

# ProToM: Promoting Prosocial Behaviour via Theory of Mind-Informed Feedback

Matteo Bortoletto<sup>1,§</sup>, Yichao Zhou<sup>2,§</sup>, Lance Ying<sup>3,4</sup>, Tianmin Shu<sup>5</sup>, Andreas Bulling<sup>1</sup>

<sup>1</sup> University of Stuttgart   <sup>2</sup> Peking University   <sup>3</sup> Harvard University

<sup>4</sup> Massachusetts Institute of Technology   <sup>5</sup> Johns Hopkins University

<sup>§</sup> Work done while at Johns Hopkins University

Correspondence to: matteo.bortoletto@vis.uni-stuttgart.de

## Abstract

While humans are inherently social creatures, the challenge of identifying when and how to assist and collaborate with others – particularly when pursuing independent goals – can hinder cooperation. To address this challenge, we aim to develop an AI system that provides useful feedback to promote prosocial behaviour – actions that benefit others, even when not directly aligned with one’s own goals. We introduce *ProToM*, a Theory of Mind-informed facilitator that promotes prosocial actions in multi-agent systems by providing targeted, context-sensitive feedback to individual agents. ProToM first infers agents’ goals using Bayesian inverse planning, then selects feedback to communicate by maximising expected utility, conditioned on the inferred goal distribution. We evaluate our approach against baselines in two multi-agent environments: Doors, Keys, and Gems, as well as Overcooked. Our results suggest that state-of-the-art large language and reasoning models fall short of communicating feedback that is both contextually grounded and well-timed – leading to higher communication overhead and task speedup. In contrast, ProToM provides targeted and helpful feedback, achieving a higher success rate, shorter task completion times, and is consistently preferred by human users.<sup>1</sup>

## 1 Introduction

Prosocial behaviour is a broad class of actions that benefit other individuals (Kakulte and Shaikh, 2023) or society as a whole (Carattini and Roesti, 2023). However, humans do not always engage in prosocial behaviour because they may not know if others need help, or how they can help. As AI systems increasingly mediate human behaviour in everyday life, they inevitably shape not only individual outcomes but the broader social fabric.

<sup>1</sup>Code: <https://git.hcics.simtech.uni-stuttgart.de/public-projects/ProToM>

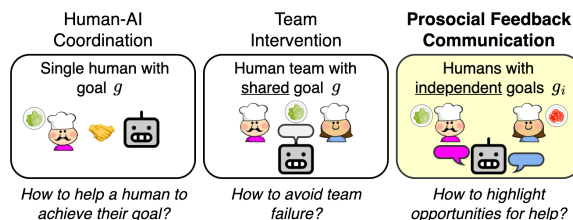


Figure 1: In contrast to the dominant paradigms of human-AI coordination and team intervention, we aim to encourage prosocial interactions among human agents pursuing independent goals. We introduce ProToM, a *facilitator* that promotes prosocial actions by providing targeted, context-sensitive feedback to individual agents.

Therefore, while much of the discussion around responsible AI focuses on preventing harm, an equally important goal is to actively promote social good, e.g., encouraging courtesy, reciprocity, or equitable resource sharing (Binns, 2018).

In this work, we argue that a promising yet under-explored direction is to leverage AI systems to encourage prosocial interactions among humans. This represents a shift from the dominant paradigm of human-AI coordination (Figure 1, left), where an AI assistant and a human operate as a joint unit with a shared goal (Carroll et al., 2019; Puig et al., 2023; Ying et al., 2024a; Zhi-Xuan et al., 2024; Chang et al., 2024) or to team intervention (Figure 1, middle) where the AI assists a team with a shared goal (Seo et al., 2023; Zhang et al., 2024b). By contrast, our proposed **prosocial facilitator paradigm** introduces an AI that operates as an observer and *facilitator* among multiple humans that pursue different goals (Figure 1, right). Our key research question is:

*Can we design an AI facilitator that infers humans’ beliefs and goals, detects opportunities for prosocial actions between them, and provides real-time feedback to promote those actions?*

Designing such a facilitator presents substantial challenges, as choosing when and what feedback to communicate highly depends both on the state of the environment and the agents’ internal states, such as beliefs and goals, which must be inferred only from available observations. Compared to human-AI coordination, an additional challenge arises from the need to model multiple agents, each with their own beliefs and goals, as well as how they interact with the environment and with one another based on those internal states.

We introduce ProToM<sup>2</sup> – a Theory of Mind-informed facilitator that observes agents in an environment and provides feedback in real-time to encourage prosocial behaviour. ProToM operates by first inferring the distribution over possible goals of each agent via Bayesian inverse planning. Conditioned on the inferred goal distributions and the current state of the environment, it then evaluates a set of possible feedback messages by computing their expected utility – measuring how helpful each feedback would be in guiding the agent toward prosocial actions that benefit another agent. It then selects the feedback with the highest expected utility and communicates it to the target agent only if the divergence between the agent’s simulated plan with or without communicating the feedback is large enough, i.e., if the agent is expected to behave differently in the absence of feedback. This ensures that ProToM delivers feedback only when it is both relevant and likely to be effective, resulting in timely and useful communication. Together with the feedback message, ProToM provides an explanation based on the inferred agents’ goals, helping the recipient understand not just what to do, but why the feedback is relevant and why it should be executed in a particular context.

We first conduct experiments with simulated human agents on two common multi-agent environments, Multi-Agent Doors, Keys, and Gems (Zhi-Xuan et al., 2024) and Overcooked (Carroll et al., 2019). Our results show that state-of-the-art large vision-language and reasoning models struggle to deliver meaningful feedback and tend to increase communication overhead. In contrast, ProToM consistently provides useful and well-timed feedback, resulting in a perfect task success rate, faster task completion, and reduced communication overhead. We further conducted a human study with real human participants. The results show that participants

perceived ProToM to have a stronger understanding of their own goals and to provide feedback that is more helpful, appropriate in both content and amount, and better explained.

In summary, we contribute: (1) a new human-AI paradigm for prosocial feedback communication; (2) a method, ProToM, that enables an AI facilitator to jointly infer agents’ mental states and generate ToM-informed feedback to promote prosocial behaviour among them; (3) a human study that evaluates the performance of models in real-time assistance with human participants, and humans’ perception of them.

## 2 Problem Formulation

We consider an AI facilitator as an omniscient observer that provides feedback to  $n$  interacting agents with the goal of promoting prosocial behaviour. Formally, we can formulate the problem as a tuple

$$\langle S, \{A^i\}_{i=1}^n, \mathcal{F}, T, \{\Omega^i\}_{i=1}^n, O, \{G^i\}_{i=1}^n \rangle \quad (1)$$

where:  $S$  is the state space,  $A^i$  is the action set of agent  $i \in \{1, \dots, n\}$ ,  $\mathcal{F}$  is the finite set of *feedback messages* the assistant can communicate,  $T(s' | s, \mathbf{a})$  is the state-transition kernel given joint action  $\mathbf{a} = (a^1, \dots, a^n) \in A^1 \times \dots \times A^n$ ,  $\Omega_i$  is the observation space of agent  $i$ ,  $O(o^i | s, f)$  is the observation kernel, and  $G^i$  is the set of possible goals of agent  $i$ . The feedback  $f \in \mathcal{F}$  is included in the observation kernel  $O$ , since feedback may convey extra information about the state that agents cannot directly observe.

**Agent Model** In partially observable environments, since agents do not know the true state  $s_t$ , they maintain a probability distribution over the possible states, called the *belief state*  $b_t$ . The probabilistic generative model for agents’ behaviour is defined as follows:

$$\begin{aligned} \text{Goal Prior:} & \quad g^i \sim P(g^i) \\ \text{Belief Update:} & \quad b_t^i \sim P(b_t^i | o_t^i, b_{t-1}^i) \\ \text{Action Selection:} & \quad a_t^i \sim P(a_t^i | b_t^i, g^i) \\ \text{State Transition:} & \quad s_{t+1} \sim T(s_{t+1} | s_t, \mathbf{a}_t^i, \mathbf{a}_t^{-i}) \end{aligned}$$

where  $\mathbf{a}_t^{-i}$  indicates other agents’ actions. To model planning, we assume that each agent selects actions  $a_t^i$  according to a policy  $\pi^i(a_t^i | b_t^i, g^i)$ , which may reflect approximate rationality (e.g., Boltzmann-rational planning) or be derived from

<sup>2</sup>All the code will be made public.

heuristic or rule-based strategies. In fully observable environments,  $b_t^i = s_t$ .

**Facilitator Model** Given an observed trajectory  $\tau = \{(s_t, \mathbf{a}_t)\}_{t=1}^T$ , the observer AI facilitator must decide whether to provide feedback, and if so, which  $f \in \mathcal{F}$  to communicate:  $f \sim P(f \mid \tau)$ . Since the agents’ internal state is not visible to the AI facilitator, the environment also is a POMDP from the facilitator’s perspective.

### 3 Method

ProToM actively assists both agents in an environment by observing their actions, inferring their goals, and providing feedback when it is expected to improve prosocial behaviour and thus overall performance. We present an overview of the method in Figure 2 and Algorithm 1.

#### 3.1 Maintaining Beliefs About Agents

To simulate and interpret the behaviour of agents, ProToM maintains a dynamic belief distribution  $b_F^i$  for each agent  $i$ , capturing both their internal belief about the environment and their underlying goal. This belief distribution is approximated using a particle filter with  $N$  particles per agent:

$$b_F^i = \{(b_{S,k}^i, b_{G,k}^i)\}_{k=1}^N, \quad (2)$$

where each particle  $k$  contains a sampled belief state  $b_{S,k}^i$  representing agent  $i$ ’s internal model of the environment, and a probability distribution  $b_{G,k}^i$  over possible goals for agent  $i$ , computed via Bayesian inverse planning.

**Agent Belief Sampling** While ProToM has access to the full environment state, it must infer agents’ goals from their actions under partial observability. To do this, for each particle  $k \in \{1, \dots, N\}$ , a sample of the agent  $i$ ’s belief state is drawn from a distribution conditioned on the agent’s current observation  $o_t^i$ :  $b_{S,k}^i \sim P(b_t^i \mid o_t^i)$ . At each timestep  $t$ , these sampled beliefs are updated using the environment’s state transition dynamics. To ensure that particles remain consistent with the agent’s actual observations, ProToM resamples any particle  $b_{S,k}^i$  that becomes incompatible with the latest observation  $o_t^i$ . This ensures that the facilitator’s estimate of the agents’ beliefs are consistent with what the agent could reasonably believe based on their observations.

**Goal Inference** Given a sampled belief state  $b_{S,k}^i$ , ProToM infers the agent’s goal using Bayesian inverse planning. Each particle maintains a goal distribution  $b_{G,k}^i(g^i)$ , which is updated by conditioning on the observed action  $a_t^i$  and the particle’s belief:

$$b_{G,k}^i(g^i) \propto b_{G,k}^i(g^i) \cdot P(a_t^i \mid g^i, b_{S,k}^i). \quad (3)$$

This update reflects the likelihood of the agent taking action  $a_t^i$  under goal  $g^i$ , assuming it holds belief  $b_{S,k}^i$ . ProToM’s overall belief about agent  $i$ ’s goal is then computed by averaging across all particles:

$$\hat{b}_G^i(g^i) = \frac{1}{N} \sum_{k=1}^N b_{G,k}^i(g^i). \quad (4)$$

#### 3.2 Feedback Selection

Given a finite set of candidate feedback messages, ProToM performs feedback selection in three steps. First, it evaluates each candidate feedback message by computing its expected utility, conditioned on the inferred goals and belief states of the agents. Second, ProToM decides whether to issue new feedback by comparing the expected utility of the best candidate against a threshold and estimating the divergence between the agent’s predicted behaviour with and without the feedback. This process ensures that feedback is only given when it is likely to meaningfully change agents’ interactions. Finally, if a feedback message is selected, ProToM generates a corresponding explanation to help the agent understand the reasoning behind the suggestion. This process is repeated at each timestep  $t$ , unless the agents are already executing feedback.

**Feedback Construction** The set of candidate feedback messages  $\mathcal{F}_t$  is constructed based on the current full environment state  $s_t$ , including object locations, agent positions, and affordances. This ensures that all feedback options are feasible. For example, if an agent cannot physically reach a particular item (e.g., due to immovable obstacles), we exclude feedback that would suggest picking up or passing that item.

**Utility Computation** Our utility function evaluates whether the feedback promotes prosocial actions that improve the overall task efficiency under sampled goal hypotheses from ProToM’s belief particles for all  $n$  agents (Algorithm 2). We denote these samples as:  $\mathbf{g} \sim \hat{\mathbf{b}}_G$  where

$$\mathbf{g} = (g^1, g^2, \dots, g^n), \quad \hat{\mathbf{b}}_G = (\hat{b}_G^1, \hat{b}_G^2, \dots, \hat{b}_G^n). \quad (5)$$

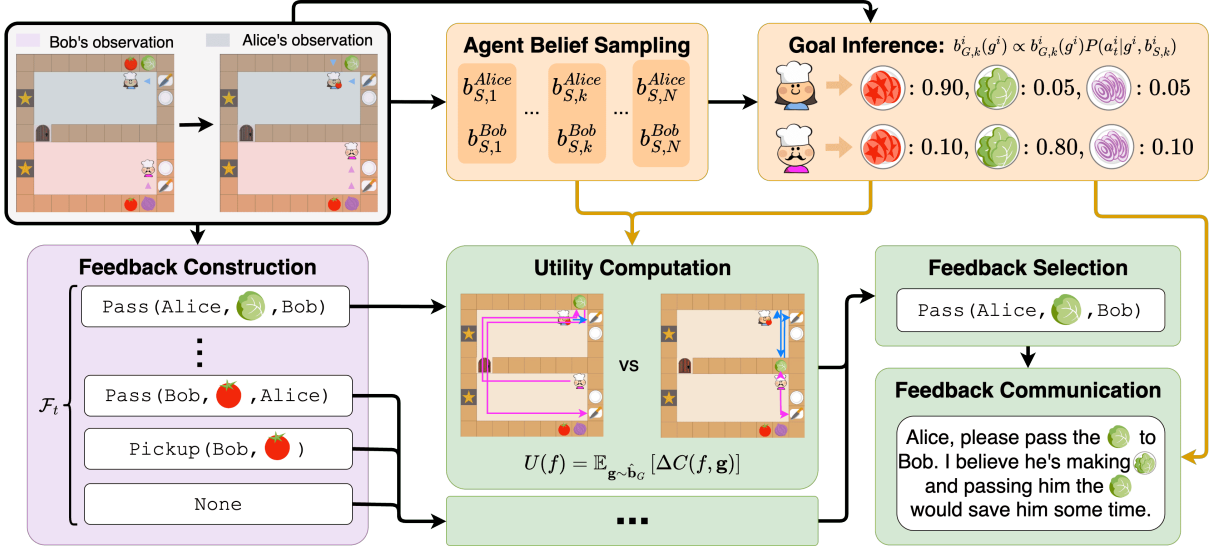


Figure 2: Overview of ProToM, a Theory of Mind-informed facilitator that promotes prosocial behaviour by communicating real-time feedback to agents pursuing independent goals. ProToM observes agents acting in a shared environment, infers their goal distributions via Bayesian inverse planning, and evaluates possible feedback messages based on expected utility. Feedback is paired with an explanation grounded in the inferred goals to clarify why the suggestion is relevant and helpful.

The utility of a feedback message  $f \in \mathcal{F}_t$  under  $\mathbf{g}$  is defined as:

$$U(f) = \mathbb{E}_{\mathbf{g} \sim \hat{\mathbf{b}}_G} [\Delta C(f, \mathbf{g})] \quad (6)$$

where we define the marginal improvement over not providing any feedback as:

$$\Delta C(f, \mathbf{g}) = C(\emptyset, \mathbf{g}) - C(f, \mathbf{g}). \quad (7)$$

Theoretically,  $C(f, \mathbf{g})$  can be any performance metric  $\ell(\cdot)$  that evaluates the agents' plans  $\pi = (\pi^1, \pi^2, \dots, \pi^n)$  generated under  $f$  and  $\mathbf{g}$ :

$$C(f, \mathbf{g}) = \ell(\pi \mid f, \mathbf{g}). \quad (8)$$

In our case,  $\ell(\cdot)$  corresponds to the number of steps to achieve  $\mathbf{g}$ , either with or without executing  $f$  beforehand. A positive utility  $U(f)$  indicates that providing feedback  $f$  is expected to help agents achieve their individual goals more efficiently – by decreasing the global number of steps required to reach task completion.

When a feedback message is targeted to agent  $i$ , it directly influences the policy  $\pi^i$  of that agent. We assume that agent  $i$  will follow the feedback upon receiving it, and continue to plan accordingly until either the feedback directive is completed, or it becomes infeasible. After that, the agent resumes planning toward its goal  $g^i$ .

**Feedback Selection** Among all candidates, we consider only feedback messages whose utility exceeds a threshold  $\phi$  and select those with the highest utility. The threshold  $\phi$  can be considered as the minimum expected benefit a feedback message must provide to even be considered. To avoid over-communicating and ensure that feedback is communicated at an appropriate time without interrupting the agent for negligible discrepancies, ProToM communicates the feedback only when the divergence from the agent's expected behaviour,  $\text{div}(f)$ , is greater than a threshold  $\epsilon$ :

$$\hat{\mathcal{F}}_t = \{f \mid U(f) = \max_{f' \in \mathcal{F}_t} U(f') > \phi, \text{div}(f) > \epsilon\}. \quad (9)$$

If  $\hat{\mathcal{F}}_t = \emptyset$ , ProToM skips communication in this timestep. If there is exactly one candidate in  $\hat{\mathcal{F}}_t$ , it is selected directly. If multiple feedback candidates satisfy these criteria, one is chosen uniformly at random:

$$\hat{f} \sim \text{Uniform}(\hat{\mathcal{F}}_t), \text{ if } \hat{\mathcal{F}}_t \neq \emptyset. \quad (10)$$

**Feedback Explanation** ProToM provides a natural-language explanation for each selected feedback by populating predefined templates  $\mathcal{T}$ . Specifically, each explanation  $e$  is determined by the agents' most probable goals in ProToM's belief particles  $\hat{\mathbf{g}}^{-i}$ , the selected feedback  $\hat{f}$ , and the

current state of the environment  $s_t$ :

$$e = \mathcal{T}(\hat{\mathbf{g}}^{-i}, \hat{f}, s_t), \quad \hat{\mathbf{g}}^{-i} = \arg \max_{\mathbf{g}^{-i}} \hat{\mathbf{b}}_G^{-i}(\mathbf{g}^{-i}). \quad (11)$$

The message highlights the inferred goals of other agents, and provides a concise reason for the feedback based on the current environment state (e.g. another agent may lack access to, or take much longer to reach a needed object), helping the recipient better understand the context and intent of the feedback.

## 4 Experiments

### 4.1 Environments

We evaluate our model against baselines in two popular multi-agent environments: Multi-Agent Doors, Keys, and Gems (Zhi-Xuan et al., 2024, mDKG) and Overcooked (Carroll et al., 2019). **mDKG** is fully observable, and the agents’ objective is to collect coloured gems, which are often located behind locked doors. Unlocking a door requires an agent to possess a key that matches the door’s colour. The possible feedback messages involve either unlocking a door or handing over a key. In **Overcooked** agents prepare and deliver meals. We adapt the environment introduced in (Wu et al., 2021) to a partially observable setting where agents can only observe the contents of the room they currently occupy. Feedback messages can either tell an agent to pass an item to the other agent, or to pick up a different but equivalent item – so the other agent can take the current one.

Importantly, we adapt these environments from their standard cooperative settings – where agents typically work together toward a shared goal – to a novel configuration in which each human pursues a separate objective. This distinction is crucial, as it allows us to investigate how the AI facilitator can promote prosocial behaviour in scenarios where agents are not explicitly required to collaborate. For each environment, we consider scenarios where (1) feedback is not needed, as agents can optimally complete their task by themselves; (2) feedback is useful to improve task completion efficiency; (3) feedback is essential to allow one of the agents to complete their task.

### 4.2 Models

**ProToM** For our experiments with mDKG, given that  $o_t^i = s_t$ , we use a single belief particle

( $N = 1$ ). As a divergence measure, we compute the probability that the two agents are already acting optimally according to a centralised planner  $\pi^*$ :  $\text{div} = P(\tau \in \pi^* \mid f)$  (Ullman et al., 2009). As thresholds, we set  $\phi = 0$  and  $\epsilon = 0.1$ . For **Overcooked**, to account for partial observability, we use  $N = 5$  belief particles per agent, and for **div** we computed the expected Jensen-Shannon divergence between the action distribution with or without  $f$ . As thresholds, we set  $\phi = 2$  and  $\epsilon = 0.3$ . These values were determined based on a small search on a held-out set (details in the Supplementary).

**Baselines** We compare ProToM against different baselines: agents that receive **no feedback**, an oracle version of ProToM that has access to agents’ ground truth goals (**ProToM-Oracle**), a random facilitator that at each timestep samples communicates a uniformly random message with probability 0.5 (**Random**), and various Large Vision-Language (**VLM**) and Reasoning Models (**RM**). Each model is provided with a description of the environment, including object types, agent action and observation space, and task rules. In addition, the model is given: an image of the current environment state, a history of recent actions, and the full set of candidate feedback messages  $\mathcal{F}_t$ . The model is then instructed to: (1) infer each agent’s goal based on the context, and (2) evaluate whether any message in  $\mathcal{F}_t$  would promote prosocial actions that positively impact task efficiency. If no feedback is deemed useful, the model is encouraged to respond with No Feedback. When assisting human participants, we also prompt the model to provide a short explanation of why it choose a specific feedback based on the inferred agent goals. We evaluate o3, o4-mini (OpenAI, 2025), GPT-4o (OpenAI, 2024), Claude 4 Sonnet (Anthropic, 2025), Gemini 2.5 Flash and Pro (Comanici et al., 2025), and Qwen 2.5 VL 72B (Bai et al., 2025). When possible, we set the model temperature to zero. The set of candidate feedback messages  $\mathcal{F}_t$  is constructed in the same way for all baselines to ensure a fair and consistent comparison.

### 4.3 Metrics

We use three metrics: success rate, speedup, and number of feedback messages. The success rate reflects the proportion of episodes in which both agents successfully complete their tasks. Speedup measures how much faster agents complete an episode with the help of a facilitator, compared

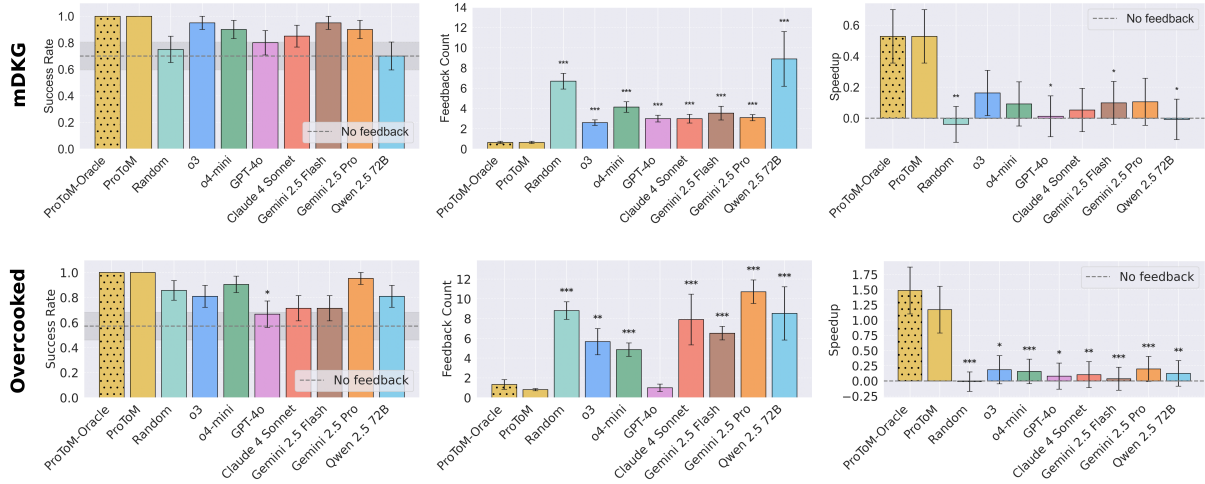


Figure 3: Simulation results comparing ProToM to baseline models on mDKG and Overcooked. ProToM achieves perfect success rates and significantly higher task speedup with minimal communication overhead. In contrast, baselines struggle to provide helpful, context-aware feedback and tend to over-communicate (\*\*\*:  $p < 0.001$ , \*\*:  $p < 0.01$ , \*:  $p < 0.05$ ).

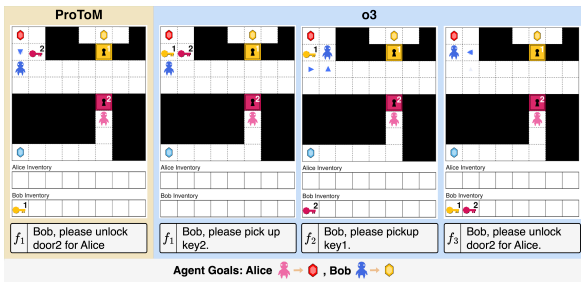


Figure 4: Example comparing ProToM and o3 on mDKG. While ProToM conveys the key feedback directly, o3 is overly detailed, resulting in the same outcome but with increased communication overhead.

to when they receive no feedback:  $\text{Speedup} = L_0/L_f - 1$ , where  $L_f$  and  $L_0$  are the average episode length with and without facilitator. The average number of feedback messages reflects how frequently a facilitator communicates with the agents – indicating whether it adopts a more parsimonious communication strategy or engages in frequent interaction.

#### 4.4 Simulation Experiments

We evaluate ProToM and the baselines with simulated human agents on both mDKG and Overcooked. In mDKG, we consider 20 scenarios: six where feedback is not needed, eight where feedback is useful, and six where feedback is necessary. The two human agents are simulated using the A\* planner from (Zhi-Xuan et al., 2024). In Overcooked, we evaluate 21 scenarios: seven

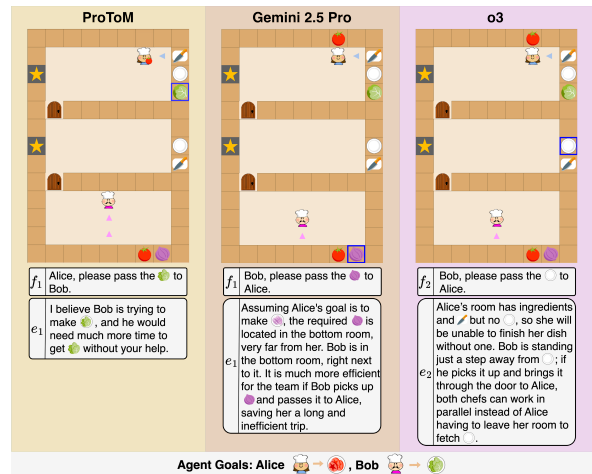


Figure 5: Example comparing ProToM, Gemini 2.5 Pro, and o3 on Overcooked. ProToM correctly infers Bob’s goal, instructing Alice to pass the ingredient he needs to save him some time. In contrast, o3 and Gemini fail to identify agents’ goals, suggesting Bob to pass irrelevant items. In this example, o3 also incorrectly states that Bob is close to a plate, suggesting a poor spatial understanding of the environment.

where feedback is not needed, seven where feedback is useful, and seven where feedback is necessary. Human agents are simulated using a stochastic heuristic-based planner.

**Results** Figure 3 shows average scores across all episodes. For statistical analysis, we used Fisher’s exact test for success rates and the Mann-Whitney U test for other metrics, applying Benjamini-Hochberg correction to control for false positives.

On mDKG, ProToM achieves a perfect success rate ( $1.00 \pm 0.00$ ), and performs on par with ProToM-Oracle, demonstrating effective goal inference. Strong reasoning models also achieve high success rates, e.g.  $0.95 \pm 0.05$  both for o3 and Gemini 2.5 Flash. However, ProToM achieves a higher average speedup of  $0.52 \pm 0.17$  compared to the no-feedback condition, while baselines struggle to improve task efficiency. Among the baselines, o3 performs the best with a speedup of  $0.16 \pm 0.14$ . Notably, ProToM is markedly more selective in its communication, issuing significantly fewer feedback messages than all other models ( $0.70 \pm 0.13$  on average). Figure 4 illustrates this qualitatively: while ProToM directly conveys the crucial feedback, o3 offers overly detailed feedback, increasing communication overhead.

Results follow similar trends in Overcooked. Here, ProToM achieves perfect success rate ( $1.00 \pm 0.00$ ) and  $1.17 \pm 0.38$  speedup, significantly outperforming all baselines – with the exception of ProToM-Oracle. Despite ProToM’s strong performance, its lower speedup and feedback count compared to ProToM-Oracle – though not statistically significant – highlight missed opportunities for useful communication, due to uncertainty about the agent’s beliefs and goals in partially observable settings. Among baselines, Gemini 2.5 Pro achieves the highest success rate ( $0.95 \pm 0.05$ ), followed by o4-mini ( $0.90 \pm 0.06$ ). However both show significantly lower speedups than ProToM:  $0.20 \pm 0.20$  and  $0.16 \pm 0.20$ , respectively. As in mDKG, ProToM communicates less frequently than all other models ( $0.81 \pm 0.11$  messages on average), with the sole exception of GPT-4o ( $1.00 \pm 0.37$ ). We find that VLMs and RMs do not always succeed in inferring agents goals, and therefore often communicate irrelevant feedback. Additionally, they sometimes demonstrate poor spatial understanding of the environment (see Figure 5).

**Ablations** To illustrate how ProToM achieves its performance, we ablate its two key components: (i) goal inference (NoGoal), where feedback is selected without conditioning on inferred goal distributions, and (ii) the divergence-based gating (NoDiv), where ProToM always communicates the highest-utility feedback without checking whether it is likely to change the recipient’s behaviour ( $\epsilon = 0$ ). Table 1 shows that both components matter, especially under partial observability. In mDKG, removing either goal reasoning or the

	mDKG			OC		
	SR	FC	SU	SR	FC	SU
NoGoal	1.00	0.85	0.47	0.85	0.48	0.35
NoDiv	1.00	0.75	0.46	1.00	2.19	0.86
ProToM	1.00	0.65	0.53	1.00	0.81	1.17

Table 1: Comparison between ProToM and ablations without goal inference (NoGoal) and divergence threshold (NoDiv). SR = success rate, FC = average feedback count, SU = speedup. Unlike ProToM, ablated versions do not achieve perfect SR and yield lower SU.

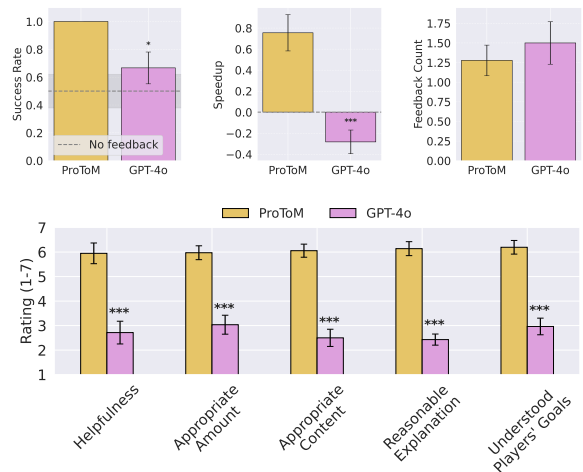


Figure 6: In our human study, ProToM significantly outperformed GPT-4o in success rate and task speedup, while requiring similar communication. Participants rated ProToM’s feedback as more helpful, appropriate, better explained and aligned with their goals (\*:  $p < 0.05$ , \*\*\*:  $p < 0.001$ ).

divergence criterion reduces speedup: compared to ProToM ( $SU = 0.53$  with  $0.65$  avg. messages), NoGoal and NoDiv achieve lower speedups ( $SU = 0.47$  and  $0.46$ ) while also increasing communication ( $FC = 0.85$  and  $0.75$ ). In Overcooked, goal inference becomes essential for reliability: NoGoal drops to  $SR = 0.85$  and yields much smaller gains ( $SU = 0.35$ ), indicating that feedback not grounded in agents’ inferred goals often fails to facilitate prosocial assistance, even causing task failures. Removing divergence gating preserves SR ( $1.00$ ) but substantially increases feedback volume ( $FC = 2.19$  vs  $0.81$ ) and still underperforms ProToM on speedup ( $SU = 0.86$  vs  $1.17$ ).

#### 4.5 Human Study

We conducted a human study in Overcooked with 18 participants grouped into nine pairs, testing three facilitator conditions: ProToM, a VLM, and

a control condition with no facilitator. Each participant pair played two trials per condition, across six scenarios: three where feedback was useful and three where it was necessary. Condition order was randomised to mitigate order effects. We selected GPT-4o as the LLM facilitator as the best trade-off between performance and response latency. We found that RMs’ longer inference times led to substantial delays at each game timestep (49s for o3, on average), which made experiments impractically long (see Supplementary). Participants were free to follow or ignore the feedback, allowing us to assess both the perceived value and reliability of the assistance. After each trial, participants rated the facilitator on four aspects: helpfulness of the feedback, appropriateness of its quantity and content, clarity of explanations, and whether they thought that the facilitator understood their goals. Responses were recorded on a 7-point Likert scale (1 = Strongly Disagree, 7 = Strongly Agree). Under this design, 18 participants provide 80% power to detect medium-to-large within-subject effects (paired Cohen’s  $d_z \approx 0.70$  at  $\alpha = 0.05$ ).

**Results** Our approach performed best in the human study, as shown in Figure 6. Compared to GPT-4o, ProToM achieved a significantly higher success rate ( $1.00 \pm 0.00$  vs  $0.66 \pm 0.11$ ), speedup ( $0.75 \pm 0.17$  vs  $-0.28 \pm 0.11$ ), while communicating a similar number of feedback messages ( $1.27 \pm 0.19$  vs  $1.50 \pm 0.27$ ). Notably, GPT-4o’s feedback negatively impacted task efficiency, suggesting that it often provided poorly timed or unhelpful feedback. We also found that human players ignored fewer messages from ProToM than GPT-4o (5 vs 10 in total), again suggesting that ProToM’s feedback was perceived as more relevant, clear, and therefore more trustworthy. This is further supported by subjective ratings, shown in Figure 6 (bottom), where participants rated ProToM significantly higher across all criteria.

#### 4.6 Discussion

The results across both simulated and human experiments paint a consistent picture: compared to other models, ProToM achieves higher success rates and task speedups, with reduced communication overhead (see Figure 3 and 6). Our human study also showed that ProToM’s feedback was perceived as more helpful, appropriate, and better explained and aligned with users’ goals. By comparison, state-of-the-art VLMs and RMs strug-

gle with communicating feedback that is useful to promote prosocial actions, often causing communication overheads. This performance gap in VLMs and RMs is due to two primary limitations of these models. First is the lack of strong Theory of Mind abilities: it is well known that existing models struggle to reliably infer mental states (Shapira et al., 2024). While modular methods like AutoToM (Zhang et al., 2025) show promise, their practical application in real-time settings is hindered by the substantial number of model calls required. Second, it is well known that current models face challenges in planning (Kambhampati et al., 2024). In our setup, these challenges are compounded as planning must incorporate inferred beliefs of multiple agents.

## 5 Related Work

Theory of Mind, the ability to attribute mental states to oneself and to others (Premack and Woodruff, 1978), has been widely studied in collaborative tasks, both using Bayesian (Ying et al., 2024a; Zhi-Xuan et al., 2024) and neural network approaches (Puig et al., 2023; Ying et al., 2024a; Zhang et al., 2024a; Bortoletto et al., 2024a,b; Cross et al., 2025; Ruhdorfer et al., 2025a). Several works focus on enhancing human-AI cooperation, where an AI assistant has to act in the environment to help a single human achieve their goal (Ying et al., 2024a; Puig et al., 2021, 2023; Buehler and Weisswange, 2020; Yu et al., 2024). Closer to our work is team assistance, where the AI assistant observes humans acting in the environment and provides feedback to them: Seo et al. (2023) use inverse reinforcement learning to instruct teams on which intent to follow; Zhang et al. (2024b) propose an assistant that, assuming a fixed plan, communicates to prevent failures or resolve deadlocks. In contrast, our novel paradigm supports sequential decision-making, relaxes the assumption that the agents share the same goal, and focuses on the broader objective of promoting prosocial actions, which can both avoid failure and improve efficiency.

## 6 Conclusion

We introduce a novel prosocial facilitator paradigm and ProToM, a method that combines Bayesian inverse planning to infer agents’ mental states and utility maximisation to select feedback that promotes prosocial behaviour among them. Our evalu-

ations show that while strong VLMs and RMs fall short of communicating effective feedback, ProToM’s is targeted and helpful, yielding a higher success rate and task speedup. Moreover, ProToM was consistently preferred by human users.

## Limitations

Our work comes with limitations. First, we have tested ProToM in simulated environments, which offer a good degree of control. mDKG and Overcooked represent a standard and widely accepted evaluation setup in multi-agent systems research (Carroll et al., 2019; Wu et al., 2021; Zhi-Xuan et al., 2024; Ying et al., 2024a; Ruhdorfer et al., 2025c,b), making it appropriate both for establishing our novel framework and for conducting controlled human studies. Evaluating ProToM in real-world settings in an interesting research direction that we aim to explore in future work. Research in cognitive science provides good evidence that people construct small state and goal spaces for planning and inference (Brooke-Wilson, 2023). In this view, Bayesian inverse planning is scalable to real-settings, as long as we retain only the relevant states and goals at any given time. This can be done using an LLM (Wong et al., 2025; Ying et al., 2024b). LLMs can also be used to amortize the policy when performing Bayesian inverse planning (Jin et al., 2024). Second, ProToM’s feedback doesn’t explicitly account for pragmatics. Future work could design a version of ProToM that not only selects what feedback to communicate but also how to convey it based on agents’ inferred beliefs and affordances.

Further extensions could model deeper recursive reasoning between agents, and improve scalability through approximate or learned inference methods. More broadly, integrating ProToM-style architectures with modern large foundation models may help bridge the gap between strong empirical performance in controlled environments and effective prosocial facilitation in real-world multi-agent systems.

## Acknowledgments

T. Shu acknowledges a grant from Amazon.

## References

Anthropic. 2025. Introducing claude 4. <https://www.anthropic.com/news/claude-4>.

Shuai Bai, Keqin Chen, Xuejing Liu, Jialin Wang, Wenbin Ge, Sibao Song, Kai Dang, Peng Wang, Shijie Wang, Jun Tang, and 1 others. 2025. Qwen2.5-v1 technical report. *arXiv preprint arXiv:2502.13923*.

Reuben Binns. 2018. Fairness in machine learning: Lessons from political philosophy. In *Conference on fairness, accountability and transparency*, pages 149–159. PMLR.

Matteo Bortoletto, Constantin Ruhdorfer, Adnen Abdesaied, Lei Shi, and Andreas Bulling. 2024a. Limits of theory of mind modelling in dialogue-based collaborative plan acquisition. In *Proc. 62nd Annual Meeting of the Association for Computational Linguistics (ACL)*, pages 1–16.

Matteo Bortoletto, Constantin Ruhdorfer, Lei Shi, and Andreas Bulling. 2024b. Explicit modelling of theory of mind for belief prediction in nonverbal social interactions. In *Proc. 27th European Conference on Artificial Intelligence (ECAI)*, pages 866–873.

Tyler Brooke-Wilson. 2023. Bounded rationality as a strategy for cognitive science. *PhD Diss., Massachusetts Institute of Technology*.

Moritz C Buehler and Thomas H Weisswange. 2020. Theory of mind based communication for human agent cooperation. In *2020 IEEE international conference on human-machine systems (ICHMS)*, pages 1–6. IEEE.

Stefano Carattini and Matthias Roesti. 2023. Trust, happiness, and pro-social behavior. *Review of Economics and Statistics*, pages 1–45.

Micah Carroll, Rohin Shah, Mark K Ho, Tom Griffiths, Sanjit Seshia, Pieter Abbeel, and Anca Dragan. 2019. On the utility of learning about humans for human-ai coordination. *Advances in neural information processing systems*, 32.

Matthew Chang, Gunjan Chhablani, Alexander Clegg, Mikael Dallaire Cote, Ruta Desai, Michal Hlavac, Vladimir Karashchuk, Jacob Krantz, Roozbeh Motlaghi, Priyam Parashar, and 1 others. 2024. Partnr: A benchmark for planning and reasoning in embodied multi-agent tasks. *arXiv preprint arXiv:2411.00081*.

Gheorghe Comanici, Eric Bieber, Mike Schaeckermann, Ice Pasupat, Noveen Sachdeva, Inderjit Dhillon, Marcel Blistein, Ori Ram, Dan Zhang, Evan Rosen, and 1 others. 2025. Gemini 2.5: Pushing the frontier with advanced reasoning, multimodality, long context, and next generation agentic capabilities. *arXiv preprint arXiv:2507.06261*.

Logan Cross, Violet Xiang, Agam Bhatia, Daniel LK Yamins, and Nick Haber. 2025. Hypothetical minds: Scaffolding theory of mind for multi-agent tasks with large language models. In *The Thirteenth International Conference on Learning Representations*.

- Marco F. Cusumano-Towner, Feras A. Saad, Alexander K. Lew, and Vikash K. Mansinghka. 2019. [Gen: A general-purpose probabilistic programming system with programmable inference](#). In *Proceedings of the 40th ACM SIGPLAN Conference on Programming Language Design and Implementation, PLDI 2019*, pages 221–236, New York, NY, USA. ACM.
- Chuangyang Jin, Yutong Wu, Jing Cao, Jiannan Xiang, Yen-Ling Kuo, Zhiting Hu, Tomer Ullman, Antonio Torralba, Joshua Tenenbaum, and Tianmin Shu. 2024. Mmtom-qa: Multimodal theory of mind question answering. In *Proceedings of the 62nd Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*, pages 16077–16102.
- Arvind Kakulte and Samreen Shaikh. 2023. Prosocial behavior, psychological well-being, positive and negative affect among young adults: A cross-sectional study. *Industrial Psychiatry Journal*, 32(Suppl 1):S127–S130.
- Subbarao Kambhampati, Karthik Valmeekam, Lin Guan, Mudit Verma, Kaya Stechly, Siddhant Bhambri, Lucas Paul Saldyt, and Anil B Murthy. 2024. [Position: LLMs can’t plan, but can help planning in LLM-modulo frameworks](#). In *Forty-first International Conference on Machine Learning*.
- OpenAI. 2024. Gpt-4o system card. <https://openai.com/index/gpt-4o-system-card/>.
- OpenAI. 2025. Introducing openai o3 and o4-mini. <https://openai.com/index/introducing-o3-and-o4-mini/>.
- David Premack and Guy Woodruff. 1978. Does the chimpanzee have a theory of mind? *Behavioral and brain sciences*, 1(4):515–526.
- Xavier Puig, Tianmin Shu, Shuang Li, Zilin Wang, Yuan-Hong Liao, Joshua B. Tenenbaum, Sanja Fidler, and Antonio Torralba. 2021. [Watch-and-help: A challenge for social perception and human-ai collaboration](#). In *International Conference on Learning Representations*.
- Xavier Puig, Tianmin Shu, Joshua B Tenenbaum, and Antonio Torralba. 2023. Nopa: Neurally-guided online probabilistic assistance for building socially intelligent home assistants. In *2023 IEEE International Conference on Robotics and Automation (ICRA)*, pages 7628–7634. IEEE.
- Constantin Ruhdorfer, Matteo Bortoletto, and Andreas Bulling. 2025a. The yokai learning environment: Tracking beliefs over space and time. *arXiv preprint arXiv:2508.12480*.
- Constantin Ruhdorfer, Matteo Bortoletto, Victor Oei, Anna Penzkofer, and Andreas Bulling. 2025b. Unsupervised partner design enables robust ad-hoc teamwork. *arXiv preprint arXiv:2508.06336*.
- Constantin Ruhdorfer, Matteo Bortoletto, Anna Penzkofer, and Andreas Bulling. 2025c. The overcooked generalisation challenge: Evaluating cooperation with novel partners in unknown environments using unsupervised environment design. *Transactions on Machine Learning Research*.
- Sangwon Seo, Bing Han, and Vaibhav Unhelkar. 2023. Automated task-time interventions to improve teamwork using imitation learning. In *Proceedings of the 2023 International Conference on Autonomous Agents and Multiagent Systems*, pages 335–344.
- Natalie Shapira, Mosh Levy, Seyed Hossein Alavi, Xuhui Zhou, Yejin Choi, Yoav Goldberg, Maarten Sap, and Vered Shwartz. 2024. Clever hans or neural theory of mind? stress testing social reasoning in large language models. In *18th Conference of the European Chapter of the Association for Computational Linguistics, EACL 2024*, pages 2257–2273. Association for Computational Linguistics (ACL).
- Tomer Ullman, Chris Baker, Owen Macindoe, Owain Evans, Noah Goodman, and Joshua Tenenbaum. 2009. Help or hinder: Bayesian models of social goal inference. *Advances in neural information processing systems*, 22.
- Lionel Wong, Katherine M Collins, Lance Ying, Cedegao E Zhang, Adrian Weller, Tobias Gerstenberg, Timothy O’Donnell, Alexander Lew, Jacob Andreas, Tyler Brooke-Wilson, and 1 others. 2025. Modeling open-world cognition as on-demand synthesis of probabilistic models. In *Proceedings of the Annual Meeting of the Cognitive Science Society*, volume 47.
- Sarah A Wu, Rose E Wang, James A Evans, Joshua B Tenenbaum, David C Parkes, and Max Kleiman-Weiner. 2021. Too many cooks: Bayesian inference for coordinating multi-agent collaboration. *Topics in Cognitive Science*, 13(2):414–432.
- Lance Ying, Kunal Jha, Shivam Aarya, Joshua B Tenenbaum, Antonio Torralba, and Tianmin Shu. 2024a. Goma: Proactive embodied cooperative communication via goal-oriented mental alignment. In *2024 IEEE/RSJ International Conference on Intelligent Robots and Systems (IROS)*, pages 7099–7106. IEEE.
- Lance Ying, Xinyi Li, Shivam Aarya, Yizirui Fang, Yifan Yin, Jason Xinyu Liu, Stefanie Tellex, Joshua B Tenenbaum, and Tianmin Shu. 2024b. Pragmatic embodied spoken instruction following in human-robot collaboration with theory of mind. *arXiv preprint arXiv:2409.10849*.
- Chuang Yu, Baris Serhan, and Angelo Cangelosi. 2024. Top-tom: Trust-aware robot policy with theory of mind. In *2024 IEEE International Conference on Robotics and Automation (ICRA)*, pages 7888–7894. IEEE.
- Ceyao Zhang, Kaijie Yang, Siyi Hu, Zihao Wang, Guanghe Li, Yihang Sun, Cheng Zhang, Zhaowei Zhang, Anji Liu, Song-Chun Zhu, and 1 others. 2024a. Proagent: building proactive cooperative

agents with large language models. In *Proceedings of the AAAI Conference on Artificial Intelligence*, volume 38, pages 17591–17599.

Yuening Zhang, Paul Robertson, Tianmin Shu, Sungkweon Hong, and Brian C Williams. 2024b. Risk-bounded online team interventions via theory of mind. In *2024 IEEE International Conference on Robotics and Automation (ICRA)*, pages 12964–12970. IEEE.

Zhining Zhang, Chuanyang Jin, Mung Yao Jia, and Tianmin Shu. 2025. Autotom: Automated bayesian inverse planning and model discovery for open-ended theory of mind. *arXiv preprint arXiv:2502.15676*.

Tan Zhi-Xuan. 2022. *Pddl.jl: An extensible interpreter and compiler interface for fast and flexible ai planning*. Ph.D. thesis, Massachusetts Institute of Technology.

Tan Zhi-Xuan, Lance Ying, Vikash Mansinghka, and Joshua B Tenenbaum. 2024. Pragmatic instruction following and goal assistance via cooperative language-guided inverse planning. In *AAMAS*.

## A Environments

### A.1 Additional Environment Details

**mDKG** (Zhi-Xuan et al., 2024) is a fully observable, 2D grid-world where the objective is to collect coloured gems, which are often located behind locked doors. Unlocking a door requires an agent to possess a key that matches the door’s colour. Agents can move in the four cardinal directions or remain stationary, while respecting movement constraints imposed by walls and locked doors. In addition to navigation, agents can pick up items, hand them over to other agents, or unlock doors when they possess a matching key. The possible feedback messages involve either unlocking a door or handing over a key. We evaluate 20 scenarios: six where feedback is not needed, eight where feedback is useful, and six where feedback is necessary. We set the maximum number of timesteps to  $T_{max} = 80$ .

**Overcooked** (Carroll et al., 2019) is a 2D grid-based kitchen in which agents prepare and deliver meals. The kitchen consists of counters that can hold either movable items, such as food ingredients and plates, or fixed stations, such as knife stations for chopping and delivery stations for serving completed dishes. Agents can move in the four cardinal directions or stay still, pick up and put down items, carry food to knife stations for chopping, and combine food with plates to assemble meals. Each agent is limited to carrying a single object at a time and cannot directly hand items to other

agents: items must be placed on shared counters for indirect transfer. We adapt the environment introduced in (Wu et al., 2021) to a partially observable setting by adding doors that divide the environment into separate rooms. Agents can only observe the contents of the room they currently occupy. The possible recipes to prepare include SimpleTomato, SimpleLettuce, and SimpleOnion – as shown in Figure A4. Each recipe follows the same three-step process:

1. Pick up the main fresh ingredient (e.g., FreshTomato).
2. Chop it at the chopping station, producing the chopped version (e.g., ChoppedTomato).
3. Combine the chopped ingredient with a plate (e.g., ChoppedTomato-Plate) and deliver it at the delivery station.

We evaluate 21 scenarios: seven where feedback is not needed, seven where feedback is useful, and seven where feedback is necessary. We set the maximum number of timesteps to  $T_{max} = 100$ .

### A.2 Feedback Types

We define a discrete set of feedback types that capture context-sensitive opportunities to promote prosocial behaviour among agents. These types are instantiated differently in each domain based on their different environment dynamics and available actions.

In mDKG:

- `Unlock(agenti, doork, agentj)`: Suggests agent<sub>i</sub> to unlock door<sub>k</sub> for agent<sub>j</sub>. If agent<sub>i</sub> doesn’t already hold a key to unlock door<sub>k</sub>, the feedback assumes that agent<sub>i</sub> will first pick it up.
- `Handover(agenti, keyk, agentj)`: Suggests agent<sub>i</sub> to hand key<sub>k</sub> over to agent<sub>j</sub>. If agent<sub>i</sub> doesn’t already key<sub>k</sub>, the feedback assumes that agent<sub>i</sub> will first pick it up.

In Overcooked:

- `Pass(agenti, itemk, agentj)`: Suggests agent<sub>i</sub> to pass a specific item<sub>k</sub> at a specified location to another agent across a counter (in Overcooked, agents cannot directly handover items to other agents). Since agents can only hold one item at the time, if agent<sub>i</sub> is currently holding an item  $\neq$  item<sub>k</sub>, they are first



Figure A1: SimpleLettuce



Figure A2: SimpleTomato



Figure A3: SimpleOnion

Figure A4: Possible recipes in Overcooked.

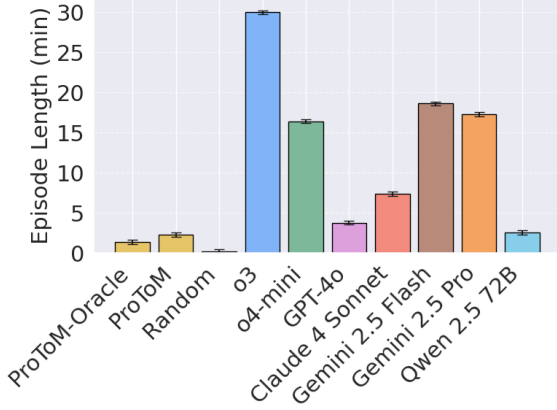


Figure A5: Average episode length in Overcooked, in minutes.

instructed to place it on the nearest available counter. Then, they pick up the target item and pass it to a shared counter that is not on the layout border and is accessible to both agents.

- Pickup(agent<sub>i</sub>, item<sub>k</sub>): Suggests agent<sub>i</sub> to pick up a specific item at a specified location. This feedback is primarily used to resolve conflicts where both agents are attempting to pick up the same item, whereas one of them may have access to an alternative.

## B Models and Implementation

### B.1 ProToM Parameters

For our experiments with mDKG, given that agents have full observability ( $o_t^i = s_t$ ), we use a single belief particle ( $N = 1$ ). For Overcooked, to account for partial observability, we use  $N = 5$  belief particles per agent. While setting  $N = 1$  for mDKG is an obvious choice, the number of particles for Overcooked is determined based on the resulting speed of the model inference. Given that we use ProToM to assist human agents in real-time, we opted for  $N = 5$ . As thresholds, we set ( $\phi = 0, \epsilon = 0.1$ ) on mDKG, and ( $\phi = 2, \epsilon = 0.3$ ) on Overcooked. The values for the thresholds  $\phi$  and  $\epsilon$  were determined based on a small search on a held-out set, with values from 0 to 1.0 increasing

---

### Algorithm 1 ProToM

---

**Require:** Goal set  $\mathcal{G} = G^1 \times \dots \times G^n$ , planners  $\pi = (\pi^1, \dots, \pi^n)$ , performance metric  $\ell(\cdot)$ , thresholds  $\phi, \epsilon$

- 1: **Initialise:** belief particles  $\{\{(b_{S,k}^i, b_{G,k}^i)\}_{k=1}^N\}_{i=1}^n$  based on initial state  $s_0, t \leftarrow 1$
- 2: **while**  $t = T_{max}$  or Done **do**
- 3: Observe state  $s_t$  and agent action  $\mathbf{a}_t = (a_t^1, \dots, a_t^n)$
- 4: **for** agent  $i \in \{1, \dots, n\}$  **do**
- 5: Get agent observation  $o_t^i$  from  $s_t$
- 6: **for** belief particle  $k \in \{1, \dots, N\}$  **do**
- 7:  $b_{S,k}^i \leftarrow \text{BELIEFUPDATE}(b_{S,k}^i, o_t^i, \mathbf{a}_t)$
- 8:  $b_{G,k}^i \leftarrow \text{GOALUPDATE}(b_{G,k}^i, b_{S,k}^i, G^i, a_t^i)$
- 9: **end for**
- 10:  $\hat{b}_G^i \leftarrow \frac{1}{N} \sum_{k=1}^N b_{G,k}^i$
- 11: **end for**
- 12:  $\mathcal{F}_t \leftarrow \text{CONSTRUCTFEEDBACK}(s_t)$
- 13:  $U(f) \leftarrow \text{COMPUTEUTILITY}(s_t, \mathcal{G}, \hat{\mathbf{b}}_G, \mathcal{F}_t, \pi, \ell)$
- 14:  $\hat{\mathcal{F}}_t \leftarrow \{f \in \mathcal{F}_t \mid U(f) = \max_{f' \in \mathcal{F}_t} U(f') > \phi, \text{div}(f) > \epsilon\}$
- 15: **if**  $\hat{\mathcal{F}}_t \neq \emptyset$  **then**
- 16:  $f \sim \text{Uniform}(\hat{\mathcal{F}}_t)$
- 17:  $e \leftarrow \text{EXPLANATION}(\mathcal{G}, \hat{\mathbf{b}}_G, s_t, f)$
- 18: Communicate feedback  $f$  with explanation  $e$
- 19: **end if**
- 20:  $t \leftarrow t + 1$
- 21: **end while**

---

by 0.1 at each search. For all our experiments, we set the random seed to 42.

### B.2 ProToM Explanation Templates

As discussed in the main paper, ProToM generates natural-language explanations for each selected feedback by instantiating predefined templates based on the inferred goal of the other agent and the current environment state. These templates are designed to provide concise and context-aware justifications for the selected feedback.

For the Pass(agent<sub>i</sub>, item<sub>k</sub>, agent<sub>j</sub>) feedback, the explanation template is structured as: Template: I believe [agent<sub>j</sub>] is trying to prepare [recipe], and...

- If agent<sub>j</sub> can access the item but it would take time to get it: ... [he/she] would need much more time to get the [item<sub>k</sub>] without your help.
- If agent<sub>j</sub> cannot access the item: ... without your help, [agent<sub>j</sub>\_article]

---

**Algorithm 2** COMPUTEUTILITY

---

**Require:** Current state  $s_t$ , goal set  $\mathcal{G} = G^1 \times \dots \times G^n$ , inferred goal probability distributions  $\hat{\mathbf{b}}_G = (\hat{b}_G^1, \dots, \hat{b}_G^n)$ , candidate feedback set  $\mathcal{F}_t$ , planners  $\pi = (\pi^1, \dots, \pi^n)$ , performance metric  $\ell(\cdot)$

- 1: Initialise  $U(f) \leftarrow 0$  for all  $f \in \mathcal{F}_t$
- 2: **for** each  $\mathbf{g} = (g^1, \dots, g^n) \in \mathcal{G}$  **do**
- 3:    $C(\emptyset, \mathbf{g}) \leftarrow \ell(\pi | \mathbf{g})$
- 4:   **for** each feedback  $f \in \mathcal{F}_t$  **do**
- 5:      $C(f, \mathbf{g}) \leftarrow \ell(\pi | f, \mathbf{g})$
- 6:      $\Delta C \leftarrow C(\emptyset, \mathbf{g}) - C(f, \mathbf{g})$
- 7:      $U(f) \leftarrow U(f) + \Delta C \cdot \prod_{i=1}^n \hat{b}_G^i(g^i)$
- 8:   **end for**
- 9: **end for**

---

wouldn't be able to get the [item<sub>k</sub>].

For the Pickup(agent<sub>i</sub>, item<sub>k</sub>) feedback, the explanation template is: I believe [agent<sub>j</sub>] is trying to prepare [recipe], and the other [item<sub>k</sub>] is easier to get for [him/her].

### B.3 LLM Details

We use the following model versions:

- o3-2025-04-16 (OpenAI API)
- o4-mini-2025-04-16 (OpenAI API)
- gpt-4o-2024-08-06 (OpenAI API)
- claude-sonnet-4-20250514 (Anthropic API)
- gemini-2.5-flash (Google API)
- gemini-2.5-pro (Google API)
- qwen2.5-vl-72b-instruct (OpenRouter API)

### B.4 LLM Prompts

We show the prompts used for evaluating large VLMs and RMs on mDKG and Overcooked in Figure A7 and Figure A8, respectively. Each model is provided with a description of the environment, including object types, agent action and observation space, and task rules. In addition, the model is given: an image of the current environment state, a history of recent actions, and the full set of candidate feedback messages  $\mathcal{F}_t$ . The model is then instructed to: (1) infer each agent's goal based on the context, and (2) evaluate whether any message in  $\mathcal{F}_t$  would promote prosocial actions that positively impact task efficiency. If no feedback is deemed useful, the model is encouraged to respond with No



Figure A6: Example in which ProToM fails to communicate feedback, and ProToM-Oracle succeeds.

Feedback. When assisting human participants, we also prompt the model to provide a short explanation of why it choose a specific feedback based on the inferred agent goals. We evaluate o3, o4-mini (OpenAI, 2025), GPT-4o (OpenAI, 2024), Claude 4 Sonnet (Anthropic, 2025), Gemini 2.5 Flash, and Gemini 2.5 Pro (Comanici et al., 2025). When possible, we set the model temperature to zero (RMs often do not allow to control this parameter).

## C Human Study

### C.1 Setup and Interface

Figure A9 shows the interface used for the human study. Each player could see the codename of the facilitator that was assisting them, their own goal recipe, their partial observation of the environment, the possible actions to take, and a legend. When a player received a feedback message, it would appear as a pop-up blue box. Players could decide to ignore feedback messages if they considered it useless or unnecessary. At the end of each episode, both players were revealed the other player's goal recipe, information about episode completion, and the full history of feedback messages given to each player, with their status ("Completed", "Ignored", "Not Executable"). We revealed this information to let both players fill out the questionnaires, even if one of the two players did not receive any feedback message during the episode.

## C.2 Participant Details and Procedure

We recruited 18 human participants: 5 female, 13 male, aged between 19 and 31 years old. The study was approved by the institutional ethics committee. Some participants were university students who received a compensation, in accordance with university regulations. The remaining participants voluntarily joined the study, without receiving any form of compensation. At the beginning of the study, participants were informed about their task, the duration of the experiment, and that their responses would be kept anonymous and used solely for research purposes. They then went through a guided tutorial that explains the rules of the game. Each full study – i.e. two participants playing together the six trials – lasted around one hour. The data we collect is anonymised and does not contain any information that names or uniquely identifies individual people.

## C.3 API Call Times and Justification for Model Choice

As we report in the main text, we selected GPT-4o as the LLM facilitator for the human study, given the good trade-off between performance and response latency. We found that RMs’ longer inference times led to substantial delays at each game timestep, which made experiments impractically long. We show the average episode length (in minutes) for simulated experiments in Overcooked in Figure A5. As one can see from the figure, models like o3, o4 or Gemini make episode trials drastically longer compared to ProToM. Therefore, using one of these models for our human experiments would: 1) be frustrating for the participants; 2) bias or distract them; 3) result in much longer study durations. Figure A5 suggests two possible alternatives to GPT-4o: Claude 4 Sonnet and Qwen 2.5 72B. However, by looking at Figure 3 (bottom middle) in the main text, we see that both Claude 4 Sonnet and Qwen 2.5 72B engage in lots of feedback communication. In preliminary experiments we noticed that such amount of messages would overwhelm the participants, that would tend to just ignore all of them. Therefore, we concluded that the best compromise is GPT-4o, which we ended up using in the human study.

## D Additional Experimental Results

### D.1 Quantitative Analysis

For both mDKG and Overcooked, we consider different scenarios where (1) feedback is not needed, as agents can optimally complete their task by themselves; (2) feedback is useful to improve task completion efficiency; (3) feedback is necessary to allow one of the agents to complete their task. While in the main text we provide scores averaged across all scenarios, here we report average scores across each type of scenario in Figure A10 for mDKG, and Figure A11 for Overcooked.

In mDKG, all as expected: feedback is communicated only in episodes where it is needed or necessary. ProToM is on par with its oracle version (ProToM-Oracle), suggesting strong goal inference.

In Overcooked, ProToM-Oracle performs better – although not statistically different – than ProToM both when feedback is necessary and useful, while communicating the same number of feedback messages (feedback necessary:  $1.12 \pm 0.12$ ; feedback useful:  $1.00 \pm 0.00$ ). This highlights that there are scenarios where ProToM does not communicate, or its feedback is not optimal or is communicated at a non-optimal timestep. We show a qualitative example in which ProToM fails in communicating feedback in Figure A6. On the other hand, when feedback is not needed, ProToM-Oracle communicates more feedback than ProToM, which yields a decrease in task speedup. This is due to the fact that, despite knowing the ground truth goals, ProToM-Oracle is still not a perfect model, as there is no ground truth values for the parameters  $\phi$  and  $\epsilon$ . For ProToM-Oracle, we used the same values as in ProToM. We expect that performing a finer grid search on a bigger training set could mitigate this issue.

## E Societal Implications

While we design our experiments in controlled environments, deploying future, more advanced versions of ProToM in real-world settings could raise a few practical societal concerns. Because it relies on inferring what people mental states, it is important to avoid reinforcing biases or misunderstanding psychological nuances. Finally, depending on how much data is collected to make these inferences, there could be privacy concerns: real-world use would benefit from clear consent, data minimization, and transparency.

## **F Infrastructure and Code**

### **F.1 Compute Resources**

We ran our model on a server running Ubuntu 22.04, Intel Xeon Platinum 8260 CPUs for a total of 96 cores. Proprietary models are used through API.

### **F.2 Code**

Our code is public under the MIT license at <https://git.hcics.simtech.uni-stuttgart.de/public-projects/ProToM>. The code for Overcooked is written in Python 3.10 and adapted from the code released by (Wu et al., 2021). The code for mDKG is written in Julia 1.11.5, and based on the Gen.jl (Cusumano-Towner et al., 2019) and PDDL.jl (Zhi-Xuan, 2022) libraries.

```

Consider a multi-agent Doors, Keys, and Gems environment described as follows:

# Entities:
- Agents: Alice is the pink human avatar, Bob is the blue human avatar. Their inventory is shown at the bottom of the image.
- Keys have colored key shapes, gems are colored hexagons.
- Doors are colored cells with a black keyhole icon.
- Walls are black cells and cannot be entered.
As shown in the image, doors and keys are numbered. The number is just for referencing the keys in the text. Keys open doors of their same color,
independently from the number.

# Properties:
- Each agent occupies a cell. Agents can occupy the same cell.
- Each item (key or gem) occupies a cell unless it is carried or off-grid.
- Each door occupies a cell and can be unlocked with a key of the same color.
- Each key and door has a color; keys can unlock doors of matching color.
- Agents act by taking turns.

# Goals:
- Each agent wants to collect a specific gem among the ones available in the environment.

# Observability:
Agents know:
- The entire grid layout (walls, doors, items, agents).
- The locked/unlocked state of doors.
- The location of agents and items.
- Their own possession of items.
- Which agent's turn it is.

# Agent Actions:
1. Move
- Agents can move up, down, left, or right, one cell at a time, if:
- There is no wall in that direction.
- There is no locked door blocking the move.

2. Pickup an item
- If an agent is on the same cell as an item, they can pick it up.
- Picking up an item removes it from the grid and adds it to the agent's inventory.

3. Unlock a door
- If an agent is adjacent to a locked door and has a key with the same color as the door, they can unlock it.
- The key is consumed in the process.

4. Handover
- If two agents are adjacent, one can give an item they hold to the other.

5. Wait
- The agent can simply pass their turn.

## Turn order:
- Only one agent is active at a time.
- After an action, the next agent becomes active.

---

# Your Task as a Social Facilitator

At each timestep, your task is to decide whether to give feedback to either agent, and what that feedback should be, in order to promote
**prosocial behavior**.

A prosocial agent:
- Helps the other agent if doing so makes it possible or more efficient to achieve their goals.
- Takes action that indicates awareness of the other's goal.

# Step-by-Step Instructions

1. Review the History
You will be given a few recent timesteps, including agent actions.

2. Infer the Agents' Goals
From item movements and paths, infer which gem each agent is targeting.

3. Evaluate if it is possible to promote prosocial behaviour
Given the action history, the current state and the inferred goal for each agent and a set of possible feedbacks, is any of the feedbacks useful
for promoting prosocial behaviour? If not, just answer "No feedback".

4. Decide Feedback
If you think that it is useful to provide a feedback given the action history and current state, then choose one feedback that (1) enables both
agents to achieve their goals and/or (2) increases overall efficiency.
Typical examples of prosocial agents are agents that take actions to unlock doors that block other agents, or that handover keys to an agent that
would need many turns to get by themselves.
Avoid giving feedback just because an action isn't optimal; only intervene when it clearly improves overall efficiency or agent success. If
this is not the case, just choose "No feedback".

## Recent Action History
$history

## Current state
The current state of the environment is shown in the image.

Based on the instructions provided above, choose which feedback to give from the possible options:
$feedback_options_str

The first element of the tuple is the target of the feedback, while the second is the feedback itself.
For example "(bob, has(alice, key1))" means "tell Bob to give key1 to Alice".

Wrap your final answer in <answer></answer> tags (for example, <answer>(bob, has(alice, key1))</answer>).

```

Figure A7: Prompt used for mDKG.

```

Consider a multi-agent 2d grid Overcooked environment described as follows.

# Entities
- Agents: Alice, Bob.
- Ingredients: Tomato, Lettuce, and Onion. Ingredients start as fresh (FreshTomato, FreshLettuce, FreshOnion). They can be picked up, put down, and chopped on a cutting board (obtaining ChoppedTomato, ChoppedLettuce, ChoppedOnion).
- Plates: A kind of objects that can be picked up, put down, and be merged with any kinds of ingredients to get a dish.
- Floors: Cells where agents can move freely. Agents can be on the same floor cell at the same time.
- Counters: Cells where agents can place objects but cannot move onto. Objects can only be picked up from or placed onto counters if the agent is adjacent to the counter
- Cutting boards: Special counter cells where agents can chop ingredients. Raw ingredients need to be chopped to be merged with plate and delivered.
- Delivery stations: Special counter cells where agents can place the dish for delivery. After both agents finish to deliver their recipe, the task is done.
- Doors: Doors allow agents to move between separate rooms. Counters, cutting boards, delivery spots, and doors can block the agents' vision, so that the agents may only observe objects in the room they locate in.

# Environment representation
You will be given a representation of the environment state as input image, where entities are represented as following:
- Alice wears a blue outfit and has a female avatar with a chef hat.
- Bob wears a purple outfit and has a male avatar with a chef hat.
- Tomatoes are round and red with a green leafy top.
- Lettuces are green and leafy, circular in shape.
- Onions are purple, bulb-shaped with layered texture.
- Plates are white, circular dishes placed on counters.
- Floors are light beige tiles that agents can walk on.
- Counters are light brown tiles
- Delivery Stations are gray tiles with a large yellow star symbol.
- Doors are wooden and arched, with vertical panels and a black keyhole on the right side.

# Properties
- Each agent and object occupies a cell.
- Agents act at the same time. At each time step, each agent takes one action.
- Each agent can hold at most one object.

# Goals
- Each agent has a goal dish they need to make and deliver. There are 3 different kinds of dishes: SimpleTomato, SimpleLettuce, and SimpleOnion.
- Each dish has a recipe:
  SimpleTomato: pickup FreshTomato, chop FreshTomato to obtain ChoppedTomato, merge ChoppedTomato and Plate, deliver ChoppedTomato-Plate
  SimpleLettuce: pickup FreshLettuce, chop FreshLettuce to obtain ChoppedLettuce, merge ChoppedLettuce and Plate, deliver ChoppedLettuce-Plate
  SimpleOnion: pickup FreshOnion, chop FreshOnion to obtain ChoppedOnion, merge ChoppedOnion and Plate, deliver ChoppedOnion-Plate
- Agents can have the same goal dish. In such cases, they must deliver two dishes separately.

# Observability:
- The environment is partially observable.
- Agents know the layout of the environment, but counters, cutting boards, delivery spots, and doors can block the agents' vision. These features divide the layout into different rooms, and agents can only see objects and other agents in the same room as themselves.

# Agent Actions
- Possible actions are: move up, down, left, or right, one cell at a time - or wait without moving. In special cases agents can also pick up or put down objects or ingredients, chopping ingredients, and merging chopped ingredients with plates (see below).
- An agent can move only if there is a floor or door in that direction. Agents cannot move on a counter.
- If an agent is holding something, you will see that as text info and also in the image.
- If an agent is adjacent to a counter with an object on it, and the agent is not holding anything, they can pick the object up.
- If an agent is adjacent to an empty counter, and the agent is holding an object, they can put the object down on the counter.
- If an agent is holding a chopped ingredient and is adjacent to a counter with a plate on it, they can pick it up and merge it with the chopped ingredient they're holding.
- If an agent is holding a plate and is adjacent to a counter with a chopped ingredient on it, they can pick it up and merge it with the plate they're holding.
- If an agent is holding a raw ingredient and they are adjacent to a cutting board, they can chop the ingredient and end up holding the chopped ingredient in their hands.
- If an agent is holding a dish ready for delivery and they are adjacent to a delivery spot, they can place the dish on the spot to complete their task.

---

# Your Task as a Social Facilitator
At each timestep, your task is to decide whether to give feedback to either agent, and what that feedback should be, in order to promote prosocial behavior.
A prosocial agent:
- Helps the other agent if doing so makes it possible or more efficient to achieve their goals.
- Takes action that indicate awareness of the other's goal.

# Step-by-Step Instructions
1. Review the History
You will be given a few recent timesteps, including agent actions.

2. Infer the Agents' Goals
Based on the observed agent behaviour, infer which dish each agent is preparing.

3. Evaluate if it is possible to promote prosocial behaviour
Given the action history, the current state and the inferred goal for each agent and a set of possible feedbacks, is any of the feedbacks useful for promoting prosocial behaviour? If not, just answer "None".

4. Decide Feedback
If you think that it is useful to provide a feedback given the action history and current state, then choose one feedback that (1) enables both agents to achieve their goals and/or (2) increases overall efficiency.
Typical examples of prosocial agents include those that pass ingredients to others or take longer routes to obtain ingredients, allowing their teammates to access the closer ones.
Avoid giving feedback just because an action isn't optimal; only intervene when it clearly improves overall efficiency or agent success. If this is not the case, just choose "None".

## Recent action history
$HISTORY$

## Current state
The current state is represented in the image given as input.

## Feedback
Given the instructions above, choose which feedback to give from the possible options:
$FEEDBACK_CHOICES$

Wrap your final answer in tags (for example Bob passes the FreshOnion in (3, 4) to Alice). Also provide an explanation in tags, with 1-2 sentences on why you chose this feedback, given the inferred goals.

```

Figure A8: Prompt used for Overcooked.

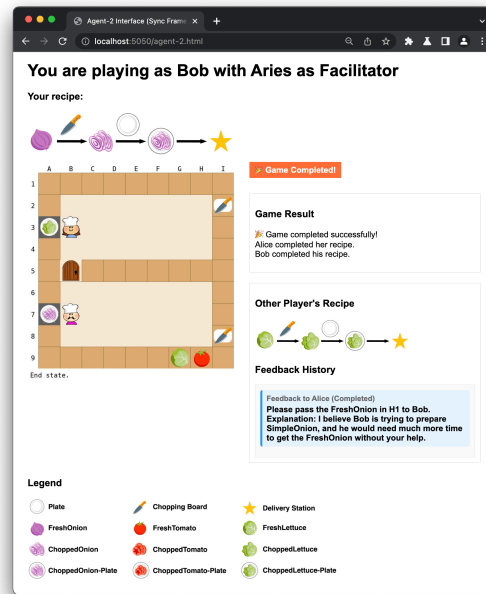
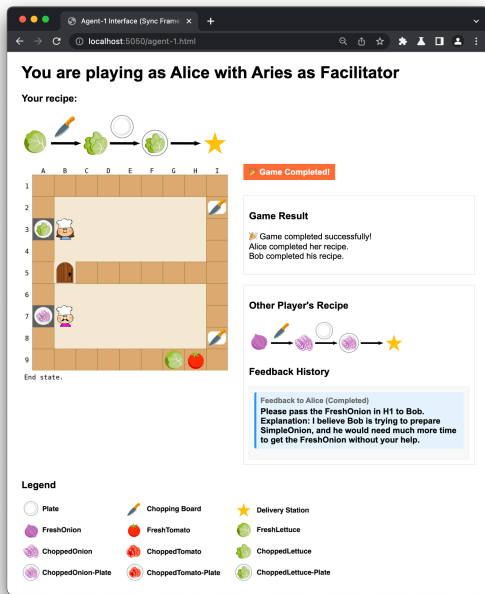
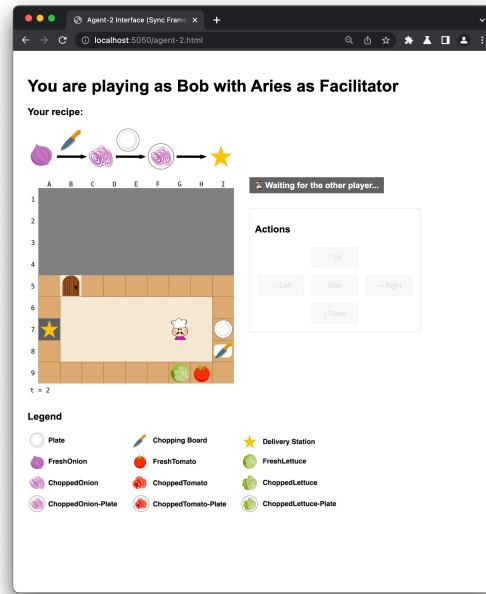
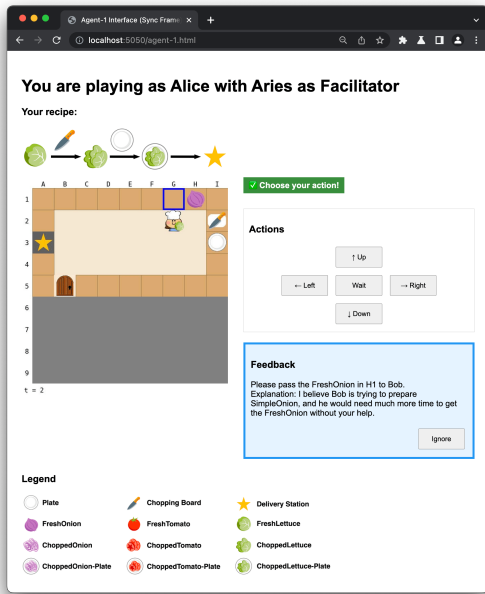


Figure A9: Interface for the human study. Each player could see their own goal recipe, their partial observation of the environment, the possible actions to take, and a legend (top). When a player received a feedback message, it would appear as a pop-up blue box, as shown on the left. Players could decide to ignore feedback messages if they considered it useless or unnecessary. At the end of each episode, both players were revealed the other player's goal recipe, information about game completion, and the full history of feedback messages given to each player, with their status ("Completed", "Ignored", "Not Executable") (bottom).

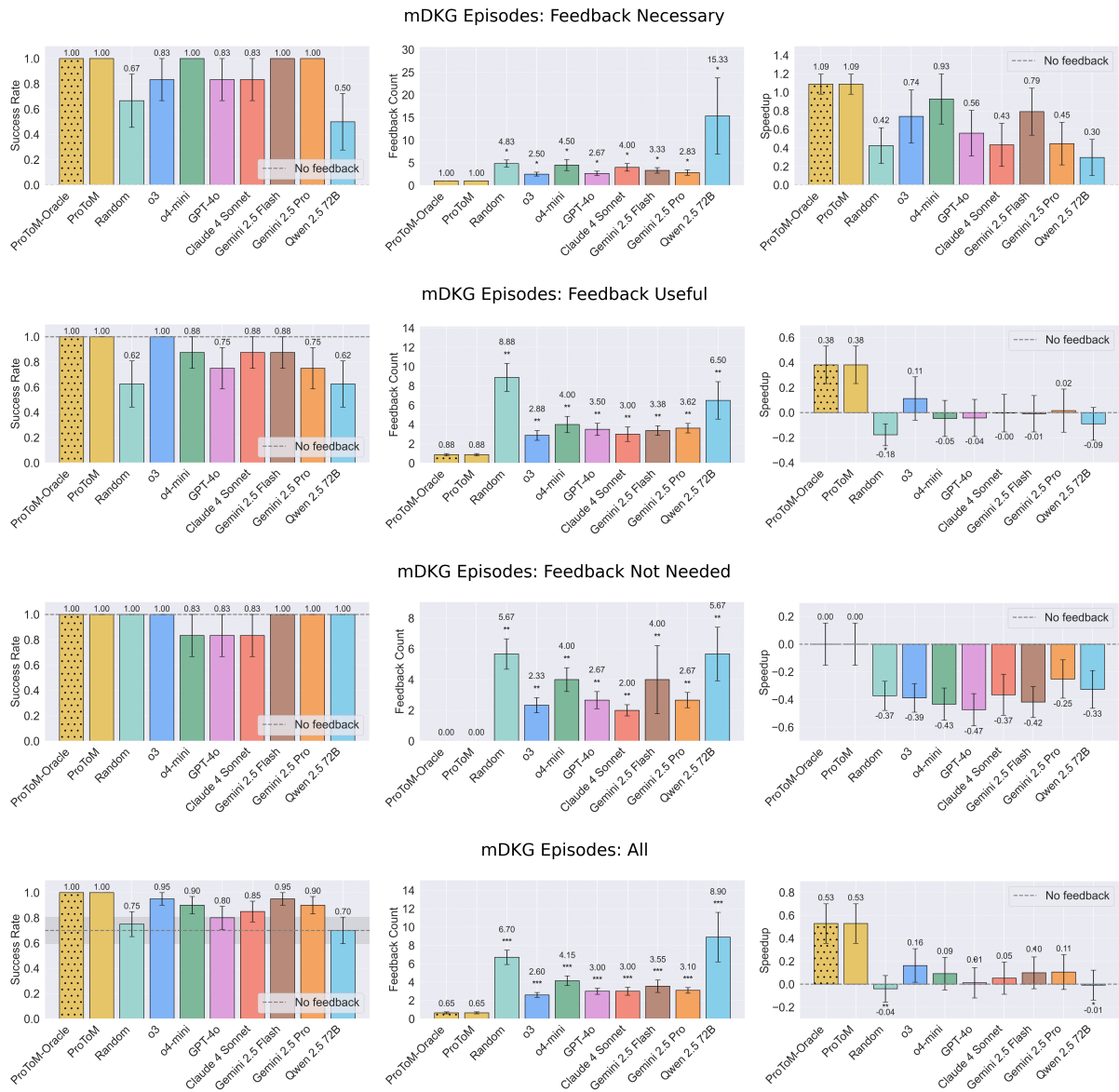


Figure A10: for different episode types in the mDKG environment.

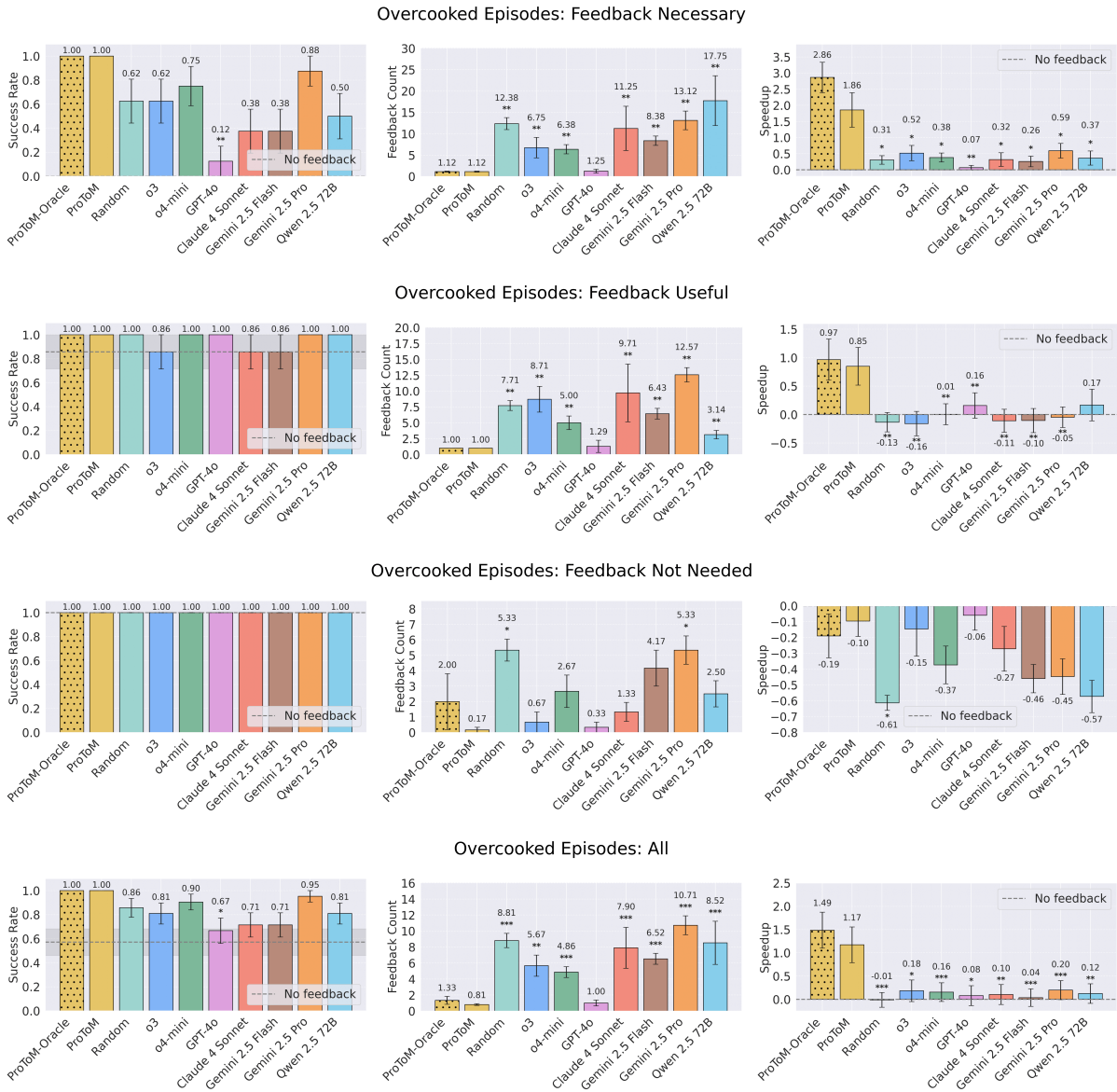


Figure A11: Results for different episode types in the Overcooked environment.