# RECOLLAB: Retrieval-Augmented LLMs for Cooperative Ad-hoc Teammate Modeling

## Conor Wallace, Umer Siddique, Yongcan Cao

Department of Electrical and Computer Engineering,
University of Texas at San Antonio
conor.wallace@my.utsa.edu, muhammadumer.siddique@my.utsa.edu,
yongcan.cao@utsa.edu

### **Abstract**

Ad-hoc teamwork (AHT) requires agents to infer the behavior of previously unseen teammates and adapt their policy accordingly. Conventional approaches often rely on fixed probabilistic models or classifiers, which can be brittle under partial observability and limited interaction. Large language models (LLMs) offer a flexible alternative: by mapping short behavioral traces into high-level hypotheses, they can serve as world models over teammate behavior. We introduce CoLLAB, a language-based framework that classifies partner types using a behavior rubric derived from trajectory features, and extend it to RECoLLAB, which incorporates retrieval-augmented generation (RAG) to stabilize inference with exemplar trajectories. In the cooperative Overcooked environment, CoLLAB effectively distinguishes teammate types, while RECoLLAB consistently improves adaptation across layouts, achieving Pareto-optimal trade-offs between classification accuracy and episodic return. These findings demonstrate the potential of LLMs as behavioral world models for AHT and highlight the importance of retrieval grounding in challenging coordination settings.

## 1 Introduction

Multi-agent reinforcement learning (MARL) has achieved remarkable success in several domains, from zero-sum games such as Go [Silver et al., 2016], to robotics [de Witt et al., 2020], autonomous driving [Zhou et al., 2020b], and cooperative control tasks [Samvelyan et al., 2019, Lin et al., 2023]. Within MARL, cooperative MARL focuses on training teams of agents to solve a common task by interacting with the environment and with each other. Although this particular setting has shown impressive results in controlled environments [Rashid et al., 2020b,a, Son et al., 2019, Yu et al., 2022], it typically assumes that all agents are trained under the same algorithm, which limits its applicability in heterogeneous and realistic settings. Moreover, despite advances in addressing challenges such as non-stationarity [Nekoei et al., 2023], credit assignment [Zhou et al., 2020a], cooperative MARL requires that all teammates be known in advance. However, in practice, many real-world settings often involve heterogeneous agents that are not jointly trained in a shared environment, such as fleets of UAVs with different specifications working together, or different autonomous vehicles sharing the same roads.

To address such scenarios, the *ad-hoc teamwork* (AHT) framework was proposed [Stone et al., 2010], which aims at enabling an autonomous agent to collaborate effectively with previously unseen teammates without prior coordination. Successful AHT requires the ability to rapidly infer the behavior of a teammate and adapt one's own strategy accordingly. This capability is especially important in MARL, where partner policies can be diverse and partially observable. A central challenge in this setting is *teammate modeling*, which refers to constructing representations or beliefs

about another agent's latent policy type based on partial observations of its behavior [Albrecht and Stone, 2018]. Prior approaches include probabilistic reasoning over discrete teammate types, as in PLASTIC [Barrett et al., 2017], and contract-based frameworks such as M³RL [Shu and Tian, 2019], which employ managerial agents to guide self-interested workers. Model-based approaches such as TEAMSTER [Ribeiro et al., 2023] further propose separate modeling of the environment and teammate behaviors in model-based ad hoc teamwork. Meanwhile, recent work on the *N*-Agent Ad Hoc Teamwork (NAHT) framework [Wang et al., 2024] extends AHT to dynamic settings with varying teammates. While effective in some domains, these methods require carefully designed state abstractions and likelihood models that may be brittle under partial observability or in more semantically complex environments.

Large Language Models (LLMs) offer a complementary approach by reasoning over natural-language summaries of observed behavior to infer teammate types [Gao et al., 2024]. LLMs can consume rich, structured prompts that describe recent interaction histories and generate semantically grounded inferences about teammate behavior. Although recent works have demonstrated LLM-based reasoning in competitive games [Bakhtin et al., 2023, Richelieu et al., 2024] and embodied decision-making [Li et al., 2024], their application as *structured world models over teammates* in cooperative MARL domains remains underexplored. Moreover, retrieval-augmented generation (RAG) [Gao et al., 2023] provides a mechanism for grounding LLM predictions in prior experience by retrieving relevant examples from offline trajectories, which can potentially enhance the robustness of LLM-based agents in ambiguous settings.

In this paper, we fill this gap by proposing a novel framework that leverages LLMs as world models for AHT. Specifically, we propose COoperative LLm-based Agent Belief or CoLLAB, which classifies the type of an unknown teammate based on their recent trajectory history and routes the interaction to a pre-trained best-response policy. Building on recent advances in retrieval augmentation, we further introduce RECoLLAB, a retrieval-augmented variant that enriches the LLM prompt with retrieved summaries of similar teammate behaviors from a labeled trajectory database. To support experimental evaluation, we also contribute a labeled dataset for the cooperative Overcooked domain [Carroll et al., 2019, Lowe et al., 2024], which contains five different teammate behavior types induced via reward shaping.

Our main contributions are as follows:

- 1. **Behavioral World Modeling:** We formulate LLM-based teammate type classification as a form of world modeling for AHT.
- 2. **Teammate Reasoning with Retrieval:** We propose COLLAB and RECOLLAB, bridging prompt-based reasoning with retrieval-augmented grounding in MARL.
- 3. Robust Early Teammate Identification: We empirically demonstrate their effectiveness in the Overcooked environment and evaluate our methods against established baselines. Our results demonstrate that LLM-based teammate world models achieve competitive or superior early-type identification compared to baselines, and that retrieval augmentation improves robustness in ambiguous or noisy settings.

#### 2 Related Work

**Ad-hoc Teamwork.** Ad-hoc teamwork (AHT) is a long-standing problem in multi-agent systems, which focuses on enabling autonomous agents to collaborate effectively with previously unseen teammates without prior coordination or shared knowledge [Stone et al., 2010]. Early approaches focused on explicit *teammate modeling*, where the agent maintains a belief over possible teammate types and adapts its policy accordingly. Barrett et al. [2013] explore robust teaming under limited information, and demonstrated applications to more complex domains such as robot soccer [Barrett and Stone, 2015]. Building on this, the PLASTIC framework [Barrett et al., 2017] represents a canonical example in the teammate modeling literature. It uses Bayesian belief updates over a discrete teammate-type space combined with best-response policy selection. Extensions include online learning of teammate models [Albrecht and Stone, 2018], policy reuse, and intention recognition. Model-based methods such as TEAMSTER [Ribeiro et al., 2023] further decouple environment and teammate modeling, while the *N*-Agent Ad Hoc Teamwork (NAHT) framework [Wang et al., 2024]

extends AHT to settings with dynamically varying teammates. While these approaches yield strong performance, they typically rely on carefully designed state abstractions or likelihood models, which may be brittle in partially observable or semantically complex environments.

Multi-Agent Reinforcement Learning. Multi-agent reinforcement learning (MARL) extends single-agent RL to settings where multiple agents interact within a shared environment. In MARL, most of the work focuses on value decomposition, including VDN [Sunehag et al., 2017] and QMIX [Rashid et al., 2020b], which factorize joint action-value functions into agent-wise utilities, and QTRAN [Son et al., 2019], which generalizes decomposition via additional value function constraints. In policy gradient methods, Independent PPO (IPPO) [De Witt et al., 2020] and Multi-Agent PPO (MAPPO) [Yu et al., 2022] train agents with decentralized policies using independent or centralized critics, respectively. In contrast to these MARL methods, which assume training of all agents, AHT considers a single agent that must collaborate with unknown teammates in control of their own actions. In this paper, we adopt IPPO as our base MARL algorithm. For empirical evaluation, we use the cooperative Overcooked environment [Carroll et al., 2019], a benchmark requiring coordination and division of labor, implemented in the JaxMARL framework [Lowe et al., 2024]. This two-player setting enables us to induce distinct teammate types for evaluating our proposed Collab and ReCollab methods.

LLMs for MARL. LLMs have recently been integrated with multi-agent decision making to enable agents that reason about partner intentions and coordinate through language. In the strategic game of Diplomacy, Bakhtin et al. [2023] combined LLM-based communication with search-based planning to achieve human-level play. In cooperative settings, ProAgent shows that LLMs can act as proactive teammates: it uses language to infer task context, anticipate partner needs, and adapt policies for zero-shot coordination, yielding strong gains across cooperative benchmarks [Zhang et al., 2024a]. A recent study on Mutual Theory of Mind (ToM) deploys an LLM agent with ToM and communication modules in a real-time shared-workspace task, finding that language-based modeling improves mutual understanding even when raw task performance gains are modest—underscoring both the promise and limits of LLMs for behavior inference in interactive collaboration [Zhang et al., 2024b].

Retrieval-Augmented Generation. Retrieval-Augmented Generation (RAG) [Gao et al., 2023] combines parametric LLM knowledge with non-parametric retrieval from a database of relevant examples. RAG has been widely used in knowledge-intensive NLP tasks and has recently been applied to embodied agents and multi-modal reasoning [Zhang et al., 2025]. In MARL, retrieval has been explored for retrieving relevant past multi-agent behaviors from a skill database to augment limited demonstrations, enabling more effective policy learning for cooperative mobile robot manipulation tasks [Kuroki et al., 2024]. However, its integration with LLM-based teammate modeling has not yet been studied. Our proposed RECOLLAB applies RAG to retrieve offline trajectory snippets to ground LLM predictions in concrete prior experience, which further improves robustness in AHT.

# **3 Background and Problem Formulation**

We model the cooperative ad-hoc teamwork scenario as a two-agent partially observable Markov game (POMG)  $\mathcal{G} = \langle \mathcal{S}, \mathcal{A}_0, \mathcal{A}_1, \Omega_0, \Omega_1, P, r_0, r_1, \gamma \rangle$  where  $\mathcal{S}$  is the set of environment states,  $\mathcal{A}_0$  and  $\mathcal{A}_1$  are the discrete action spaces for the *teammate* and the *controlled agent*, respectively. In this model,  $P: \mathcal{S} \times \mathcal{A}_0 \times \mathcal{A}_1 \to \Delta(\mathcal{S})$  defines the state transition function,  $\Omega_i: \mathcal{S} \to \Delta(\mathcal{O}_i)$  is the observation function for agent  $i \in \{0,1\}$  producing partial observations  $o_t^i \sim \Omega_i(s_t)$ , rewards are  $r_i: \mathcal{S} \times \mathcal{A}_0 \times \mathcal{A}_1 \to \mathbb{R}$ , and discount factor is  $\gamma \in (0,1]$ . Since we consider the cooperative setting, we assume  $r_0 = r_1 = r$ .

In this model, agents do not have access to the true transition function P (and may not know the exact teammate policy), and must act based on partial observations. The goal of the agent i is to learn a policy  $\pi$  (possibly history dependent), which is a mapping  $\pi^i:\mathcal{H}^i\to\Delta(\mathcal{A}_i)$ , where  $\mathcal{H}^i$  denotes the space of local action-observation histories for agent i. Unless otherwise stated, we consider stochastic Markov policies for the controlled agent and the teammate may exhibit diverse behaviors captured by the considered model type. In this mode, the teammate's policy  $\pi^{0,\tau}$  which is unknown to the controlled agent and is drawn (per trajectory) from a finite set of M possible teammate types  $\mathcal{T}=\{\tau_1,\tau_2,\ldots,\tau_M\}$  where each type represents a distinct behavioral strategy, potentially induced by different training curricula or reward shaping. The controlled agent has access to a corresponding

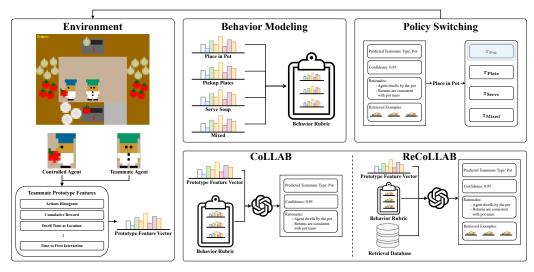


Figure 1: **System diagram of** COLLAB **and** RECOLLAB. The Overcooked environment produces observed trajectories from a controlled agent interacting with a teammate. These trajectories are transformed into prototype feature vectors (e.g., action histograms, dwell times, cumulative reward), which are matched against behavior rubrics to model teammate types. In COLLAB, an LLM classifies the teammate type directly from the behavior rubric. In RECOLLAB, the LLM additionally conditions on retrieved exemplar trajectories from a database, grounding its classification. Both approaches output a predicted teammate type with associated confidence and rationale, which is used to select the best-response policy from a library of trained policies.

set of best-response policies:  $\Pi = \{\pi^{1,\tau_1}, \pi^{1,\tau_2}, \dots, \pi^{1,\tau_M}\}$  where  $\pi^{1,\tau_m}$  is optimized to maximize expected return when paired with teammate of type  $\tau_m$ .

**Problem statement.** We consider the cooperative ad-hoc teamwork scenario where the goal is to infer the teammate type  $\tau$  as quickly and accurately as possible from the changing history  $h_t$  and select actions using the corresponding best-response policy  $\pi^{1,\hat{\tau}}$  to maximize the team's discounted return. We formalize this problem of **teammate type classification** as learning a mapping function:

$$f_{\theta}: \mathcal{H} \to \mathcal{T}, \quad \hat{\tau}_t = f_{\theta}(h_t),$$
 (1)

where  $\mathcal{H}$  is the space of possible histories. In this paper, we model  $f_{\theta}$  as either a prompt-based LLM (CoLLAB) or a retrieval-augmented LLM (RECoLLAB), where  $h_t$  is supplemented with retrieved summaries from an offline trajectory database to improve robustness.

### 4 Methods

To solve the problem of teammate type classification (1), we frame this as a world modeling problem wherein the controlled agent, which is restricted from a fundamental source of information, must infer the teammate's behavior type  $\tau \in \mathcal{T}$  from partial observations in order to form the best-response policy. To achieve this, we leverage LLMs as lightweight world models that reason over structured descriptions of teammate behavior. We propose two methods: CoLLAB, which classifies using statistical prototypes summarized in a behavior rubric, and RECoLLAB, which augments CoLLAB with retrieval from an offline trajectory database to resolve ambiguities and enhance robustness.

### 4.1 Behavior Modeling and Rubric Construction

In complex multi-agent environments such as Overcooked [Carroll et al., 2019], raw trajectories are typically high-dimensional and thus difficult for LLMs to interpret directly. This occurs for two main reasons: 1) LLMs struggle to extract relevant features from episodic trajectories because they lack an understanding of the underlying game mechanics, and 2) they lack knowledge of the specific behavioral characteristics of each teammate type. To address these limitations, we construct

a higher-level rubric of teammate behaviors from an early probing window that captures the most relevant features required to discriminate between teammate types and is easy for an LLM to ingest. Formally, let  $\mathbf{f} = (f_1, f_2, \dots, f_m)$  denote the feature vector computed over the first P steps of an episode probing period, where each feature  $f_j$  encodes a behavioral statistic such as dwell time near a station, number of interactions with pots, or cumulative reward.

**Feature selection.** To retain only discriminative features, we estimate the mutual information (MI), Shannon [1948], between each feature  $f_i$  and the teammate type label  $\tau \in \mathcal{T}$ :

$$I(f_j; \tau) = \sum_{f_j, \tau} p(f_j, \tau) \log \frac{p(f_j, \tau)}{p(f_j) p(\tau)}.$$

This ranking highlights which behavioral dimensions most strongly reduce uncertainty about  $\tau$ . Empirically, features such as dwell time at the window or plate pile exhibit high MI, whereas handoff events and blocked events are less informative. We keep the top-r ranked features (typically  $r \leq 20$ ) to define the rubric.

**Rubric construction.** For each teammate type  $\tau \in \mathcal{T}$ , we compute summary statistics (mean and standard deviation) of the selected features across multiple offline episodes:

$$\mu_{j,\tau} = \mathbb{E}[f_j \mid \tau], \quad \sigma_{j,\tau} = \sqrt{\operatorname{Var}[f_j \mid \tau]}.$$

The rubric function  $r(\mathcal{T}) = \mathtt{Rubric}(\mathcal{T})$  thus consists of behavior prototypes  $\{(\mu_{j,\tau}, \sigma_{j,\tau}) : j = 1, \ldots, r\}$  for each type  $\tau \in \mathcal{T}$ . These serve as reference points describing in natural language how a typical teammate allocates its time and actions in the probe window.

#### 4.2 COLLAB

We first introduce Collab, our base method for teammate type classification in AHT. Collab leverages an LLM as an implicit world model over teammates and classifies the teammate by comparing observed behavior fingerprints against rubric prototypes. The observed behavior prototypes  $\mathbf{f} = (f_1, \dots, f_r)$  from the probe phase are converted into a structured natural language description,  $d(\mathbf{f}) = \mathtt{Describe}(\mathbf{f})$ , which is presented to the LLM along with the behavior rubric. The model outputs the predicted type as:

$$\hat{\tau} = f_{\theta}(d(\mathbf{f}), r(\mathcal{T})),$$

where  $f_{\theta}$  denotes the LLM,  $d(\mathbf{f})$  is the language-based representation of the prototype features, and  $r(\mathcal{T})$  is the behavioral rubric to follow. CoLLAB exploits the LLM's reasoning ability to interpret structured cues without training a task-specific classifier.

### 4.3 RECOLLAB

Although CoLLAB classifies teammates by comparing observed fingerprints against rubric prototypes, it can fail miserably in settings when multiple teammate types yield overlapping statistics (e.g., in the Overcooked environment Plate vs. Mixed). To address this issue, we introduce RECoLLAB, a retrieval-augmented variant that grounds LLM reasoning in concrete trajectory exemplars.

We build an offline database  $\mathcal{D}=\{(h_P^{(i)},\tau)\}_{i=1}^N$  from rollouts up until the probing time P of controlled agents paired with each teammate type  $\tau\in\mathcal{T}$ . We collect N probing trajectories for each teammate type, each of which utilizes a different random seed to ensure a diverse search database. The behavior feature vector  $\mathbf{f}^{(i)}$  is computed from trajectory  $h_P^{(i)}$ , which is then converted to a natural language description using  $d(\mathbf{f}^{(i)})=\mathrm{Describe}(\mathbf{f}^{(i)})$ . This is then embedded into a vector space using an encoder  $E_P=f_\phi(d(\mathbf{f}^{(i)}))$ , where  $f_\phi$  is a language embedding model. These embeddings form the keys for similarity search. At inference, given an observed probe trajectory  $h_P$  and its corresponding behavior prototype description  $d(\mathbf{f})$ , we compute its embedding  $E_P$  and retrieve the top-k nearest neighbors from  $\mathcal{D}$ :

$$\mathcal{R}(d(\mathbf{f})) = \{(d(\mathbf{f}), d(\mathbf{f}^{(i)})) : i \in \mathsf{TopK}_k \big[ s(\phi(d(\mathbf{f})), \phi(d(\mathbf{f}^{(i)}))) \big] \},$$

where  $s(\cdot, \cdot)$  is a cosine similarity function. The LLM is then prompted with the structured description of the observed fingerprint  $d(\mathbf{f})$ , the rubric  $r(\mathcal{T})$ , and the retrieved exemplars  $\mathcal{R}(d(\mathbf{f}))$ . This enhances

Table 1: Classification accuracy across three Overcooked layouts (mean±std over 5 seeds).

Method	Cramped Room	Asymmetric Advantage	Coordination Ring
Random	$0.20{\pm}0.01$	$0.20{\pm}0.01$	$0.20{\pm}0.01$
Logistic Regression	<b>0.96</b> ±0.00	$0.69 \pm 0.15$	$0.81{\pm}0.14$
PLASTIC	$0.81 {\pm} 0.08$	$0.69 \pm 0.15$	$0.58 {\pm} 0.17$
CoLLAB	$0.66 \pm 0.19$	$0.39 \pm 0.00$	$0.35{\pm}0.14$
RECOLLAB ( $k = 5$ )	$0.92{\pm}0.08$	$0.77 \pm 0.12$	<b>0.96</b> ±0.00

Table 2: Cumulative returns across three Overcooked layouts (mean±std over 5 seeds).

Method	Cramped Room	Asymmetric Advantage	Coordination Ring
Oracle	$188.0{\pm}1.6$	$272.0 {\pm} 51.5$	$188.0 {\pm} 32.5$
Static	$45.6 \pm 3.2$	$189.6 \pm 7.4$	$44.0 \pm 0.00$
Random	$58.4 \pm 17.8$	$116.0 \pm 14.3$	$86.4 \pm 20.3$
Logistic Regression	<b>129.6</b> ±10.9	<b>200.8</b> ±12.0	$140.0 \pm 41.4$
PLASTIC	$119.2 \pm 21.1$	$200.8 \pm 12.0$	$133.6 \pm 39.6$
CoLLAB	$103.2 \pm 23.7$	$149.6 \pm 14.2$	$66.4 \pm 21.0$
ReCoLLAB ( $k = 5$ )	$120.8{\pm}15.7$	$181.6 \pm 22.1$	$146.4 \pm 36.5$

the prompt with grounded evidence of how real teammates of each type behave under similar probe conditions and outputs:

$$\hat{\tau} = f_{\theta}(d(\mathbf{f}), r(\mathcal{T}), \mathcal{R}(d(\mathbf{f}))).$$

By grounding abstract rubric features with retrieval, RECOLLAB mitigates classification errors in ambiguous settings and stabilizes predictions.

## 4.4 Policy Adaptation

**Policy adaptation.** Once a type  $\hat{\tau}$  is predicted, the controlled agent selects the corresponding best-response policy  $\pi^{0,\hat{\tau}}$  from the policy library  $\Pi$ . To prevent instability from repeated switching, we route policies only once, immediately after the probe window. This means that a major hyperparameter in our framework is the probe length P of this early trajectory window. The longer this period, the more stable the prediction will be, but at the cost of potential gains in cumulative reward. Both Collab and Recollab integrate information-theoretic fingerprints with language-model reasoning to serve as lightweight behavior models. Rather than fitting a parametric classifier, our methods exploit the LLM's ability to interpret structured cues against a rubric of known behavioral prototypes, functioning as a world model that infers latent teammate types during the probe phase.

# 5 Experimental Setup and Metrics

Environments and Layouts. We evaluate all methods in the cooperative Overcooked environment using the JaxMARL framework [Lowe et al., 2024]. To capture diverse coordination challenges, we select three layouts with distinct structural properties. Our first layout Cramped Room, is a small, narrow kitchen where agents must share a tight space to access ingredients, pots, plates, and the serving window, making blocking inevitable and coordination essential. Our second layout from Overcooked is Asymmetric Advantage, which is an asymmetric kitchen layout where one agent has direct access to key stations (e.g., onion and pot), while the other must navigate a longer path, creating strong role specialization pressures. Our last considered environment is Coordination Ring, which is a circular kitchen where onions, pots, plates, and the serving window are arranged sequentially around a loop, requiring agents to coordinate their direction of movement to avoid blocking. These environments present increasing levels of difficulty, ranging from relatively constrained coordination to highly structured role specialization.

**Teammates and Best Response Policies.** Episodes are 400 timesteps long, with synchronous actions. We train five fixed teammate policies  $\pi^{0,\tau}$ , one for each behavior type induced via reward

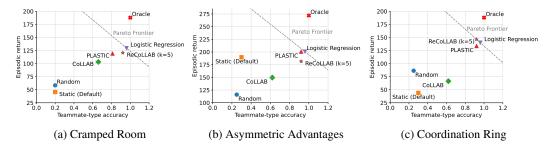


Figure 2: **Pareto Frontier Study.** We plot the teammate-type classification accuracy vs. the obtained cumulative reward and estimate the Pareto optimal frontier. RECOLLAB consistently lies near or directly on the Pareto frontier.

shaping. These include *default* (no reward shaping), *pot-focused* (reward for placing onions in the pot), *plate-focused* (reward for picking up plates), *serve-focused* (reward for delivering completed soups), and *mixed* (combined shaping to encourage proficiency across subtasks). Each teammate policy is trained with PPO for 6 million timesteps. For each type  $\tau \in \mathcal{T}$ , the controlled agent policy  $\pi^{1,\tau}$  is trained with standard rewards while paired with  $\pi^{0,\tau}$ , yeilding a set of best-response policies.

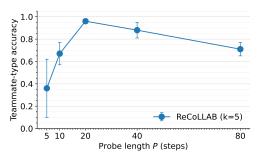
Offline Dataset for RAG. We collect a dataset of 10 evaluation episodes per teammate type per layout. During the probing phase of length P=20 timesteps, we always use the best-response policy to the *default* teammate. Once the probe trajectory has been converted to natural language, we embed the trajectory using the text-embedding-3-large language embedding model from OpenAI [OpenAI, 2024]. This dataset is used exclusively for retrieval in RECoLLAB and is disjoint from evaluation episodes.

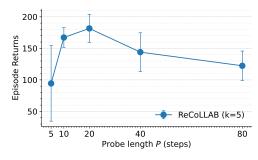
**Evaluation and Metrics.** At the beginning of every episode, we sample a teammate type  $\tau \sim \mathcal{T}$ . As described in the previous section, the best response policy to the *default* teammate is used during the probing phase. We report the teammate type classification as well as the cumulative reward over an episode. For both Collab and ReCollab, we utilize GPT-5 as the base LLM [OpenAI, 2025].

**Baselines.** We compare Collab and ReCollab against several baselines. These include **Oracle**, which has access to the ground-truth teammate type at t=0, **Static**, which maintains the same default policy throughout the episode, and **PLASTIC** [Barrett et al., 2017] baseline, which performs Bayesian belief update with handcrafted likelihood functions over teammate actions. As a sanity check, we also compare our methods with the **Random** baseline. This baseline randomly chooses the best-response policy at every time step t. Finally, the **Logistic Regression** baseline performs the logistic regression classifier fit to the features from the behavior rubric  $r(\mathcal{T})$ .

## 6 Results and Discussion

**Teammate type classification.** Table 1 reports teammate-type classification accuracy across three Overcooked layouts. We observe that a simple Logistic Regression baseline achieves surprisingly strong performance, attaining the highest accuracy in *Cramped Room* and remaining competitive in both *Asymmetric Advantage* and *Coordination Ring*. This finding highlights that the engineered fingerprint features already provide significant discriminative power, such that a lightweight linear classifier can separate teammate types effectively. By contrast, PLASTIC displays variable performance: it achieves strong results in *Cramped Room* but struggles in the other two layouts, reflecting the brittleness of Bayesian updating when teammate behaviors are overlapping or sparsely expressed. CoLLAB, which relies only on rubric-based prompting of an LLM, consistently underperforms across layouts. This suggests that the rubric alone provides insufficient grounding for LLM classification. In contrast, RECollAB substantially improves robustness by incorporating retrieval, outperforming PLASTIC in all three layouts. These results indicate that retrieval grounding stabilizes the LLM's performance, particularly in more challenging layouts where teammate behaviors are less easily separable.





- (a) Classification Accuracy vs. Probe-length
- (b) Episodic Returns vs. Probe-length

Figure 3: **Probe-length ablations.** Teammate-type classification accuracy and episodic returns as a function of probe length P.

**Cumulative returns.** Table 2 compares the cumulative rewards achieved under each method when the controlled agent adapts its best-response policy based on the predicted teammate type. As expected, the Oracle establishes the performance ceiling, while random and static policies tend to stand as the performance floor. Among adaptive methods, Logistic Regression, PLASTIC, and RECOLLAB all produce substantially higher returns than COLLAB, confirming the importance of accurate teammate modeling for ad-hoc teamwork. Notably, RECOLLAB edges out PLASTIC in *Coordination Ring* and surpasses it as well as Logistic Regression in *Coordination Ring*. Logistic Regression again performs very competitively, particularly in *Asymmetric Advantage*, suggesting that in certain layouts simple classifiers suffice for effective adaptation.

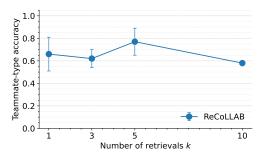
**Discussion.** Figure 2 plots each method's cumulative reward and teammate-type classification accuracy against each other and estimates the Pareto optimal frontier. RECOLLAB consistently lies near, or directly on this frontier. Overall, the results reveal three main insights. First, the strong performance of Logistic Regression highlights the surprising discriminative power of fingerprint features, reinforcing the need to compare against lightweight baselines. Second, RECOLLAB consistently improves upon COLLAB and in many cases rivals or surpasses PLASTIC, especially in more difficult layouts. This demonstrates that retrieval grounding provides stability and robustness beyond what rubric-based prompting alone can achieve. Third, the divergence between classification accuracy and cumulative return underscores that accurate teammate-type inference is necessary but not sufficient; robustness to misclassification and stability of policy switching also play critical roles in maximizing reward. Taken together, these findings suggest that retrieval-grounded LLMs offer an interpretable and extensible alternative to Bayesian methods, while also revealing key limitations of current trajectory representations and prompting strategies. We view the development of richer fingerprint features, hybridized approaches combining retrieval with probabilistic priors, and scaling to more diverse teammate behaviors as promising directions for future work.

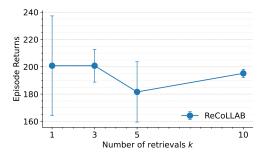
#### 7 Ablation Studies

In addition to the main results, we conduct ablation studies to analyze two key factors in teammate-type classification and adaptive performance: the length of the probe phase (P) and the number of retrieval exemplars (k). We focus on the *Coordination Ring* layout which is the most difficult of the three in terms of teammate behavior complexity. These studies provide insight into the adaptation speed and robustness of each method.

#### 7.1 Effect of Probe Length P

Figure 3 reports classification accuracy and cumulative reward as a function of probe length  $P \in \{5, 10, 20, 40, 80\}$ . We find that both classification accuracy and cumulative returns are maximized at probe length P=20. Accuracy rises from near-random at P=5 to nearly perfect at P=20, after which performance begins to diminish. The dropoff in returns is due to longer probe phases increasing the time of sub-optimal cooperation which can be more difficult to break out of. These results highlight a trade-off between adaptation speed and performance: shorter probes allow faster





- (a) Classification Accuracy vs. Number of Retrievals
- (b) Episodic Returns vs. Number of Retrievals

Figure 4: **Number of retrievals ablations.** Teammate-type classification accuracy and episodic returns as a function of the number of retrieved exemplars k.

decision-making but risk misclassification, while a probe at length P=20 strikes the best balance by achieving high accuracy and near-optimal returns without excessive delay.

## 7.2 Effect of Retrieval Exemplars k

We next vary the number of retrieved exemplar trajectories  $k \in \{1, 3, 5, 10\}$  to evaluate the impact of retrieval grounding. Results are shown in Figure 4. Varying the number of retrieved exemplars yielded only modest differences in both teammate-type classification accuracy and episodic returns. Accuracy remains relatively stable across values of k, and returns do not show consistent improvements with larger retrieval sets. This suggests that while retrieval grounding is important for RECOLLAB 's overall performance, the precise number of exemplars plays a less critical role. In practice, small values such as k=3-5 appear sufficient to stabilize performance, and additional retrievals offer diminishing returns. These results highlight that the primary benefits of RECOLLAB arise from conditioning on retrieved examples at all, rather than from scaling the retrieval set size.

## 8 Conclusion

We introduced Collab, an LLM-based framework for teammate type classification and policy routing in ad-hoc teamwork, and ReCollab, its retrieval-augmented variant that grounds predictions in prior trajectory data. By framing type classification as a form of world modeling over teammates, our approach leverages the reasoning capabilities of LLMs while maintaining structured outputs for downstream control.

In the cooperative Overcooked domain, RECOLLAB achieves competitive or superior early classification accuracy and team reward compared to the PLASTIC and Logistic Regression baselines. These findings suggest that LLM-based world models, when paired with targeted retrieval, can serve as effective agents in structured multi-agent RL settings without extensive fine-tuning. Our work opens several directions for future exploration, including scaling to continuous teammate behavior spaces, integrating online policy adaptation, and extending retrieval to multi-modal or human-agent teaming scenarios.

By bridging structured LLM reasoning, retrieval, and best-response policy selection, we take a step toward more flexible, generalizable agents capable of robust ad-hoc teamwork in complex environments.

## Acknowledgments and Disclosure of Funding

This work was supported by the Office of Naval Research under Grant N000142412405 and the Army Research Office under Grants W911NF2110103 and W911NF2310363.

### References

- Stefano V Albrecht and Peter Stone. Autonomous agents modelling other agents: A comprehensive survey and open problems. *Artificial Intelligence*, 258:66–95, 2018.
- Anton Bakhtin, Noam Wu, Adam Lerer, et al. Human-level play in the game of diplomacy by combining language models with strategic reasoning. In *Nature*, volume 620, pages 526–531, 2023.
- Samuel Barrett and Peter Stone. Cooperating with unknown teammates in complex domains: A robot soccer case study of ad hoc teamwork. In *Proceedings of the Twenty-Ninth Conference on Artificial Intelligence (AAAI)*, page 943–949, 2015.
- Samuel Barrett, Peter Stone, Sarit Kraus, and Avi Rosenfeld. Teamwork with limited knowledge of teammates. In *Proceedings of the Twenty-Seventh AAAI Conference on Artificial Intelligence (AAAI)*, pages 1219–1225, 2013.
- Samuel Barrett, Avi Rosenfeld, Sarit Kraus, and Peter Stone. Making friends on the fly: Cooperating with new teammates. *Artificial Intelligence*, 242:132–171, 2017.
- Micah Carroll, Rohin Shah, Mark Ho, Thomas L Griffiths, Sanjit A Seshia, Pieter Abbeel, and Anca D Dragan. Utility learning from trajectories for cooperative games. In *Advances in Neural Information Processing Systems*, volume 32, 2019.
- Christian Schroeder De Witt, Tarun Gupta, Denys Makoviichuk, Viktor Makoviychuk, Philip HS Torr, Mingfei Sun, and Shimon Whiteson. Is independent learning all you need in the starcraft multi-agent challenge? *arXiv preprint arXiv:2011.09533*, 2020.
- Christian Schroeder de Witt, Bei Peng, Pierre-Alexandre Kamienny, Philip Torr, Wendelin Böhmer, and Shimon Whiteson. Deep multi-agent reinforcement learning for decentralized continuous cooperative control. *arXiv preprint arXiv:2003.06709*, 19, 2020.
- Chen Gao, Xiaochong Lan, Nian Li, Yuan Yuan, Jingtao Ding, Zhilun Zhou, Fengli Xu, and Yong Li. Large language models empowered agent-based modeling and simulation: A survey and perspectives. *Humanities and Social Sciences Communications*, 11(1):1–24, 2024. doi: 10.1057/s41599-024-03611-3.
- Luyu Gao et al. Retrieval-augmented generation for knowledge-intensive nlp tasks: A survey. *arXiv* preprint arXiv:2312.10997, 2023.
- So Kuroki, Mai Nishimura, and Tadashi Kozuno. Multi-agent behavior retrieval: Retrieval-augmented policy training for cooperative push manipulation by mobile robots. In 2024 IEEE/RSJ International Conference on Intelligent Robots and Systems (IROS). IEEE, 2024.
- Manling Li, Shiyu Zhao, Qineng Wang, Kangrui Wang, Yu Zhou, Sanjana Srivastava, Cem Gokmen, Tony Lee, Erran Li, Ruohan Zhang, Weiyu Liu, Percy Liang, Fei-Fei Li, Jiayuan Mao, and Jiajun Wu. Embodied agent interface: Benchmarking LLMs for embodied decision making. In *Proceedings of the 38th International Conference on Neural Information Processing Systems (NeurIPS)*, Datasets and Benchmarks Track, 2024. URL https://arxiv.org/abs/2410.07166.
- Fanqi Lin, Shiyu Huang, Tim Pearce, Wenze Chen, and Wei-Wei Tu. Tizero: Mastering multi-agent football with curriculum learning and self-play. *arXiv preprint arXiv:2302.07515*, 2023.
- Ryan Lowe, Frans Oliehoek, et al. Jaxmarl: Multi-agent reinforcement learning in jax. https://github.com/FLAIROx/JaxMARL, 2024.
- Hadi Nekoei, Akilesh Badrinaaraayanan, Amit Sinha, Mohammad Amini, Janarthanan Rajendran, Aditya Mahajan, and Sarath Chandar. Dealing with non-stationarity in decentralized cooperative multi-agent deep reinforcement learning via multi-timescale learning. In *Conference on Lifelong Learning Agents*, pages 376–398. PMLR, 2023.
- OpenAI. text-embedding-3-large. https://platform.openai.com/docs/models/text-embedding-3-large, 2024. Accessed date.

- OpenAI. GPT-5. https://openai.com/gpt-5/, 2025. Released August 7, 2025.
- Tabish Rashid, Gregory Farquhar, Bei Peng, and Shimon Whiteson. Weighted qmix: Expanding monotonic value function factorisation for deep multi-agent reinforcement learning. *Advances in neural information processing systems*, 33:10199–10210, 2020a.
- Tabish Rashid, Mikayel Samvelyan, Christian Schroeder De Witt, Gregory Farquhar, Jakob Foerster, and Shimon Whiteson. Monotonic value function factorisation for deep multi-agent reinforcement learning. *Journal of Machine Learning Research*, 21(178):1–51, 2020b.
- João G. Ribeiro, Gonçalo Rodrigues, Alberto Sardinha, and Francisco S. Melo. TEAMSTER: Model-based reinforcement learning for ad hoc teamwork. *Artificial Intelligence*, 324:104013, 2023. doi: 10.1016/j.artint.2023.104013.
- Guillaume Richelieu et al. Self-evolving llm-based agents for ai diplomacy. In *Advances in Neural Information Processing Systems*, 2024.
- Mikayel Samvelyan, Tabish Rashid, Christian Schroeder De Witt, Gregory Farquhar, Nantas Nardelli, Tim GJ Rudner, Chia-Man Hung, Philip HS Torr, Jakob Foerster, and Shimon Whiteson. The starcraft multi-agent challenge. *arXiv* preprint arXiv:1902.04043, 2019.
- Claude E. Shannon. A mathematical theory of communication. *Bell System Technical Journal*, 27(3): 379–423, 1948.
- Tianmin Shu and Yuandong Tian. M<sup>3</sup>RL: Mind-aware multi-agent management reinforcement learning. In *International Conference on Learning Representations (ICLR)*, 2019.
- David Silver, Aja Huang, Chris J Maddison, Arthur Guez, Laurent Sifre, George Van Den Driessche, Julian Schrittwieser, Ioannis Antonoglou, Veda Panneershelvam, Marc Lanctot, et al. Mastering the game of go with deep neural networks and tree search. *nature*, 529(7587):484–489, 2016.
- Kyunghwan Son, Daewoo Kim, Wan Ju Kang, David Earl Hostallero, and Yung Yi. Qtran: Learning to factorize with transformation for cooperative multi-agent reinforcement learning. In *International conference on machine learning*, pages 5887–5896. PMLR, 2019.
- Peter Stone, Gal A. Kaminka, Sarit Kraus, and Jeffrey S. Rosenschein. Ad hoc autonomous agent teams: Collaboration without pre-coordination. In *Proceedings of the Twenty-Fourth AAAI Conference on Artificial Intelligence (AAAI)*, pages 1504–1509, 2010.
- Peter Sunehag, Guy Lever, Audrunas Gruslys, Wojciech Marian Czarnecki, Vinicius Zambaldi, Max Jaderberg, Marc Lanctot, Nicolas Sonnerat, Joel Z Leibo, Karl Tuyls, et al. Value-decomposition networks for cooperative multi-agent learning. *arXiv preprint arXiv:1706.05296*, 2017.
- Caroline Wang, Muhammad Arrasy Rahman, Ishan Durugkar, Elad Liebman, and Peter Stone. Nagent ad hoc teamwork. *Advances in Neural Information Processing Systems*, 37:111832–111862, 2024.
- Chao Yu, Akash Velu, Eugene Vinitsky, Jiaxuan Gao, Yu Wang, Alexandre Bayen, and Yi Wu. The surprising effectiveness of ppo in cooperative multi-agent games. *Advances in neural information processing systems*, 35:24611–24624, 2022.
- Ceyao Zhang, Kaijie Yang, Siyi Hu, Zihao Wang, Guanghe Li, Yihang Sun, Cheng Zhang, Zhaowei Zhang, Anji Liu, Song-Chun Zhu, Xiaojun Chang, Junge Zhang, Feng Yin, Yitao Liang, and Yaodong Yang. ProAgent: Building proactive cooperative agents with large language models. In *Proceedings of the AAAI Conference on Artificial Intelligence*, volume 38, pages 17591–17599, 2024a. doi: 10.1609/aaai.v38i16.29710.
- Shao Zhang, Xihuai Wang, Wenhao Zhang, Yongshan Chen, Landi Gao, Dakuo Wang, Weinan Zhang, Xinbing Wang, and Ying Wen. Mutual theory of mind in human-ai collaboration: An empirical study with llm-driven ai agents in a real-time shared workspace task. *arXiv preprint arXiv:2409.08811*, 2024b.
- X Zhang et al. A visual rag pipeline for few-shot fine-grained product classification. In *Proceedings* of the IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR), 2025.

Meng Zhou, Ziyu Liu, Pengwei Sui, Yixuan Li, and Yuk Ying Chung. Learning implicit credit assignment for cooperative multi-agent reinforcement learning. *Advances in neural information processing systems*, 33:11853–11864, 2020a.

Ming Zhou, Jun Luo, Julian Villella, Yaodong Yang, David Rusu, Jiayu Miao, Weinan Zhang, Montgomery Alban, Iman Fadakar, Zheng Chen, et al. Smarts: Scalable multi-agent reinforcement learning training school for autonomous driving. *arXiv preprint arXiv:2010.09776*, 2020b.