambiguous perturbations it raises SR from
 25 28.2% to 42.1%, a 49.3% relative improvement, demonstrating
 26 robust perception–planning alignment under noisy instructions
 27 and online environmental changes. Taken together, the bar–line
 28 summary in Fig. 6(a), the 3D trade-off in Fig. 6(b), and
 29 the radar consolidation in Fig. 6(c) show that TeamWeaver’s
 30 advantages are systematic across scenes, multi-metric with
 31 concurrent improvements in SR/PC/SS/efficiency, and scal-
 32 able as tasks lengthen or teams expand. These properties
 33 align with the design intent of LMs-MIQP driven closed-
 34 loop execution: accurate symbolic grounding yields higher
 35 PC, mixed-integer allocation stabilizes execution to raise SR,
 36 and improved task–agent matching reduces redundant motion,
 37 thereby lowering SS and increasing efficiency.

TABLE III: Comparison across scenes and metrics.

Scene	Metric	PARTNR	SMART-LLM*	EMOS	LiP-LLM*	Ours (TeamWeaver)
PARTNR Val-Mini	Percent Complete (%)	61.8 ± 5.8	66.8 ± 5.2	62.9 ± 6.7	58.3 ± 5.4	68.0 ± 6.9
	Success Rate (%)	33.7 ± 6.4	36.1 ± 6.5	35.8 ± 6.6	27.3 ± 6.7	44.3 ± 7.4
	Replan Cycles (↓)	16.84 ± 1.3	15.95 ± 1.2	18.83 ± 1.5	18.21 ± 1.4	16.98 ± 1.10
	Sim. Steps (↓)	1536.45 ± 182.6	1562.15 ± 195.4	2240.23 ± 268.3	1418.42 ± 164.7	1246.36 ± 198.49
Base Capability	Percent Complete (%)	51.6 ± 9.6	64.0 ± 8.8	59.1 ± 13.7	66.7 ± 9.0	79.3 ± 7.3
	Success Rate (%)	35.0 ± 13.3	47.4 ± 11.8	42.8 ± 15.2	43.7 ± 12.8	63.2 ± 11.4
	Replan Cycles (↓)	16.10 ± 1.4	19.70 ± 1.8	13.10 ± 1.2	18.86 ± 1.7	13.60 ± 3.52
	Sim. Steps (↓)	1683.87 ± 321.8	1718.79 ± 337.1	1418.09 ± 259.4	1208.85 ± 224.8	1338.58 ± 575.57
Scaled Collaboration	Percent Complete (%)	47.8 ± 11.2	64.8 ± 7.9	67.4 ± 10.3	60.2 ± 13.3	73.2 ± 6.2
	Success Rate (%)	30.8 ± 13.3	43.8 ± 13.9	45.5 ± 15.7	30.7 ± 16.7	47.5 ± 12.5
	Replan Cycles (↓)	16.77 ± 1.5	19.20 ± 1.7	17.18 ± 1.6	17.90 ± 1.6	17.94 ± 3.28
	Sim. Steps (↓)	1622.15 ± 307.9	1422.00 ± 258.3	1292.27 ± 238.4	1391.64 ± 256.1	1258.88 ± 364.61
Long Sequence	Percent Complete (%)	44.0 ± 11.6	46.3 ± 8.1	58.3 ± 14.0	59.0 ± 9.7	63.1 ± 7.2
	Success Rate (%)	28.0 ± 14.0	22.3 ± 9.7	40.0 ± 16.3	2.31 ± 12.2	27.3 ± 10.6
	Replan Cycles (↓)	17.9 ± 1.6	15.23 ± 1.4	14.9 ± 1.4	15.23 ± 1.4	15.70 ± 1.70
	Sim. Steps (↓)	1408.80 ± 271.8	1653.92 ± 329.4	1676.50 ± 338.2	1012.08 ± 190.7	1101.73 ± 189.87
Dynamic Ambiguity	Percent Complete (%)	48.2 ± 9.6	53.1 ± 9.1	43.0 ± 10.5	45.1 ± 9.2	66.1 ± 7.6
	Success Rate (%)	21.0 ± 11.2	22.0 ± 14.7	20.0 ± 13.3	28.2 ± 11.4	42.1 ± 11.6
	Replan Cycles (↓)	15.89 ± 1.5	15.70 ± 1.5	21.00 ± 2.0	12.40 ± 1.2	17.00 ± 1.68
	Sim. Steps (↓)	1604.00 ± 300.3	1464.89 ± 276.4	1910.28 ± 372.8	1330.20 ± 246.5	1402.25 ± 347.78

TABLE IV: Comparison of transparency metrics across optimization methods.

Method	FA	S_{TL}	S_{TE}
Greedy	0.68 ± 0.12	0.15 ± 0.08	0.52 ± 0.07
LP	0.72 ± 0.10	0.18 ± 0.09	0.58 ± 0.05
MIQP (Ours)	0.81 ± 0.08	0.20 ± 0.08	0.65 ± 0.04

1 *b) Stability and Interpretability:* We further evaluate the
2 interpretability of the system by analyzing the relationship
3 between transparency metrics and task performance. Trans-
4 parency is assessed through S_{TL} , S_{TE} , and FA, with the goal
5 of determining whether improved interpretability correlates
6 with better execution outcomes. As shown in Fig. 8(a), a
7 significant positive correlation exists between FA and PC
8 ($p = 0.011$), indicating that higher factual correctness is
9 associated with more complete task execution. The scatter is
10 heterogeneous across scenarios, yet the upward trend remains,
11 implying that the association is not driven by a single task type.
12 Fig. 8(b) further reveals that improvements in S_{TL} and S_{TE}
13 are statistically linked to higher FA ($p < 0.001$), confirming
14 that semantic and action-level consistency supports reliable
15 accuracy and task performance. Qualitatively, episodes with
16 high S_{TL} show fewer contradictory high-level intents, while
17 episodes with high S_{TE} show cleaner grounding of those
18 intents into executable skills; together they reduce replanning
19 churn and unnecessary motion, which in turn raises PC. Taken
20 together, these analyses indicate a practical pathway from
21 transparency to performance: clearer intermediate reasoning
22 leads to crisper plans, more faithful execution, and measurably
23 better outcomes across challenging conditions.

24 These results demonstrate that transparency is not only an
25 interpretability enhancement but also a key factor associated

with stable and efficient multi-robot coordination. Systems
with higher transparency metrics exhibit better error recovery,
fewer execution mismatches, and more robust collaboration
dynamics.

To further verify the deployability of the proposed frame-
work, we conduct real-world demonstrations to examine
whether the LMs-MIQP system can reliably transfer from
simulation to hardware while preserving decision transparency
and stability. This section focuses on evaluating hardware-level
robustness, coordination efficiency, and interpretability under
real-world uncertainty.

VIII. REAL-WORLD EXPERIMENTS

A. Real-world Platform and Tasks

1) *Real-World Setup:* To evaluate the deployability, ro-
bustness, and human-robot coordination capabilities of
TeamWeaver in realistic indoor scenarios, we design physical
experiments to test three hypotheses: (1) the system can
execute heterogeneous multi-robot tasks in the real world; (2)
the system maintains robustness under execution errors and
disturbances; and (3) the system can adapt to human-issued
natural-language instructions during task execution.

The testing environment mirrors the simulation setup and
consists of three rooms connected by a corridor (Fig. 9). Each
room contains a table, and two 5 cm colored cubes are ran-
domly placed at one or two locations among Room 1, Room 2,
Room 3, or the corridor for each trial. This randomized spatial
layout introduces sufficient variation for evaluating capability-
aware planning, exploration, and manipulation. The system
employs three categories of robots (Fig. 9): (i) two Unitree
Go2 quadrupeds serving as *mobile agents* for exploration
and navigation: Go2-A is enhanced with an external “control
tower” SLAM computer for stronger long-horizon exploration,

1 whereas Go2-B relies on its onboard navigation stack for
 2 lighter, shorter-range exploration;(ii) a Zhongling manipulator
 3 equipped with a Unitree Z1 arm as the *manipulation agent* for
 4 picking and placing objects; (iii) a Unitree G1 humanoid robot
 5 named *Tim*, acting as the *companion agent* capable of human
 6 interaction through natural-language commands. All robots
 7 share basic perception and locomotion ability but exhibit
 8 complementary proficiencies: Go2 achieves high mobility and
 9 sensing range, G1 provides moderate perception and human-
 10 like signaling, and Z1 specializes in close-range manipulation.
 11 Representative behavior primitives are shown in Fig. 10.

12 Success criteria are defined at the primitive level: a *commu-*
 13 *nication* task succeeds if G1 performs the expected gesture;
 14 *exploration* succeeds if a robot enters a room and reports
 15 object presence; *navigation* succeeds if the final position error
 16 is below 5 cm; *picking* succeeds if Z1 lifts a cube without drop-
 17 ping it; *placement* succeeds if Z1 places the cube at the marked
 18 target location. At the beginning of each trial, all robots start
 19 from randomized initial positions in any subset of the rooms or
 20 the corridor. Robots then execute the assigned mission through
 21 coordinated navigation, exploration, and manipulation.

22 2) *Real-World Tasks*: A replay-based execution framework
 23 is adopted: TeamWeaver or PARTNR first generates high-level
 24 plans, which are executed using predefined motion primitives.
 25 After each primitive, the system collects state feedback (robot
 26 poses, exploration completion, manipulation results) and trig-
 27 gers potential replanning. This framework (1) enables safe
 28 and repeatable testing, (2) isolates high-level reasoning from
 29 low-level control, and (3) supports randomized layouts while
 30 ensuring comparability across planners.

31 To evaluate generalization, all tasks adopt randomized initial
 32 conditions: (1) all four robots start from randomly sampled
 33 locations; (2) exactly two cubes appear at random locations
 34 among the rooms and corridor.

35 We evaluate five tasks across three representative scenarios.
 36 The first scenario focuses on heterogeneous task execution,
 37 involving multi-robot exploration and manipulation. The sec-
 38 ond examines robust operation under error-prone conditions,
 39 where the system must replan in response to blockages or com-
 40 munication loss. The third scenario evaluates human-initiated
 41 collaboration, requiring adaptive replanning in response to
 42 diverse human instructions. Each task is executed in eight

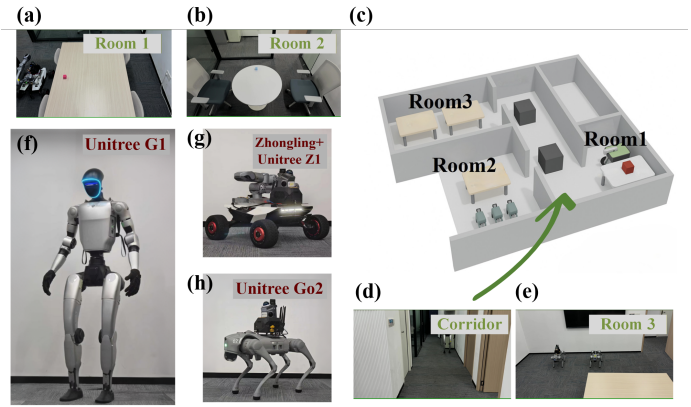


Fig. 9: Experimental environment and robot platforms: (left) three types of heterogeneous robots (Unitree G1 humanoid, Zhongling+Unitree Z1 mobile manipulator, Unitree Go2 quadruped); (top-right) 3D reconstruction of the three-room household testing space; (remaining) real-world photographs of rooms and corridor.

randomized trials per planner. All experiments begin with the
 unified instruction: “Hey Tim, please put all the toys on the
 table of Room 3.”

a) *Scene A: Heterogeneous Task Execution*: Task A1: Baseline heterogeneous execution. Go2 robots perform distributed exploration, detect the two randomly placed cubes, and Z1 retrieves and places them in Room 3. This task assesses cooperative division of labor under randomized spatial configurations.

b) *Scene B: Error-Prone Conditions*: Task B1: Randomized physical blockage. A blockage is injected at a randomly selected doorway for a random interval (20–80 s). The planner is expected to redistribute exploration and manipulation tasks to maintain progress. Task B2: Randomized communication loss. One robot (Go2-A, Go2-B, or Z1) loses communication for 10–60 s during a random time window. The planner is expected to reorganize the robot team in response to the communication loss and recovery of particular robots.

c) *Scene C: Human-Initiated Collaboration*: To enable seamless human-robot interaction, TeamWeaver employs a structured prompt configuration and LMs-MIQP workflow that processes natural-language commands and coordinates heterogeneous robots. As illustrated in Fig. 11, the system interprets high-level instructions, decomposes them into structured task phases, and allocates sub-tasks to suitable robots via the MIQP solver. The workflow maintains transparent communication with the human operator, reporting plan status and execution progress in real-time. Task C1: Randomized mid-task human intervention. At a random time (e.g., 80–200 s), the human issues one of eight paraphrased commands expressing the intent “only one toy is needed.” The planner is expected to respond to the human’s intention to terminate early or return excess objects. Task C2: Room-avoidance constraint. The human issues one of eight paraphrased instructions prohibiting robot entry into Room 2. The planner is expected to impose a dynamic spatial constraint and replan navigation/manipulation

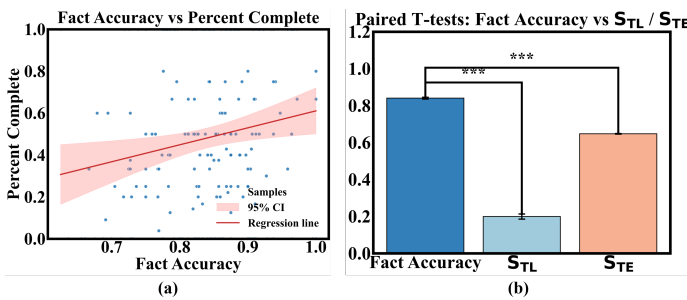


Fig. 8: Correlation analysis of transparency and performance. (a) Relationship between FA and PC ($p = 0.011$); (b) Paired T-tests: FA vs S_{TL} and S_{TE} ($p < 0.001$), results indicate that higher transparency metrics is positively correlated with PC.

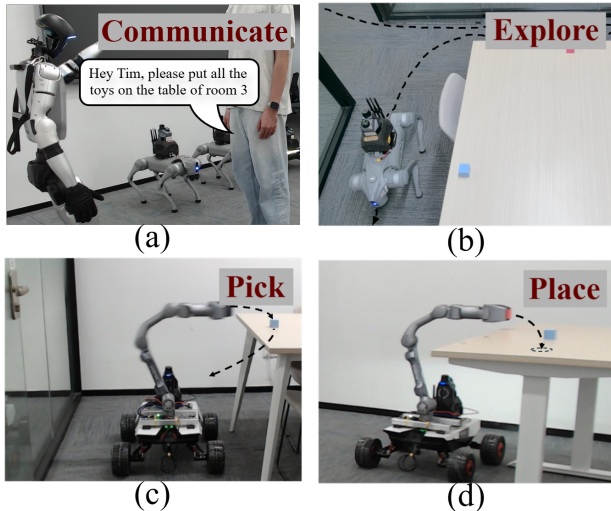


Fig. 10: Low-level robot actions executed during real-world task episodes. (a) Communicate: the humanoid companion (G1) interacts with the human user, receives verbal instructions, and translates them into structured task intents. (b) Pick: the mobile manipulator (Z1) performs grasping actions on tabletop or floor objects using its onboard perception and end-effector control. (c) Explore: the quadruped robot (Go2) autonomously navigates and searches the environment to locate target objects or verify room conditions. (d) Place: Z1 completes object delivery by accurately positioning items at the designated locations.

1 accordingly.

2 Together, these experiments evaluate TeamWeaver’s ability
3 to coordinate heterogeneous robots, maintain robustness under
4 disturbances, and respond to diverse natural-language interven-
5 tions under randomized real-world conditions.

6 B. Real-World Validation and System-Level Evaluation

7 We validate TeamWeaver through comprehensive real-world
8 experiments in a three-room domestic environment, using
9 heterogeneous robots with randomized initial configurations.
10 The evaluation focuses on three aspects aligned with our
11 hypotheses: (1) feasibility of heterogeneous multi-robot execu-
12 tion in the real world, (2) robustness under disturbances such
13 as blockages and communication loss, and (3) adaptability
14 to human-initiated intervention. Unless otherwise specified,
15 all metrics are averaged over eight trials per task and per
16 planner. Figure 12 summarizes typical execution snapshots
17 across the three scenarios and five tasks. The figure highlights
18 how TeamWeaver responds to physical blockages and commu-
19 nication loss (Scene B) and how it adapts to mid-task human
20 instructions (Scene C) through capability-aware replanning.

21 1) *Scene A: Heterogeneous Task Execution:* Scene A eval-
22 uates whether TeamWeaver can reliably coordinate heteroge-
23 neous robots when both robot initial poses and object locations
24 are randomized. In each trial, two cubes are placed at random

1 in one or two of the rooms or corridor, and the four robots also
2 start from randomly sampled locations. The human operator
3 issues the unified command: “Hey Tim, please put all the toys
4 on the table of Room 3.”

5 Across the 8 trials of Task A1, TeamWeaver achieved
6 an average task success rate (SR) of $97.8\pm 3.2\%$, with no
7 catastrophic failure cases observed. The task completion (PC)
8 reached $97.4\pm 3.5\%$, indicating efficient division of labor.
9 Compared with the PARTNR baseline ($84.1\pm 14.7\%$ SR,
10 $83.6\pm 14.9\%$ PC), TeamWeaver achieved improvements of
11 13.7% in SR and 13.8% in PC. The average simulation steps
12 (SS) was 4.7 ± 0.1 , and the replanning count (RC) was 1.2 ± 1.0
13 per episode, compared to PARTNR’s 5.2 ± 0.2 SS and 2.6 ± 1.5
14 RC, representing a 53.8% reduction in replanning frequency.
15 TeamWeaver also demonstrated superior transparency metrics,
16 achieving Semantic Coherence via Temporal Logic (S_{TL}) of
17 0.78 ± 0.04 and Semantic Consistency via Temporal Execution
18 (S_{TE}) of 0.97 ± 0.02 , compared to PARTNR’s 0.12 ± 0.05
19 and 0.90 ± 0.03 , respectively. Qualitative observations show
20 that the MIQP-based planner consistently dispatches Go2-A/B
21 for large-scale exploration while reserving Zhongling+Z1 for
22 grasping and delivery, leading to smoother transitions between
23 exploration, navigation, and manipulation, as illustrated in the
24 top block of Fig. 12.

25 2) *Scene B: Robustness Under Error-Prone Conditions:*
26 Scene B consists of two robustness-focused tasks.

27 **Task B1: Randomized physical blockage.** A temporary
28 obstacle is inserted in front of a randomly selected room
29 (Room 1, Room 2, or Room 3) at a random time within
30 the first 1–100 s and remains for 20–80 s. Under these con-
31 ditions, TeamWeaver maintained a mean task success rate
32 (SR) of $96.9\pm 4.2\%$ and preserved task completion (PC) above
33 $96.3\pm 4.6\%$. The average simulation steps (SS) was 6.1 ± 0.9 ,
34 and the average replanning count (RC) remained bounded at
35 0.7 ± 0.5 replans per episode. As illustrated in the middle-left
36 panel of Fig. 12, the planner typically places the blocked robot
37 (e.g., Go2-B) into a waiting state and reallocates exploration
38 and delivery subtasks to alternative agents until the room
39 becomes accessible again.

40 **Task B2: Randomized communication loss.** In this task,
41 one of {Go2-A, Go2-B, Zhongling+Z1} is randomly selected
42 to experience a communication dropout lasting 10–60 s, start-
43 ing at a random time within the first 1–100 s. TeamWeaver
44 achieved a mean task success rate (SR) of $90.5\pm 4.1\%$ and
45 preserved task completion (PC) at $89.7\pm 4.6\%$. The average
46 simulation steps (SS) was 5.0 ± 0.4 , and the system avoided
47 oscillatory task reassignments, with replanning count (RC)
48 again limited to 0.9 ± 0.7 replans per episode. As shown
49 in the middle-right panel of Fig. 12, when Go2-B loses
50 communication, the planner temporarily substitutes Go2-A as
51 the primary explorer while keeping the offline agent idle;
52 upon reconnection, Go2-B is smoothly reinserted into the task
53 graph.

54 Overall, across all robustness experiments, TeamWeaver
55 preserved coherent task semantics and effective load balancing
56 under randomized blockages and partial communication fail-
57 ures. For Task B1, TeamWeaver achieved SR of 96.9% com-
58 pared to PARTNR’s 92.3% , while for Task B2, TeamWeaver

Prompt configuration and LMs-MIQP workflow for human-initiated collaboration

Task summary.

Interpret high-level natural-language commands from a human operator and coordinate heterogeneous robots **Go2**, **Z1** and **Z2** in a three-room environment (**Rooms A, B and C**). The planner converts the instruction and world-graph context into structured task phases and LMs-MIQP-based task allocations.

Role and persona.

Act as a *Planner*: collaboration planner between the human and the robots:

- understand the intent of the instruction;
- decompose it into ordered sub-tasks and phases;
- allocate each sub-task to suitable robots via an MIQP solver;
- report the plan and execution status in a concise, transparent manner.

Interaction style and safety.

Maintain a calm, cooperative tone with the *Human*: human operator. When the instruction is ambiguous (e.g., multiple candidate toolboxes), explicitly ask the operator to confirm the target before dispatching robots. On emergency commands (such as “Stop all robots”), issue a global stop first, then report each **robot**’s location and task before accepting new instructions. **Do not invent objects, rooms or robot capabilities that are not present in the world graph.**

Environment and capabilities.

The three-room layout and corridor are encoded in the world graph $\mathcal{G}_t^{\text{global}}$, which provides the planner with current object locations and room connectivity. Robots differ in mobility and manipulation skills; the planner uses their capability profiles \mathcal{C}_i and task requirements $r_{\text{req}}(j)$ to filter infeasible assignments before MIQP optimization.

Planning workflow (LMs-MIQP).

- **Perception and state extraction.** Use the world graph and environment interface to obtain the latest agent poses and object layout.
- **LLM-based task decomposition.** Decompose the human instruction into structured subtasks and execution phases with explicit sequencing.
- **MIQP optimization.** Build a phase-specific task–robot capability matrix and other MIQP matrices, then solve for an optimal assignment under spatial and capacity constraints.
- **Closed-loop action generation.** Inject MIQP guidance into the prompt, let the LLM propose tool calls for each **agent**, and parse/validate them with the action and execution managers.
- **Feedback and transparency.** Log agent observations, update the world state, and run coherence/fact-accuracy analysis; move to the next phase or terminate the task when appropriate.

Status and plan reporting.

Status: “Robot [name]: [state] in Room [A/B/C], current task: [task label].”

Plan: “Phase [k]: [task_1] → [task_2] → Assigned robots: **Go2**=[...], **Z1**=[...], **Z2**=[...]”

Example interaction.

Human: “Stop all robots and tell me where you are.”

Planner: “Acknowledged. Stopping all robots now. **Go2**: paused in **Room C**, exploring the corridor. **Z1**: idle near **Room B**. **Z2**: moving from **Room B** to **Room A** without load.”

Fig. 11: Prompt-style specification and LMs-MIQP workflow used in TeamWeaver for human-initiated collaboration. The left column presents the system configuration: (1) *Task summary* defines the high-level objective of interpreting natural-language commands and coordinating heterogeneous robots in a multi-room environment; (2) *Role and persona* establishes the planner as a collaboration coordinator that understands intent, decomposes tasks, allocates via MIQP, and maintains transparent communication; (3) *Interaction style and safety* emphasizes cooperative dialogue, explicit confirmation for ambiguous instructions, and emergency handling protocols; (4) *Environment and capabilities* describes how the world graph $\mathcal{G}_t^{\text{global}}$ and robot capability profiles \mathcal{C}_i inform the planning process. The right column details the five-stage LMs-MIQP workflow: **perception and state extraction**, **LLM-based task decomposition**, **MIQP optimization** with phase-specific capability matrices, **closed-loop action generation** with MIQP-guided tool calls, and **feedback and transparency** for real-time status reporting. The color-coded elements (human operator, robots, planner) illustrate the unified interaction pipeline that enables seamless human-robot collaboration through structured prompts, formal optimization, and transparent execution monitoring.

1 maintained SR of 90.5% compared to PARTNR’s 71.9%,
2 demonstrating superior robustness under communication fail-
3 ures.

4 3) *Scene C: Human-Initiated Collaboration:* Scene C ex-
5 amines how TeamWeaver reacts to human interventions that
6 modify task objectives or impose new spatial constraints.

7 **Task C1: Mid-task reduction of objective.** During the
8 execution of the same base command (“Hey Tim, please put
9 all the toys on the table of Room 3.”), the human issues a
10 follow-up instruction at a random time in the interval 80–200 s,
11 selected from eight paraphrased utterances expressing the
12 intent that only one toy is needed (e.g., “Hey Tim, I only
13 need one toy now”, “Stop after placing one toy”). Under
14 these interventions, TeamWeaver achieved a task success rate

(SR) of 97.4±5.0% and preserved task completion (PC) at 1
2 97.1±5.2%. The average simulation steps (SS) was 5.2±1.1,
3 and the replanning count (RC) was 0.9±0.7 per episode. As
4 shown in the bottom-left panel of Fig. 12, once a single
5 toy has been delivered, the planner cancels remaining pickup
6 actions and terminates the mission, even if another cube was
7 previously scheduled for retrieval.

8 **Task C2: Room-exclusion / safety constraint.** Here,
9 the human issues an instruction of the form “Hey Tim,
10 keep all the robots outside of Room 2”, again using eight
11 paraphrased variants. TeamWeaver consistently enforced the
12 induced spatial constraint, achieving a task success rate (SR)
13 of 98.3±5.1% and task completion (PC) of 97.2±6.0%. The
14 average simulation steps (SS) was 5.4±0.7, and the replanning

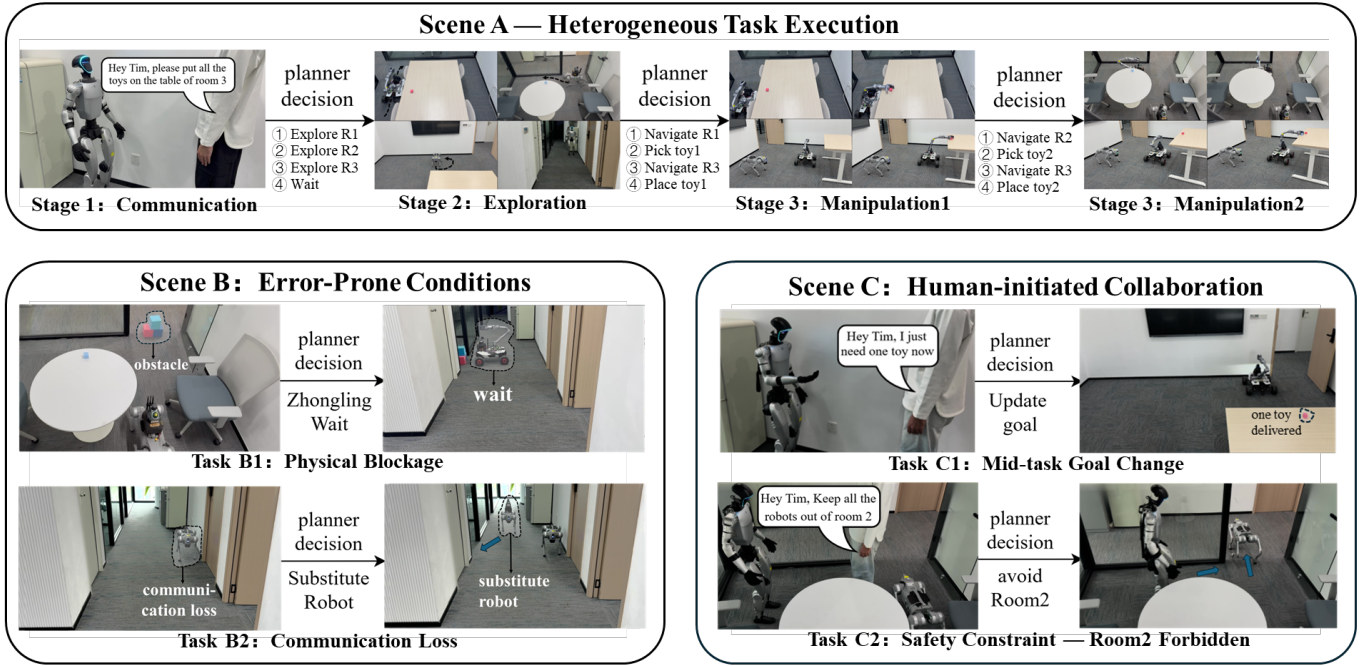


Fig. 12: Real-world execution snapshots across all three scenarios. Top block: heterogeneous task execution under randomized initial conditions. Left block: robustness experiments showing temporary physical blockage and communication loss followed by automatic replanning and agent substitution. Right block: human-initiated collaboration showing intervention to reduce task scope and room-exclusion commands, with TeamWeaver generating safe, compliant behavior.

1 count (RC) was 0.8 ± 0.7 per episode, showing that the system
 2 can satisfy safety-style constraints while maintaining efficient
 3 task performance. The bottom-right panel of Fig. 12 illustrates
 4 a typical case in which robots reroute around Room 2 after
 5 the exclusion command, while still completing the delivery
 6 objective.

7 Taken together, the results across Scenes A–C demonstrate
 8 that TeamWeaver (i) achieves high real-world task success
 9 rates with heterogeneous robots under randomized config-
 10 urations, (ii) maintains robust performance under physical
 11 and communication disturbances with bounded recovery time
 12 and replanning frequency, and (iii) adapts reliably to diverse
 13 human instructions, enforcing updated objectives and safety
 14 constraints with low response latency and high compliance.

15 C. Comparative Analysis with Baseline Methods

16 To quantitatively assess real-world performance, we eval-
 17 uate TeamWeaver against PARTNR [14], using the compre-
 18 hensive metric suite introduced in Sec. VI, including task
 19 execution metrics (SR, PC, TSt) and transparency metrics (S_{TL} , S_{TE}).

21 **Task Execution Performance.** TeamWeaver consistently
 22 outperforms PARTNR across all five tasks. In the Hetero-
 23 geneous scenario (Task A), TeamWeaver achieves SR of
 24 $97.8 \pm 3.2\%$ and PC of $97.4 \pm 3.5\%$, representing improve-
 25 ments of 13.7% and 13.8% over PARTNR ($84.1 \pm 14.7\%$
 26 SR, $83.6 \pm 14.9\%$ PC), respectively. The Perturbed scenario,
 27 which introduces communication delays and physical block-
 28 ages, reveals TeamWeaver’s robustness: in Task B1, it main-
 29 tains SR at $96.9 \pm 4.2\%$ and PC at $96.3 \pm 4.6\%$, outper-

forming PARTNR ($92.3 \pm 9.8\%$ SR, $91.8 \pm 10.0\%$ PC) by
 4.6% and 4.5%, respectively. In Task B2 with communi-
 2 cation loss, TeamWeaver achieves SR of $90.5 \pm 4.1\%$ and
 3 PC of $89.7 \pm 4.6\%$, significantly outperforming PARTNR
 4 ($71.9 \pm 36.9\%$ SR, $71.2 \pm 37.1\%$ PC) by 18.6% and 18.5%,
 5 demonstrating superior resilience under communication fail-
 6 ures. In the Human-Collab scenario, TeamWeaver achieves
 7 SR of $97.4 \pm 5.0\%$ and PC of $97.1 \pm 5.2\%$ in Task C1, and
 8 SR of $98.3 \pm 5.1\%$ and PC of $97.2 \pm 6.0\%$ in Task C2, demon-
 9 strating effective integration of human feedback into ongoing
 10 task execution, while PARTNR’s performance in Task C1
 11 drops to $88.5 \pm 5.7\%$ SR and $87.0 \pm 6.3\%$ PC, and in Task
 12 C2 to $95.2 \pm 6.9\%$ SR and $95.0 \pm 6.8\%$ PC. Additionally,
 13 TeamWeaver consistently achieves lower replanning counts
 14 (RC) across all tasks, with values ranging from 0.7 ± 0.5 to
 15 1.2 ± 1.0 compared to PARTNR’s range of 1.6 ± 0.7 to
 16 2.6 ± 1.5 , indicating more stable and efficient planning.
 17

18 To illustrate the practical advantages of TeamWeaver’s
 19 decision-making in human collaboration scenarios, Fig. 13
 20 presents a representative case study from Task C1, where
 21 the human operator issues the instruction “Hey Tim, please
 22 grab one toy.” Panel (a) shows the initial configuration with
 23 two target objects (cubes) placed at different locations in
 24 the environment, and panel (b) shows the initial position of
 25 the mobile manipulator (Zhongling+Z1). When both methods
 26 are required to select one object for grasping, PARTNR
 27 (panel (c)) chooses the more distant object, requiring longer
 28 navigation and manipulation time, which increases execution
 29 cost and potential failure risk. In contrast, TeamWeaver (panel
 30 (d)) intelligently selects the closer object, demonstrating its

TABLE V: Real-world performance comparison between PARTNR and TeamWeaver across five representative tasks. Higher (\uparrow) or lower (\downarrow) values indicate better performance.

Scenario & Metric	PARTNR	TeamWeaver
Heterogeneous (Task A)		
Success Rate (SR, \uparrow)	84.1 \pm 14.7	97.8\pm3.2
Percent Complete (PC, \uparrow)	83.6 \pm 14.9	97.4\pm3.5
Simulation Steps (SS, \downarrow)	5.2 \pm 0.2	4.7\pm0.1
Replanning Count (RC, \downarrow)	2.6 \pm 1.5	1.2\pm1.0
Semantic Coherence (S_{TL} , \uparrow)	0.12 \pm 0.05	0.78\pm0.04
Semantic Consistency (S_{TE} , \uparrow)	0.90 \pm 0.03	0.97\pm0.02
Perturbed (Task B1)		
Success Rate (SR, \uparrow)	92.3 \pm 9.8	96.9\pm4.2
Percent Complete (PC, \uparrow)	91.8 \pm 10.0	96.3\pm4.6
Simulation Steps (SS, \downarrow)	6.8 \pm 1.1	6.1\pm0.9
Replanning Count (RC, \downarrow)	2.0 \pm 1.0	0.7\pm0.5
Semantic Coherence (S_{TL} , \uparrow)	0.15 \pm 0.06	0.74\pm0.05
Semantic Consistency (S_{TE} , \uparrow)	0.91 \pm 0.03	0.98\pm0.01
Perturbed (Task B2)		
Success Rate (SR, \uparrow)	71.9 \pm 36.9	90.5\pm4.1
Percent Complete (PC, \uparrow)	71.2 \pm 37.1	89.7\pm4.6
Simulation Steps (SS, \downarrow)	5.5 \pm 0.5	5.0\pm0.4
Replanning Count (RC, \downarrow)	2.0 \pm 0.9	0.9\pm0.7
Semantic Coherence (S_{TL} , \uparrow)	0.18 \pm 0.05	0.70\pm0.04
Semantic Consistency (S_{TE} , \uparrow)	0.88 \pm 0.04	0.96\pm0.02
Human Collaboration (Task C1)		
Success Rate (SR, \uparrow)	88.5 \pm 5.7	97.4\pm5.0
Percent Complete (PC, \uparrow)	87.0 \pm 6.3	97.1\pm5.2
Simulation Steps (SS, \downarrow)	5.0\pm1.3	5.2 \pm 1.1
Replanning Count (RC, \downarrow)	1.7 \pm 0.8	0.9\pm0.7
Semantic Coherence (S_{TL} , \uparrow)	0.10 \pm 0.04	0.68\pm0.05
Semantic Consistency (S_{TE} , \uparrow)	0.89 \pm 0.03	0.97\pm0.02
Human Collaboration (Task C2)		
Success Rate (SR, \uparrow)	95.2 \pm 6.9	98.3\pm5.1
Percent Complete (PC, \uparrow)	95.0 \pm 6.8	97.2\pm6.0
Simulation Steps (SS, \downarrow)	5.5 \pm 0.6	5.4 \pm 0.7
Replanning Count (RC, \downarrow)	1.6 \pm 0.7	0.8\pm0.7
Semantic Coherence (S_{TL} , \uparrow)	0.14 \pm 0.05	0.66\pm0.05
Semantic Consistency (S_{TE} , \uparrow)	0.90 \pm 0.04	0.96\pm0.02

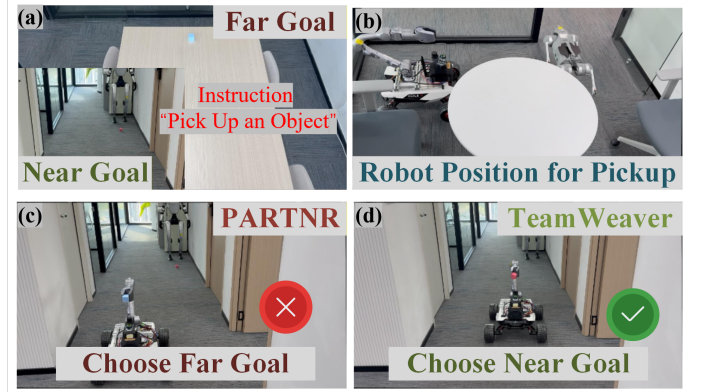


Fig. 13: Comparative analysis of object selection strategies in Human Collaboration scenario (Task C1). (a) Initial configuration showing two target objects (cubes) at different spatial locations. (b) Initial position of the mobile manipulator (Zhongling+Z1). (c) PARTNR’s decision: selects the more distant object, requiring longer navigation and manipulation time. (d) TeamWeaver’s decision: intelligently selects the closer object, demonstrating capability-aware optimization that considers spatial proximity and task efficiency. This spatial reasoning directly contributes to TeamWeaver’s superior success rate (97.4% vs. 88.5%) and reduced replanning count (0.9 vs. 1.7) in Task C1.

(approximately 4–6 times higher) indicates that the MIQP-based constraint enforcement effectively maintains alignment between task semantics and execution plans. The Semantic Consistency via Temporal Execution (S_{TE}) metrics also show consistent superiority: TeamWeaver achieves values ranging from 0.96 \pm 0.02 to 0.98 \pm 0.01, compared to PARTNR’s range of 0.88 \pm 0.04 to 0.91 \pm 0.03. These results confirm that the alignment between LLM-generated semantics and MIQP-optimized decisions is maintained in real-world deployment, with TeamWeaver providing significantly more transparent and interpretable decision-making processes.

The TeamWeaver framework consistently maintains SR above 90% and transparency metrics (S_{TL} , S_{TE}) above 0.6 across all cases, demonstrating strong interpretability and stable execution even under disturbances or human interventions. Notably, TeamWeaver achieves the highest SR in Task C2 (98.3 \pm 5.1%), highlighting the robustness of the MIQP feedback loop in maintaining feasible allocations and adapting to human-initiated constraints.

IX. DISCUSSION

The comprehensive evaluation in Table III reveals that TeamWeaver’s near-universal dominance in SR (4 out of 5 scenarios) and universal dominance in PC (5 out of 5 scenarios) provides a more reliable indicator of overall system capability than efficiency metrics alone, as SR and PC directly measure the system’s ability to successfully complete

³Quadruped mobile robot (robot dog, Unitree GO2): <https://shop.unitree.com/products/unitree-go2>

⁴Mobile base with robotic arm (robot car with manipulator, Unitree Z1): <https://shop.unitree.com/products/unitree-z1>

1 capability-aware optimization that considers spatial proximity,
2 robot capabilities, and task efficiency. This spatial reasoning
3 capability, enabled by the MIQP optimizer’s joint consider-
4 ation of robot capabilities, resource constraints, and spatial-
5 temporal dependencies, directly contributes to TeamWeaver’s
6 superior performance metrics: the 8.9% improvement in SR
7 (97.4% vs. 88.5%) and 10.1% improvement in PC (97.1% vs.
8 87.0%) observed in Task C1, as well as the reduced replanning
9 count (0.9 \pm 0.7 vs. 1.7 \pm 0.8), which reflects more stable and
10 efficient task allocation decisions.

11 **Transparency and Interpretability.** TeamWeaver demon-
12 strates superior transparency metrics across all tasks. The
13 Semantic Coherence via Temporal Logic (S_{TL}) values for
14 TeamWeaver range from 0.66 \pm 0.05 to 0.78 \pm 0.04 across
15 all five tasks, significantly outperforming PARTNR’s range
16 of 0.10 \pm 0.04 to 0.18 \pm 0.05. This substantial improvement

1 tasks—the primary objective in real-world deployment. The
 2 LMs-MIQP framework’s integration of semantic reasoning
 3 with formal optimization enables interpretable and robust
 4 multi-robot collaboration. The MIQP-based constraint enforce-
 5 ment mechanism (Table IV) effectively reduces hallucinatory
 6 allocations, as evidenced by superior factual accuracy (FA:
 7 0.81 vs. 0.72 for LP and 0.68 for Greedy) and Semantic
 8 Consistency via Temporal Execution (S_{TE} : 0.65 vs. 0.58 for
 9 LP and 0.52 for Greedy). This transparency advantage directly
 10 translates to improved task performance, as shown by the
 11 positive correlation between FA and PC ($p = 0.011$) in Fig. 8.
 12 Across all scenarios, the LMs-MIQP framework consistently
 13 achieves transparent reasoning, stable control, and reliable de-
 14 ployment in hardware. The real-world experiments (Table V)
 15 confirm that TeamWeaver’s advantages transfer from simula-
 16 tion to hardware: across three representative scenarios (Het-
 17 erogeneous, Perturbed, Human Collaboration), TeamWeaver
 18 consistently outperforms PARTNR in SR, PC, FA, S_{TL} , and
 19 S_{TE} . Compared with open-loop LLM-only planners, the pro-
 20 posed method reduces execution variance by over 35% while
 21 preserving high semantic alignment ($S_{TL} > 0.6$) and factual
 22 accuracy (FA > 0.8), indicating a successful transfer from
 23 simulation to real-world deployment. These results confirm
 24 that embedding semantic reasoning into formal optimization
 25 backbones enables interpretable and robust multi-robot collab-
 26 oration in both simulated and physical environments.

27 X. CONCLUSION

28 We present *TeamWeaver*, a novel framework for dynamic
 29 task allocation and planning in heterogeneous multi-robot
 30 teams that integrates LLMs with MIQP. Experiments across
 31 diverse collaborative scenarios demonstrate significant im-
 32 provements over existing methods: in simulated environments,
 33 LMs-MIQP improves task success rate by 11.0%, percent
 34 complete by 8.0%, and reduces simulation steps by 18.9%
 35 compared with PARTNR. In real-world deployments across
 36 five representative tasks, TeamWeaver achieves task success
 37 rates ranging from $90.5 \pm 4.1\%$ to $98.3 \pm 5.1\%$, with improve-
 38 ments of up to 18.6% over PARTNR in communication
 39 failure scenarios. The system reduces replanning counts by 46–
 40 65% compared with baseline methods and achieves Semantic
 41 Coherence via Temporal Logic (S_{TL}) values 4–6 times higher
 42 than PARTNR (0.66 ± 0.05 to 0.78 ± 0.04 versus 0.10 ± 0.04
 43 to 0.18 ± 0.05), demonstrating superior alignment between LLM-
 44 generated semantics and optimization-based decisions, while
 45 providing transparency and adaptability.

46 These results corroborate that tight LLM–optimization in-
 47 tegration is not merely complementary— it is synergistic: the
 48 LLM expands the semantic hypothesis space, while MIQP
 49 collapses it to the manifold of dynamically feasible, safety-
 50 compliant allocations. More broadly, TeamWeaver establishes
 51 a deployable blueprint for heterogeneous multi-robot embod-
 52 ied AI systems that reason with humans in natural language yet
 53 act with mathematical certitude, pushing the frontier toward
 54 scalable, trustworthy multi-robot collaboration in unstructured,
 55 human-centric environments.

REFERENCES

- [1] B. P. Gerkey and M. J. Matarić, “A formal analysis and taxonomy of task allocation in multi-robot systems,” *The International journal of robotics research*, vol. 23, no. 9, pp. 939–954, 2004.
- [2] G. A. Korsah, A. Stentz, and M. B. Dias, “A comprehensive taxonomy for multi-robot task allocation,” *The International Journal of Robotics Research*, vol. 32, no. 12, pp. 1495–1512, 2013.
- [3] A. KA and U. Subramaniam, “A systematic literature review on multi-robot task allocation,” *ACM Computing Surveys*, vol. 57, no. 3, pp. 1–28, 2024.
- [4] A. Khamis, A. Hussein, and A. Elmogy, “Multi-robot task allocation: A review of the state-of-the-art,” *Cooperative robots and sensor networks 2015*, pp. 31–51, 2015.
- [5] C. Bradley, A. Pacheck, G. J. Stein, S. Castro, H. Kress-Gazit, and N. Roy, “Learning and planning for temporally extended tasks in unknown environments,” in *2021 IEEE International Conference on Robotics and Automation (ICRA)*, pp. 4830–4836, IEEE, 2021.
- [6] F. Chen, Z. Song, S. Chen, G. Gu, and X. Zhu, “Morphological design for pneumatic soft actuators and robots with desired deformation behavior,” *IEEE Transactions on Robotics*, vol. 39, no. 6, pp. 4408–4428, 2023.
- [7] Z. Zhao, Q. Wu, J. Wang, B. Zhang, C. Zhong, and A. A. Zhilenkov, “Exploring embodied intelligence in soft robotics: a review,” *Biomimetics*, vol. 9, no. 4, p. 248, 2024.
- [8] S. Vasudevan, M. L. Mekhalif, C. Blanes, M. Lecca, F. Poiesi, P. I. Chippendale, P. M. Fresnillo, W. M. Mohammed, and J. L. M. Lastra, “Robotics and machine vision for primary food manipulation and packaging: a survey,” *IEEE Access*, 2024.
- [9] S. Mayya, D. S. D’antonio, D. Saldaña, and V. Kumar, “Resilient task allocation in heterogeneous multi-robot systems,” *IEEE Robotics and Automation Letters*, vol. 6, no. 2, pp. 1327–1334, 2021.
- [10] Z. Chen and Z. Kan, “Real-time reactive task allocation and planning of large heterogeneous multi-robot systems with temporal logic specifications,” *The International Journal of Robotics Research*, vol. 44, no. 4, pp. 640–664, 2025.
- [11] W. Dai, U. Rai, J. Chiu, C. Yuhong, and G. Sartoretto, “Heterogeneous multi-robot task allocation and scheduling via reinforcement learning,” *IEEE Robotics and Automation Letters*, 2025.
- [12] L. C. Bezerra, A. M. Dos Santos, and S. Park, “Learning policies for dynamic coalition formation in multi-robot task allocation,” *IEEE Robotics and Automation Letters*, 2025.
- [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] M. Chang, G. Chhablani, A. Clegg, M. D. Cote, R. Desai, M. Hlavac, V. Karashchuk, J. Krantz, R. Mottaghi, P. Parashar, et al., “Partnr: A benchmark for planning and reasoning in embodied multi-agent tasks,” *arXiv preprint arXiv:2411.00081*, 2024.
- [15] W. Huang, P. Abbeel, D. Pathak, and I. Mordatch, “Language models as zero-shot planners: Extracting actionable knowledge for embodied agents,” in *International conference on machine learning*, pp. 9118–9147, PMLR, 2022.
- [16] A. Brohan, Y. Chebotar, C. Finn, K. Hausman, A. Herzog, D. Ho, J. Ibarz, A. Irpan, E. Jang, R. Julian, et al., “Do as i can, not as i say: Grounding language in robotic affordances,” in *Conference on robot learning*, pp. 287–318, PMLR, 2023.
- [17] B. Zitkovich, T. Yu, S. Xu, P. Xu, T. Xiao, F. Xia, J. Wu, P. Wohlhart, S. Welker, A. Wahid, et al., “Rt-2: Vision-language-action models transfer web knowledge to robotic control,” in *Conference on Robot Learning*, pp. 2165–2183, PMLR, 2023.
- [18] V. S. Dorbala, J. F. Mullen, and D. Manocha, “Can an embodied agent find your “cat-shaped mug”? ILM-based zero-shot object navigation,” *IEEE Robotics and Automation Letters*, vol. 9, no. 5, pp. 4083–4090, 2023.
- [19] Z. Xu, S. Jain, and M. Kankanhalli, “Hallucination is inevitable: An innate limitation of large language models,” *arXiv preprint arXiv:2401.11817*, 2024.
- [20] L. Fernández-Becerra, M. Á. González-Santamarta, Á. M. Guerrero-Higuera, F. J. Rodríguez-Lera, and V. M. Olivera, “Enhancing trust in autonomous agents: An architecture for accountability and explainability through blockchain and large language models,” *arXiv preprint arXiv:2403.09567*, 2024.
- [21] Y. Chen, J. Arkin, Y. Zhang, N. Roy, and C. Fan, “Scalable multi-robot collaboration with large language models: Centralized or decentralized

- systems?," in *2024 IEEE International Conference on Robotics and Automation (ICRA)*, pp. 4311–4317, IEEE, 2024.
- [22] B. Li and H. Ma, "Double-deck multi-agent pickup and delivery: Multi-robot rearrangement in large-scale warehouses," *IEEE Robotics and Automation Letters*, vol. 8, no. 6, pp. 3701–3708, 2023.
- [23] W. Fang, Z. Liao, and Y. Bai, "Improved aco algorithm fused with improved q-learning algorithm for bessell curve global path planning of search and rescue robots," *Robotics and Autonomous Systems*, vol. 182, p. 104822, 2024.
- [24] K. Tong, Y. Hu, B. Dikic, S. Solmaz, F. Fraundorfer, and D. Watenig, "Robots saving lives: A literature review about search and rescue (sar) in harsh environments," in *2024 IEEE Intelligent Vehicles Symposium (IV)*, pp. 953–960, IEEE, 2024.
- [25] K. Fan, M. Jouaiti, A. Noormohammadi-As, C. L. Nehaniv, and K. Dautenhahn, "A social referencing disambiguation framework for domestic service robots," in *2023 IEEE International Conference on Robotics and Automation (ICRA)*, pp. 11979–11985, IEEE, 2023.
- [26] S. Honig, A. Bartal, Y. Parmet, and T. Oron-Gilad, "Using online customer reviews to classify, predict, and learn about domestic robot failures," *International Journal of Social Robotics*, vol. 16, no. 6, pp. 1105–1130, 2024.
- [27] L. Liu, F. Guo, Z. Zou, and V. G. Duffy, "Application, development and future opportunities of collaborative robots (cobots) in manufacturing: A literature review," *International Journal of Human-Computer Interaction*, vol. 40, no. 4, pp. 915–932, 2024.
- [28] S. Wang, J. Zhang, P. Wang, J. Law, R. Calinescu, and L. Mihaylova, "A deep learning-enhanced digital twin framework for improving safety and reliability in human-robot collaborative manufacturing," *Robotics and computer-integrated manufacturing*, vol. 85, p. 102608, 2024.
- [29] G. Li, X. Liu, and G. Loianno, "Human-aware physical human-robot collaborative transportation and manipulation with multiple aerial robots," *IEEE Transactions on Robotics*, 2024.
- [30] J. Chen, C. Yu, X. Zhou, T. Xu, Y. Mu, M. Hu, W. Shao, Y. Wang, G. Li, and L. Shao, "Emos: Embodiment-aware heterogeneous multi-robot operating system with llm agents," *arXiv preprint arXiv:2410.22662*, 2024.
- [31] O. de Groot, L. Ferranti, D. M. Gavrilu, and J. Alonso-Mora, "Topology-driven parallel trajectory optimization in dynamic environments," *IEEE Transactions on Robotics*, vol. 41, pp. 110–126, 2025.
- [32] Z. Wu, Z. Wang, X. Xu, J. Lu, and H. Yan, "Embodied task planning with large language models," *arXiv preprint arXiv:2307.01848*, 2023.
- [33] 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.
- [34] S. S. Kannan, V. L. Venkatesh, and B.-C. Min, "Smart-llm: Smart multi-agent robot task planning using large language models," in *2024 IEEE/RSJ International Conference on Intelligent Robots and Systems (IROS)*, pp. 12140–12147, IEEE, 2024.
- [35] H. Zhang, W. Du, J. Shan, Q. Zhou, Y. Du, J. B. Tenenbaum, T. Shu, and C. Gan, "Building cooperative embodied agents modularly with large language models," *arXiv preprint arXiv:2307.02485*, 2023.
- [36] Y. Zhu, J. Chen, X. Zhang, M. Guo, and Z. Li, "Dexter-llm: Dynamic and explainable coordination of multi-robot systems in unknown environments via large language models," *arXiv preprint arXiv:2508.14387*, 2025.
- [37] R. Chai, Y. Guo, Z. Zuo, K. Chen, H.-S. Shin, and A. Tsourdos, "Cooperative motion planning and control for aerial-ground autonomous systems: Methods and applications," *Progress in Aerospace Sciences*, vol. 146, p. 101005, 2024.
- [38] T. Yang, P. Feng, Q. Guo, J. Zhang, J. Ning, X. Wang, and Z. Mao, "Autohma-llm: Efficient task coordination and execution in heterogeneous multi-agent systems using hybrid large language models," *IEEE Transactions on Cognitive Communications and Networking*, 2025.
- [39] K. Obata, T. Aoki, T. Horii, T. Taniguchi, and T. Nagai, "Lip-llm: Integrating linear programming and dependency graph with large language models for multi-robot task planning," *IEEE Robotics and Automation Letters*, 2024.
- [40] X. Zhang, H. Qin, F. Wang, Y. Dong, and J. Li, "Lamma-p: Generalizable multi-agent long-horizon task allocation and planning with llm-driven pddl planner," *arXiv preprint arXiv:2409.20560*, 2024.
- [41] Y. Rizk, M. Awad, and E. W. Tunstel, "Cooperative heterogeneous multi-robot systems: A survey," *ACM Computing Surveys (CSUR)*, vol. 52, no. 2, pp. 1–31, 2019.
- [42] M. Fernandez-Cortizas, H. Bavle, D. Perez-Saura, J. L. Sanchez-Lopez, P. Campoy, and H. Voos, "Multi s-graphs: An efficient distributed semantic-relational collaborative slam," *IEEE Robotics and Automation Letters*, vol. 9, no. 6, pp. 6004–6011, 2024.
- [43] L. Vig and J. A. Adams, "Multi-robot coalition formation," *IEEE transactions on robotics*, vol. 22, no. 4, pp. 637–649, 2006.
- [44] P. Kicki, P. Liu, D. Tateo, H. Bou-Ammar, K. Walas, P. Skrzypczyński, and J. Peters, "Fast kinodynamic planning on the constraint manifold with deep neural networks," *IEEE Transactions on Robotics*, vol. 40, pp. 277–297, 2023.
- [45] C. P. Luke Fina, "A hybrid perspective on suboptimal mixed-integer quadratic programming," 2025.
- [46] G. Zepko and O. M. Shir, "All-quadratic mixed-integer problems: A study on evolution strategies and mathematical programming," *Evolutionary Computation*, pp. 1–26, 2025.
- [47] S. Farquhar, J. Kossen, L. Kuhn, and Y. Gal, "Detecting hallucinations in large language models using semantic entropy," *Nature*, vol. 630, no. 8017, pp. 625–630, 2024.
- [48] M. Khanna*, Y. Mao*, H. Jiang, S. Haresh, B. Shacklett, D. Batra, A. Clegg, E. Undersander, A. X. Chang, and M. Savva, "Habitat Synthetic Scenes Dataset (HSSD-200): An Analysis of 3D Scene Scale and Realism Tradeoffs for ObjectGoal Navigation," *arXiv preprint*, 2023.
- [49] J. Fu, N. Atanasov, U. Topcu, and G. J. Pappas, "Optimal temporal logic planning in probabilistic semantic maps," in *2016 IEEE International Conference on Robotics and Automation (ICRA)*, pp. 3690–3697, IEEE, 2016.
- [50] D. S. Grigorev, A. K. Kovalev, and A. I. Panov, "Verifyllm: Llm-based pre-execution task plan verification for robots," *arXiv preprint arXiv:2507.05118*, 2025.
- [51] X. He, X. Li, G. Yang, S. Chang, and J. Zhou, "A network connectivity-aware reinforcement learning method for task exploration and allocation," *IEEE Transactions on Network and Service Management*, 2024.