# FLEET: Formal Language-Grounded Scheduling for Heterogeneous Robot Teams

Corban Rivera[†1,2], Grayson Byrd[1,2], Meghan Booker[1], Bethany Kemp[1], Allison Gaines[1],
Emma Holmes[1], James Uplinger[3], Celso M de Melo[3], David Handelman[1]
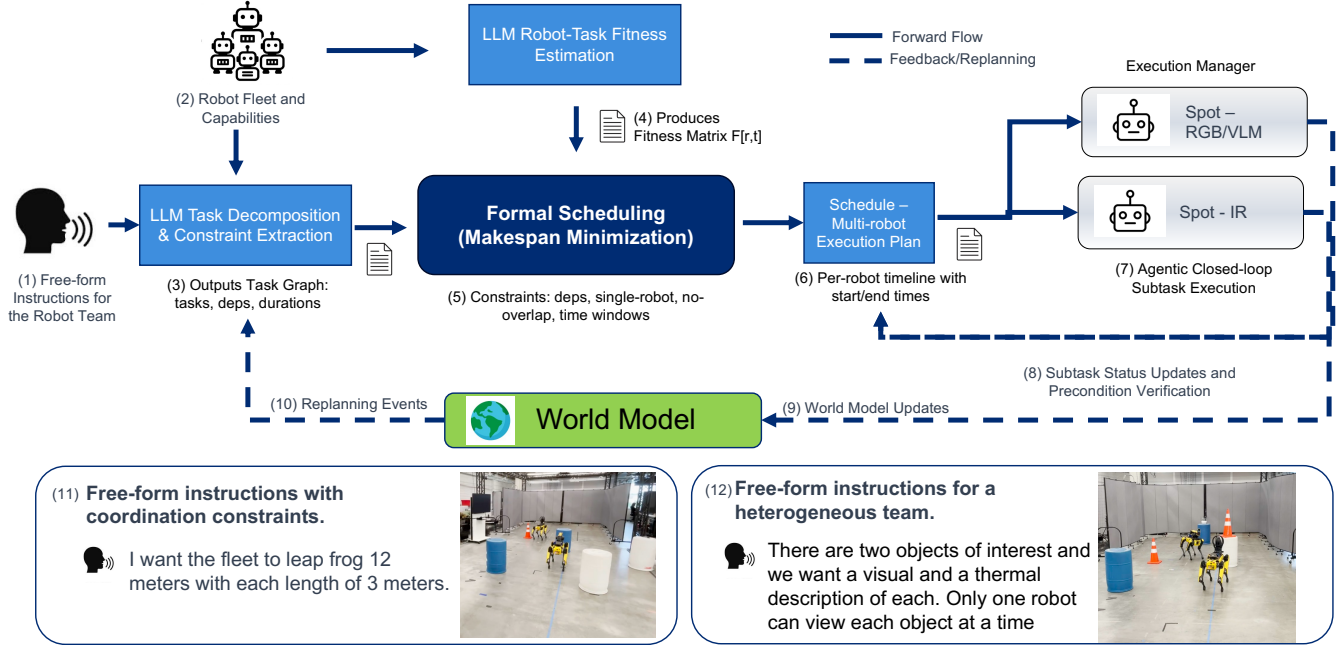[1]JHU APL, [2]JHU, [3]DEVCOM ARL

Fig. 1: **FLEET — Formal Language-grounded Execution and Efficient Teaming** is a hybrid generative–formal framework for natural-language multi-robot tasking. Free-form operator instructions (bottom examples) are ingested by an LLM that (3) decomposes the command into a task graph and constraints and (4) estimates a robot–task fitness matrix. A formal mixed-integer linear programming (MILP) scheduler solves a makespan-minimization problem under precedence, capacity, and spatial constraints to produce a multi-robot schedule (6). Robots execute the plan (7) while streaming status and perception to a World Model (8-9); deviations (delays, failures, new detections) trigger closed-loop replanning back to the LLM and scheduler (10). The architecture supports heterogeneous teams (e.g., IR and RGB/VLM Spots) and yields interpretable artifacts—task graph, fitness matrix, and schedule—that explain decisions.

*Abstract*—Coordinating heterogeneous robot teams from free-form natural-language instructions is hard: language-only planners struggle with long-horizon coordination and hallucination, while purely formal methods require closed-world models. We present *FLEET*, a hybrid decentralized framework that turns language into optimized multi-robot schedules. An LLM front-end produces (i) a task graph with durations and precedence and (ii) a capability-aware robot–task fitness matrix; a formal back-end solves a makespan-minimization problem while the underlying robots execute their free-form subtasks with agentic closed-loop control. Across multiple free-form language-guided autonomy coordination benchmarks, *FLEET* improves success over state of the art generative planners on two-agent teams across heterogeneous tasks. Ablations show that mixed integer linear programming (MILP) primarily improves temporal structure, while LLM-derived fitness is decisive for capability-coupled tasks; together they deliver the highest overall performance. We demonstrate the translation

to real world challenges with hardware trials using a pair of quadruped robots with disjoint capabilities.

## I. INTRODUCTION

Coordinating heterogeneous robot teams in open-world environments remains a central challenge in robotics. Unlike structured factories or warehouses, homes, offices, and disaster sites are dynamic, partially observed, and under-specified at design time. A planner must translate free-form instructions into executable multi-robot strategies, respect capability and resource constraints, and adapt online to delays, perception errors, and newly discovered goals. Purely formal methods provide guarantees but assume closed-world models with carefully engineered predicates and costs; purely generative methods (e.g., LLM planners) offer semantic flexibility and rapid iteration from language, but struggle with long-horizon coordination, precedence tracking, and

hallucination.

We subscribe to the recent trend that language models and formal optimization are complementary. The contributions of hybrid approaches have been primarily demonstrated on single robots. In this work, we introduce hybrid generative and formal optimization concepts for multi-robot coordination. LLMs excel at *semantic front-end* tasks—decomposing natural language into subtasks, exposing commonsense ordering, and explaining why a robot is (or is not) suitable for a role—while a formal *back-end* can guarantee feasibility and optimize the team schedule under explicit constraints. This paper presents *FLEET*, illustrated in Figure 1, a hybrid framework that embeds LLM-derived artifacts into a makespan-minimizing scheduler. Given a free-form instruction and a short profile of each robot, the system produces: (i) a task graph with durations and precedence, (ii) a robot–task fitness matrix aligned with capabilities, and (iii) a multi-robot schedule with start times and assignments. The scheduler enforces precedence, non-overlap on each robot. It runs in *anytime* mode with tight caps and falls back to fast Auction allocators when needed, yielding interpretable plans that can be executed by lightweight reasoning agents.

Empirically, over multiple free-form language-guided autonomy coordination benchmarks [1], our method improves success on heterogeneous team tasks over strong language-model planners, and ablations show that the combination of capability-aware fitness and formal scheduling is crucial. Hardware trials with two Boston Dynamics Spots (IR vs. RGB/VLM) highlight safety benefits (deconfliction in the center field) and reduced idle during cross-modal inspections.

## II. RELATED WORK

### A. Formal Methods for Multi-Robot Coordination

Classic Multi-Robot Task Allocation (MRTA) [2], [3] includes exact assignment with Hungarian [4], market-based auctions and contracts for scalability [3], [5]–[7], and greedy/list-scheduling heuristics with provable bounds in simplified settings [8]. Beyond one-shot pairing, temporal/resource constraints are handled via MILP or constraint programming [9], yielding optimal or bounded-suboptimal schedules with explicit feasibility guarantees. These approaches rely on structured models and carefully tuned costs, which limits application to open-world, language-specified missions.

### B. Generative Planning for Single Robots

Language-guided decision making has advanced via LLM-grounded planners and vision-language action models, e.g., SayCan/SayNav [10], [11], RT-1/RT-2 and related policy models [12]–[14], and tool-use agents such as ReAct and successors [15], [16]. Low-level control is often delegated to trajectory optimizers or learned diffusion/flow policies [17], [18]. While these systems can follow open-ended instructions and improvise with tools, they typically plan for *one* robot and degrade on long-horizon, tightly constrained tasks.

### C. LLM-based Centralized Multi-Robot Planning

LLMs have been used to map team-level goals to subtask assignments and dialogue among agents [16], [19]–[21]. SMART-LLM [22] is a representative centralized planner in which a single LLM decomposes and allocates subtasks to robots. Benchmarks such as PARTNR [1] (built on Habitat [23]) expose persistent weaknesses of purely generative planners on multi-agent tasks: coordination errors, precedence violations, and poor recovery from partial failure.

### D. Hybrid Formal–Generative Single Robot Execution

Recent work uses LLMs to produce symbolic structures for classical solvers—PDDL domain/goal synthesis [24], temporal logic specifications [25], [26], precondition/postcondition checking [27]–[29], factor-graph formulations [30], or linear programs [31]. Two-stage pipelines show that LLMs can propose subgoals which are then solved by a combinatorial back-end [32]. Our approach differs by focusing on multi-robot coordination and *jointly* using (i) an LLM-derived task graph and (ii) an LLM-derived, capability-aware fitness matrix as *inputs* to a makespan-minimizing scheduler. This yields feasible, optimized schedules that remain semantically aligned with the original instruction and robot capabilities, and it enables clear ablations (MILP only, fitness only) and fast fallback to auction assignment under strict latency caps.

### E. Contributions

We contribute a practical and interpretable pipeline for language-guided multi-robot coordination in open-world settings:

- **Hybrid planning architecture.** An LLM front-end produces a task graph and a robot–task fitness matrix; a formal back-end (MILP with anytime settings) computes a schedule that enforces precedence, non-overlap.
- **Allocator plug-ins and anytime operation.** The scheduler supports MILP, Auction back-ends behind a shared interface; MILP runs with strict time/gap caps and warm-starts, and the system falls back to heuristics to guarantee progress.
- **Execution with reasoning agents.** Each robot executes assigned subtasks via a constrained tool-use loop (navigation, perception, manipulation) and streams status to a shared world model, triggering event-based replanning when needed.
- **Comprehensive evaluation.** We compare against state-of-the-art language-model planners and run ablations isolating the roles of fitness and formal scheduling across multiple PARTNR benchmark categories. Hardware trials with two heterogeneous Spots validate safety (deconfliction) and efficiency (reduced idle) in cross-modal tasks.

## III. METHODS

Our framework blends the semantic flexibility of LLMs with the coordination guarantees of formal optimization. A single pipeline converts a free-form instruction into three

Fig. 2: **PARTNR Free-form Language-guided Benchmarks** Partnr free-form language mulit-agent benchmarks builds on habitat-sim and introduces several categories of free-form language-guided tasks to be completed by one or more agents. The categories include "Constraint Free" where the subtasks are separable and do not necessarily depend on each other, "Heterogeneous" where the agents have disjoint capabilities that must be leveraged correctly to complete the tasks, and "Temporal" where the tasks have an implied dependency structure among the subtasks. This Figure illustrates a task from the Heterogeneous task set where agents are acting on the command "Take all of the glasses from the bedroom to the kitchen and wash them". In this example, the human agent can clean and the quadruped robots can not.
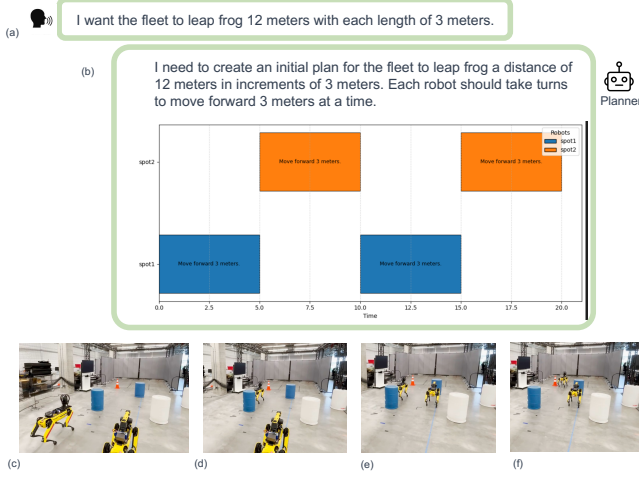


Fig. 3: Hardware trial: **Maneuver with implied dependencies** (a) Operator instruction. (b) Planner output: schedule with 3 m segments and enforced alternation. (c–f) Execution frames showing alternating advances.

concrete artifacts: (i) a *task graph* with durations and precedence, (ii) a *robot–task fitness matrix* that captures capabilities and preferences, and (iii) a *multi-robot schedule* (start times and assignments).

### A. Problem Setup and Artifacts

We consider a heterogeneous team $R = \{1, \ldots, n\}$ and a set of tasks $T = \{1, \ldots, m\}$. Each task $j$ has duration $d_j > 0$, optional precedence constraints $\text{Pred}(j) \subseteq T$, and optional time windows $[r_j, \ell_j]$. Robots have capabilities $\mathcal{C}_i$ (e.g., thermal QA, VLM QA), which induce a feasibility mask $g_{ij} \in \{0, 1\}$ (task $j$ is feasible for robot $i$ iff required capabilities are in $\mathcal{C}_i$). The LLM also returns a normalized fitness score $f_{ij} \in [0, 1]$ indicating how suitable robot $i$ is for task $j$ (higher is better).

*a) Artifacts produced by the LLM.:* Given a natural-language instruction and short robot profiles, the LLM outputs a JSON task list:

- `id`, `description` (subtask string), `duration` $d_j$,
- `dependencies` (list of predecessors), and optional `constraints`: `location` (for travel accounting).

A second prompt maps robot profiles to a scalar $f_{ij}$, forming a matrix $F \in [0, 1]^{n \times m}$. In practice we use few-shot, chain-of-thought style grading and min–max normalization across robots per task to reduce bias. When $F$ is unavailable, we default to uniform scores.

### B. Scheduling Back-end: MILP Formulation

Let $R = \{1, \ldots, n\}$ denote robots and $T = \{1, \ldots, m\}$ tasks. Each task $j$ has duration $d_j > 0$ and precedence edges $E \subseteq T \times T$; $(k, j) \in E$ means $k$ must finish before $j$ starts. Feasibility $g_{ij} \in \{0, 1\}$ encodes capabilities; fitness $f_{ij} \in [0, 1]$ induces a linear cost $c_{ij} = \frac{1}{1 + \gamma f_{ij}} + \tau \, \text{travel}_{ij}$ (set $\tau = 0$ if travel is unused). Pick $M \geq \sum_{j \in T} d_j$. Decision variables: $x_{ij} \in \{0, 1\}$ (assign $j$ to $i$), $s_j \geq 0$ (start time), $y^i_{jk} \in \{0, 1\}$ (on robot $i$, $j$ before $k$), $C_i \geq 0$ (robot completion), and $C_{\max} \geq 0$ (makespan).

$$\min \; \alpha \, C_{\max} + \beta \sum_{i \in R} C_i + \lambda \sum_{i \in R} \sum_{j \in T} c_{ij} \, x_{ij} \qquad (1)$$

*Prioritizes makespan ($\alpha$) with secondary load balance ($\beta$) and a soft preference for capability/fit and optional travel cost ($\lambda$).*

$$\sum_{i \in R} x_{ij} = 1 \qquad \forall j \in T \qquad (2)$$

$$x_{ij} \leq g_{ij} \qquad \forall i \in R, \; \forall j \in T \qquad (3)$$

$$s_j \geq s_k + d_k \qquad \forall (k, j) \in E \qquad (4)$$

*(2) assigns each task to exactly one robot; (3) enforces capability feasibility; (4) respects all precedence edges.*
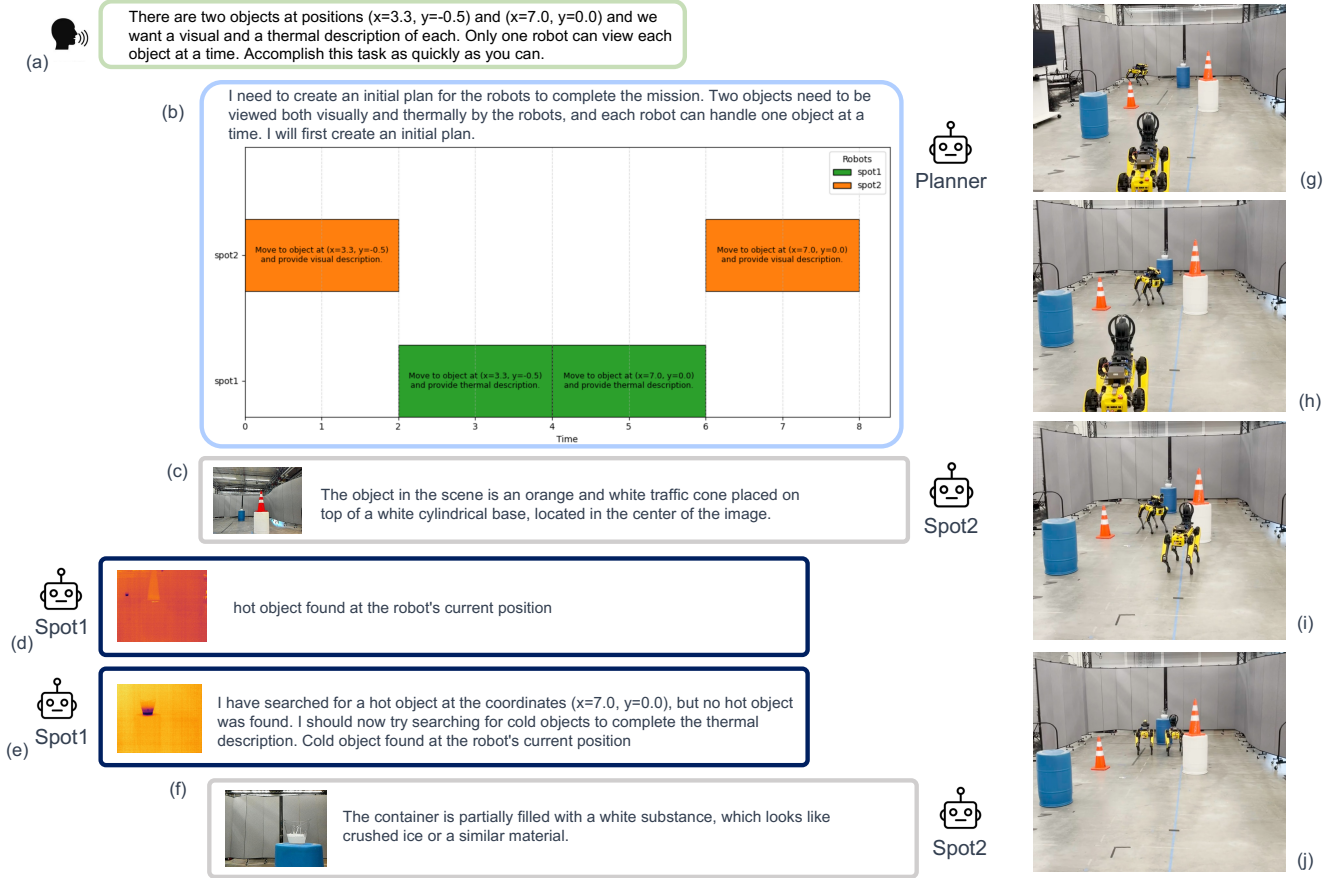
Fig. 4: Hardware trial: **Cross-modal inspection.** The environment included a heat pad under a traffic cone at point-of-interest 1 and a bucket of ice at point of interest 2. The robot team was asked to visually and thermally characterize both points of interest with the additional constraint that only one robot could analyse a point of interest at a time. The robots have disjoint capabilities where Spot-IR can only provide thermal analysis and Spot-RGB/VLM is the only robot that can provide visual question answering. (a) Natural-language instruction. (b) Planner schedule with AND-dependencies (visual+thermal). (c–f) robot QA responses (RGB/IR). (g–j) execution frames. The formal scheduler releases steps on dependency completion, reducing idle time and handoff latency.

$$s_j + d_j \leq s_k + M(1 - y^i_{jk}) \qquad \forall i \in R, \ \forall j < k \quad (5)$$
$$+ M(2 - x_{ij} - x_{ik})$$

*if $x_{ij}=x_{ik}=1$ and $y^i_{jk}=1$, then $j$ must finish before $k$ on robot $i$; otherwise the big-M terms relax the constraint.*

$$s_k + d_k \leq s_j + M y^i_{jk} \qquad \forall i \in R, \ \forall j < k \quad (6)$$
$$+ M(2 - x_{ij} - x_{ik})$$

*Symmetric branch—if $x_{ij}=x_{ik}=1$ and $y^i_{jk}=0$, then $k$ must finish before $j$ on robot $i$; otherwise it relaxes. Together with (5) this prevents overlap on the same robot.*

$$C_i \geq s_j + d_j \qquad \forall i \in R, \ \forall j \in T \quad (7)$$
$$- M(1 - x_{ij})$$

*Robot completion time $C_i$ lower-bounds the finish time of every task assigned to robot $i$.*

$$C_{\max} \geq s_j + d_j \qquad \forall j \in T. \quad (8)$$

*Makespan $C_{\max}$ lower-bounds the finish time of every task (overall completion).*

*Anytime MILP and fallbacks. – FLEET run CBC with a wall-clock cap and a gap stop (e.g., `timeLimit=120 s`, `gapRel=1%`); on timeout, the incumbent plan is used. If no incumbent exists, FLEET falls back to an Auction allocator (same I/O schema). This ensures progress during evaluation and enables quality/latency trade-offs.*

### C. Pluggable Allocators (Auction, MILP)

For scalability studies we swap the MILP with: **Auction:** robots bid for ready tasks using costs $c_{ij}$, iterating until $\epsilon$-optimality; ties break by earliest finish time and robot ID. Dependencies gate task readiness [7]. All allocators emit the same schedule dictionary (agent ID, start, end, metadata), which we visualize as a Gantt chart and push to the execution layer.

### D. Closed-loop Execution with Reasoning Agents

The global schedule specifies *who* does *what* and *when*. Each robot executes its assigned free-form language subtasks with *ConceptAgent* [27]: the LLM can call a small set of verified tools (navigation, perception, manipulation), must return JSON under a fixed schema, and logs status back to a shared world model (task state, detections, poses). *Triggers* for replanning include (i) task completion, (ii) delay beyond a threshold, (iii) perception contradictions (e.g., VLM/IR mismatch), or (iv) newly discovered obstacles. On replanning, we re-score fitness only for impacted tasks and resolve with the same allocator (MILP/Auction), warm-starting from the current partial schedule to preserve stability. Prompt formats are derived from PARTNR [1] decentralized baselines.

*Implementation details.* – We use fixed prompt templates (few-shot for both decomposition and fitness scoring), normalize $f_{ij}$ per task, and clamp durations to positive values. For CBC we set `presolve=true`, enable cuts, and cap threads to available cores. $M$ is chosen as $\sum_j d_j$, and time windows are omitted if absent. When locations are present we include travel either in $c_{ij}$ (cost term) or by augmenting $d_j$ with the preceding leg's travel time; both variants are supported by our code.

## IV. RESULTS

### A. Evaluation in Simulation (PARTNR)

We evaluate on PARTNR [1], illustrated in Figure 2, which scores success on free-form, language-guided multi-robot tasks. Following its taxonomy, we report three categories: (i) *constraint-free* (subtasks can be completed in any order by any agent), (ii) *temporal* (explicit/implicit precedence constraints), and (iii) *heterogeneous* (capability-restricted subtasks). PARTNR provides a verified evaluation function per instruction (goal propositions and constraints), enabling automatic scoring. We omit strictly spatial-relation tasks and focus on categories most aligned with capability- and order-aware scheduling.

*a) Protocol and metrics:* Unless noted, we measure task *success rate* (fraction completed; higher is better; binary variable). We evaluate 2- and 3-agent (when possible) teams. Prompts, seeds, and full configuration are in the supplement. We report macro-averages across PARTNR categories; all methods share the same instruction sets and LLM back-ends.

### B. Comparison to State-of-the-Art Planners (Two Agents)

In Table I, we compare FLEET against a decentralized baseline from PARTNR and the centralized LLM planner SMART-LLM. Scores are averaged over runs with both `gpt-4o` and an open-weight `gpt-oss-20b`.

*Findings.* – (1) *Against decentralized planning.* FLEET improves substantially over the PARTNR decentralized baseline across all categories: +0.24 (constraint-free), +0.42 (heterogeneous), +0.27 (temporal), and +0.31 overall (0.59 vs. 0.28). This reflects the benefit of explicit precedence/resource constraints and capability-aware assignment.

(2) *Against SMART-LLM.* FLEET is comparable with SMART - LLM on constraint-free tasks (0.56 vs. 0.56) and

TABLE I: PARTNR-sim success rate (↑) for two agents. Best per column in **bold**.

| Method | Constraint Free | Hetero-geneous | Temporal | Average |
|---|---|---|---|---|
| FLEET (Ours) | **0.56** | **0.67** | 0.53 | **0.59** |
| PARTNR [1] | 0.32 | 0.25 | 0.26 | 0.28 |
| SMART-LLM [22] | **0.56** | 0.60 | **0.56** | 0.57 |

TABLE II: Ablations on FLEET. Success rate (↑). Avg = macro-average. Best per column in **bold**.

| Method | Ablations | Constraint Free | Hetero-geneous | Temporal | Average |
|---|---|---|---|---|---|
| FLEET | [+MILP +LLM fitness] | **0.56** | **0.67** | **0.53** | **0.59** |
| | [+auction +LLM fitness] | **0.56** | 0.50 | 0.44 | 0.50 |
| | [+MILP] | 0.41 | 0.27 | 0.49 | 0.39 |
| | [base] | 0.32 | 0.25 | 0.26 | 0.28 |

temporal tasks (0.53 vs. 0.56), is stronger on heterogeneous tasks (0.67 vs. 0.60, +0.07), yielding a higher overall average (0.59 vs. 0.57). The heterogeneous gain is consistent with our fitness-guided allocation.

A hybrid approach that pairs LLM-derived task graphs and fitness with a formal scheduler yields competitive constraint-free performance, stronger capability-coupled coordination, and state-of-the-art average success with two agents.

### C. Ablation Studies

We ablate two components of FLEET: the *formal scheduler* (MILP) and the *LLM-based robot–task fitness*. The base variant uses uniform fitness with no formal optimizer; FLEET (+MILP) adds the scheduler; FLEET (+MILP + LLM fitness) adds capability-aware fitness on top of MILP. Scores are macro-averaged across task categories.

*Findings* – (1) **MILP mainly improves temporal structure.** Adding MILP to the base raises *temporal* success from $0.26 \rightarrow 0.49$ (+0.23, $\approx +88\%$ rel.), consistent with enforcing precedence. It also modestly helps *constraint-free* tasks ($0.32 \rightarrow 0.41$, +0.09). The auction assignment strategy matched the performance of MILP for constraint free tasks, but had lower performance than MILP for heterogeneous and temporal tasks.

(2) **LLM fitness is decisive for heterogeneous tasks.** With uniform fitness, MILP barely changes *heterogeneous* performance ($0.25 \rightarrow 0.27$). Injecting LLM-derived fitness with the same MILP lifts it to 0.67 (+0.40 over MILP; +0.42 over base), indicating that capability-aware scoring is the primary driver of correct assignments.

(3) **Components are complementary.** The base averages 0.28; MILP alone reaches 0.39 (+0.11, +39% rel.), and MILP+fitness reaches 0.59 (+0.31, +111% rel.; +0.20 over MILP). Constraint-free tasks also benefit from the full model ($0.32 \rightarrow 0.56$, +0.24), reflecting better load balance and reduced idle from explicit scheduling guided by capability fit. Formal scheduling and capability-aware fitness address

complementary failure modes—temporal feasibility vs. role selection—and together yield the highest performance.

### D. Hardware Trials

We tested *FLEET* using two Boston Dynamic Spot robots with complementary but disjoint sensing capabilities: Spot-IR (thermal QA) and Spot-RGB/VLM (visual QA), both with waypoint navigation and `home`. Each robot executes locally (ConceptAgent [27]); *FLEET* computes the schedule. Figure 4 illustrates a task with heterogeneous and temporal constraints. In *Cross-modal POIs*, two points of interest must be analyzed visually and thermally under a single-server constraint (only one robot may analyze a POI at a time). In *Maneuver with implied dependencies*, the team advances in fixed segments with alternating motion. Figure 3 illustrates *FLEET* performing the trial. Compared to the PARTNR decentralized baseline, *FLEET* encodes precedence and schedules sensing to avoid mid-field conflicts; empirically this increases minimum inter-robot separation and reduces handoff latency and makespan on multi-modal tasks, while matching performance on leap-frog where structure is simple.

## V. CONCLUSIONS

We presented *FLEET*, a hybrid framework that turns free-form language into optimized multi-robot schedules by pairing an LLM front-end (task graph + capability-aware fitness) with a formal back-end (MILP/Auction). The approach enforces precedence, non-overlap constraints and operates in an anytime mode for dependable progress. In simulation (PARTNR), *FLEET* improves success over strong language-model planners—especially on heterogeneous tasks—and ablations show that formal scheduling and LLM-derived fitness are complementary. Hardware trials with two heterogeneous Spots demonstrate practical benefits: deconflicted execution and shorter makespan on cross-modal inspections, while matching baselines on simpler maneuvers.

## VI. ACKNOWLEDGEMENTS

## REFERENCES

[1] 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.

[2] 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.

[3] 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.

[4] H. W. Kuhn, "The hungarian method for the assignment problem," *Naval research logistics quarterly*, vol. 2, no. 1-2, pp. 83–97, 1955.

[5] B. P. Gerkey and M. J. Mataric, "Sold!: Auction methods for multi-robot coordination," *IEEE transactions on robotics and automation*, vol. 18, no. 5, pp. 758–768, 2002.

[6] M. B. Dias, R. Zlot, N. Kalra, and A. Stentz, "Market-based multirobot coordination: A survey and analysis," *Proceedings of the IEEE*, vol. 94, no. 7, pp. 1257–1270, 2006.

[7] D. P. Bertsekas, "A distributed algorithm for the assignment problem," *Lab. for Information and Decision Systems Working Paper, MIT*, vol. 3, 1979.

[8] R. L. Graham, "Bounds on multiprocessing timing anomalies," *SIAM journal on Applied Mathematics*, vol. 17, no. 2, pp. 416–429, 1969.

[9] M. Gombolay, R. Wilcox, and J. Shah, "Fast scheduling of multi-robot teams with temporospatial constraints," 2013.

[10] M. Ahn, A. Brohan, N. Brown, Y. Chebotar, O. Cortes, B. David, C. Finn, C. Fu, K. Gopalakrishnan, K. Hausman *et al.*, "Do as i can, not as i say: Grounding language in robotic affordances," *arXiv preprint arXiv:2204.01691*, 2022.

[11] A. Rajvanshi, K. Sikka, X. Lin, B. Lee, H.-P. Chiu, and A. Velasquez, "Saynav: Grounding large language models for dynamic planning to navigation in new environments," in *Proceedings of the International Conference on Automated Planning and Scheduling*, vol. 34, 2024, pp. 464–474.

[12] A. Brohan, N. Brown, J. Carbajal, Y. Chebotar, J. Dabis, C. Finn, K. Gopalakrishnan, K. Hausman, A. Herzog, J. Hsu *et al.*, "Rt-1: Robotics transformer for real-world control at scale," *arXiv preprint arXiv:2212.06817*, 2022.

[13] A. Brohan, N. Brown, J. Carbajal, Y. Chebotar, X. Chen, K. Choromanski, T. Ding, D. Driess, A. Dubey, C. Finn *et al.*, "Rt-2: Vision-language-action models transfer web knowledge to robotic control," *arXiv preprint arXiv:2307.15818*, 2023.

[14] L. Chen, K. Lu, A. Rajeswaran, K. Lee, A. Grover, M. Laskin, P. Abbeel, A. Srinivas, and I. Mordatch, "Decision transformer: Reinforcement learning via sequence modeling," *Advances in neural information processing systems*, vol. 34, pp. 15 084–15 097, 2021.

[15] S. Yao, J. Zhao, D. Yu, N. Du, I. Shafran, K. Narasimhan, and Y. Cao, "React: Synergizing reasoning and acting in language models," *arXiv preprint arXiv:2210.03629*, 2022.

[16] 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*. PMLR, 2022, pp. 9118–9147.

[17] C. Chi, Z. Xu, S. Feng, E. Cousineau, Y. Du, B. Burchfiel, R. Tedrake, and S. Song, "Diffusion policy: Visuomotor policy learning via action diffusion," *The International Journal of Robotics Research*, p. 02783649241273668, 2023.

[18] K. Nguyen, A. T. Le, T. Pham, M. Huber, J. Peters, and M. N. Vu, "Flowmp: Learning motion fields for robot planning with conditional flow matching," *arXiv preprint arXiv:2503.06135*, 2025.

[19] 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.

[20] 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)*. IEEE, 2024, pp. 286–299.

[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)*. IEEE, 2024, pp. 4311–4317.

[22] 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)*. IEEE, 2024, pp. 12 140–12 147.

[23] X. Puig, E. Undersander, A. Szot, M. D. Cote, T.-Y. Yang, R. Partsey, R. Desai, A. W. Clegg, M. Hlavac, S. Y. Min *et al.*, "Habitat 3.0: A co-habitat for humans, avatars and robots," *arXiv preprint arXiv:2310.13724*, 2023.

[24] J. Xiang, T. Tao, Y. Gu, T. Shu, Z. Wang, Z. Yang, and Z. Hu, "Language models meet world models: Embodied experiences enhance language models," *Advances in neural information processing systems*, vol. 36, pp. 75 392–75 412, 2023.

[25] Z. Wei, X. Luo, and C. Liu, "Hierarchical temporal logic task and motion planning for multi-robot systems," *arXiv preprint arXiv:2504.18899*, 2025.

[26] X. Lin and R. Tron, "Adaptive bi-level multi-robot task allocation and learning under uncertainty with temporal logic constraints," *arXiv preprint arXiv:2502.10062*, 2025.

[27] C. Rivera, G. Byrd, W. Paul, T. Feldman, M. Booker, E. Holmes, D. Handelman, B. Kemp, A. Badger, A. Schmidt *et al.*, "Conceptagent: Llm-driven precondition grounding and tree search for robust task planning and execution," 2025.

[28] S. Mukherjee, C. Paxton, A. Mousavian, A. Fishman, M. Likhachev, and D. Fox, "Reactive long horizon task execution via visual skill and precondition models," in *2021 IEEE/RSJ International Conference on Intelligent Robots and Systems (IROS)*. IEEE, 2021, pp. 5717–5724.

[29] S. S. Raman, V. Cohen, I. Idrees, E. Rosen, R. Mooney, S. Tellex, and D. Paulius, "Cape: Corrective actions from precondition errors using large language models," in *2024 IEEE International Conference on Robotics and Automation (ICRA)*. IEEE, 2024, pp. 14 070–14 077.

[30] U. A. Mishra, Y. Chen, and D. Xu, "Generative factor chaining: Coordinated manipulation with diffusion-based factor graph," in *ICRA 2024 Workshop {\textemdash} Back to the Future: Robot Learning Going Probabilistic*, 2024.

[31] M. Peng, Z. Chen, J. Yang, J. Huang, Z. Shi, Q. Liu, X. Li, and L. Gao, "Automatic milp model construction for multi-robot task allocation and scheduling based on large language models," *arXiv preprint arXiv:2503.13813*, 2025.

[32] D. Bai, I. Singh, D. Traum, and J. Thomason, "Twostep: Multi-agent task planning using classical planners and large language models," *arXiv preprint arXiv:2403.17246*, 2024.