LEARNING ROBUST SOCIAL STRATEGIES WITH LARGE LANGUAGE MODELS

Anonymous authors

000

001

002 003 004

010 011

012

013

014

015

016

017

018

019

021

025

026

027

028

029

031

032033034

035

037

040

041

042

043

044

045

046

047

048

051

052

Paper under double-blind review

ABSTRACT

As agentic AI becomes more widespread, agents with distinct and possibly conflicting goals will interact in complex ways. These multi-agent interactions pose a fundamental challenge, particularly in social dilemmas, where agents' individual incentives can undermine collective welfare. While reinforcement learning (RL) has been effective for aligning large language models (LLMs) in the single-agent regime, prior small-network results suggest that standard RL in multi-agent games often converges to defecting, self-interested policies. We show the same effect in LLMs: despite cooperative priors, RL-trained LLM agents develop opportunistic behavior that can exploit even advanced closed-source models. To address this tendency of RL to converge to poor equilibria, we build on an opponentlearning awareness algorithm, Advantage Alignment, to fine-tune LLMs toward multi-agent cooperation and non-exploitability. Specifically, we derive a novel variant of Advantage Alignment under the assumption of non-observability of other players' actions on the current time step, resulting in jit-Advantage Alignment. We further introduce a group-relative baseline that simplifies advantage computation, enabling multi-agent training at LLM scale. Agents fine-tuned with our method learn the well-known tit-for-tat strategy in the classic Iterated Prisoner's Dilemma. In complex environments, our method achieves higher collective payoffs while remaining robust against exploitation by greedy agents. Finally, we contribute a suite of social dilemma benchmarks to advance the study of cooperation in agentic AI.

1 Introduction

LLMs undergo large-scale pretraining, instruction tuning, and reinforcement learning, and continue to exhibit increasingly advanced capabilities (Guo et al., 2025). Coupled with decreasing deployment costs and improved adaptability to downstream tasks, these trends enhance the commercial and practical viability of LLM agents across a wide range of applications. Recent efforts are already translating this potential into concrete systems. Anthropic's Model Context Protocol (MCP; Anthropic, 2024) enables an LLM to interact with external systems and become more capable as an autonomous decision-making agent. CICERO (FAIR et al., 2022) demonstrates strategic, human-level play with LLMs in the complex board game Diplomacy. Voyager (Wang et al., 2023) leverages Minecraft to illustrate the rising potential of LLMs as agents for open-ended exploration and skill acquisition. As LLM-agents become commonplace, new infrastructure is emerging to support agent-agent interaction, e.g. Google's Agent2Agent protocol (Agent2AgentProtocol, 2024) enabling collaboration between LLM-based agents with varying capabilities, potentially from different organizations.

Despite rapid progress, LLM behavior in multi-agent settings remains poorly understood. One common scenario involves agents with conflicting goals that discourage cooperation, even when cooperation would lead to better outcomes for all. These situations, known as *social dilemmas* (Rapoport & Chammah, 1965), frequently arise in real-world contexts where agents face a tension between individual gain and collective welfare. They appear in everyday scenarios such as navigating traffic, as well as in more complex settings such as business negotiations or international policy coordination. A recent example is the case of many LLM crawlers downloading training data from small code-hosting websites, causing them to be overwhelmed with DDoS-like traffic (SourceHut, 2025). Such interactions are analogous to the famous *tragedy of the commons*, a social dilemma concerning

the maintenance of public goods, where self-interested behavior leads to resource depletion. Such cases illustrate the types of social dilemmas that may arise in complex environments where LLMs are increasingly expected to act and interact autonomously.

Learning to resolve these scenarios is typically framed within multi-agent reinforcement learning (MARL). Unlike single-agent RL, where an agent improves an objective in a static environment, in MARL each agent must adapt to the strategies of other agents, who can also be learning over time. This leads to non-stationarity, since the policy of each learning agent affects the collective outcome. Social dilemmas can be treated as a specific subclass of MARL problems that are mixed-motive; neither fully cooperative nor fully competitive. Initial attempts using MARL to play social dilemmas were unsuccessful. Training agents based on small neural networks with standard RL resulted in sub-optimal greedy strategies (Sandholm & Crites, 1996). To account for other learning agents, Foerster et al. (2018) introduced Opponent Shaping (OS), an RL paradigm that explicitly considers agent interactions in hopes of steering their dynamics towards mutually beneficial outcomes. LOLA, the first OS algorithm, is capable of finding the pareto-optimal strategy of *tit-for-tat* in simple social dilemmas like the Iterated Prisoner's Dilemma.

However, modern agents are not small, randomly initialized networks. Prior work largely focused on teaching such tabula-rasa agents reciprocity—punishing greed and rewarding cooperation—where the central obstacle was that uninformed policies gravitated toward short-sighted, self-interested strategies. By contrast, LLMs arrive with rich priors and human-like social norms induced by pretraining and post-training (instruction tuning/RLHF) (Ross et al., 2024), potentially altering the learning dynamics and failure modes in multi-agent settings. This raises a key question: when fine-tuned with standard RL, do LLMs reproduce the same failure modes as small networks, or do their human-biased priors mitigate them? Given that LLM agents already interact in the wild, we investigate whether naïve RL fine-tuning systematically erodes collective welfare.

To evaluate this, we introduce a novel testbed for social dilemmas in the LLM setting. The testbed integrates both small-scale social dilemma frameworks (Duque et al., 2025) which we extend into the textual domain and our newly proposed Split Games, designed to measure cooperation as well as exploitability and require agents to gain and maintain trust. Using this environment, we conduct extensive experiments across a range of modern LLMs and find that standard RL consistently produces greedy behavior across all settings. Probing further, we show that even state-of-the-art closed-source models are exploitable when facing agents trained with standard RL. These results underscore that current LLMs are not yet prepared to robustly operate in real-world multi-agent settings and highlight a novel risk of current agentic AI.

Next, building on recent advances (Duque et al., 2025), we introduce jit-Advantage Alignment, an opponent-shaping algorithm that trains agents to align their incentives with their opponent and learn both cooperation and non-exploitability. When trained with jit-Advantage Alignment, agents learn the non-exploitable and effective *tit-for-tat* strategy in the classic Iterated Prisoner's Dilemma. In our novel testbed, jit-Advantage Alignment agents learn to cooperate with cooperative players as well as themselves, while remaining robust against greedy players.

In summary, our key contributions are:

- Re-deriving an improved formulation of the Advantage Alignment algorithm (Duque et al., 2025), creating jit-Advantage Alignment
- Creating a novel testbed of social dilemmas for LLM agents, both local and API models. We include both standard games and extended games with extra communication channels.
- Demonstrating that standard RL leads to greedy, suboptimal agents across a range of open-source LLM models
- Discovering that state-of-the-art closed-source LLMs are susceptible to being exploited by greedy agents
- Applying jit-Advantage Alignment to achieve cooperative and non-exploitatable agents in two games of our novel testbed.

2 BACKGROUND

2.1 MARKOV GAMES

A n-agents Markov Game can be defined as a tuple $(\Pi, \mathcal{S}, \mathcal{A}, \mathcal{R}, \mathcal{P}, \gamma)$. \mathcal{S} is a set of possible states. \mathcal{A} is a set of functions $\mathcal{A}_1, \ldots, \mathcal{A}_n$ where $\mathcal{A}_i(\mathcal{S})$ gives the set of possible actions of agent i at state S. \mathcal{R} is the set of reward functions $\{R_1, \ldots, R_n\}$ where $R_i: \mathcal{S} \to \mathbb{R}$ is the reward function of agent i. \mathcal{P} is the transition function, mapping a probability distribution over each transition $\mathcal{P}(S \to S')$. Π is the set of policies $\{\pi_1, \ldots, \pi_n\}$, each π_i mapping any state S to a probability distribution over $\mathcal{A}_i(S)$. γ is the discount factor on the returns. The expected discounted return of player i is $V_i(\pi) = \mathbb{E}_{\tau \sim \Pr_{\mu}^{\Pi}} \left[\sum_{t=0}^{\infty} \gamma^t R_i(s_t, \mathbf{a}_t) \right]$, where \Pr_{μ}^{Π} is distribution of trajectories induced by the initial state distribution μ and the set of policies Π .

For any Markov Decision Process 1 the REINFORCE algorithm uses unbiased estimates of the gradient of J with respect to the policy parameters π in order to perform gradient ascent on $J(\pi)$ (Zhang et al., 2020). We refer to applying REINFORCE independently to each agent optimizing their objectives $V_i^*(\pi_i)$, while considering the other agents' policies as fixed at each learning step as standard MARL. In contrast to Standarad MARL, where each agent optimizes its own return, Cooperative MARL optimizes for the sum of expected returns across all agents. Specifically, We define Cooperative MARL as applying REINFORCE on each agent with objective modified as $V_i^*(\pi_i) := \sum_{j=1}^n V_j(\pi_j)$. This formulation encourages agents to learn policies that maximize overall welfare rather focusing on individual benefits.

2.2 ROBUST SOCIAL STRATEGIES

In a zero-sum game, the agents' payoffs always add up to zero; every gain for one side is matched by an equal loss for the other. Consequently, in a two-player zero-sum setting, cooperation offers no benefit. In this work, we focus on general-sum games, where total payoffs is not fixed, and agents may improve their outcomes without necessarily diminishing those of others, thereby creating the possibility of beneficial cooperation. However, in social dilemmas, agents face a tension between their short-term individual benefit and possible long-term collective welfare. Each agent has a short-term incentive to act selfishly (i.e., not cooperate), but if all agents do so, the resulting outcome leads to reduced overall welfare i.e. a lower total sum of discounted returns for all agents. A socially robust strategy is one that encourages other rational agents to behave in the collective interest. Such a strategy does not maximize collective welfare at its own expense, but rather optimizes its own long-term self-interest by optimizing that of the group. Formally, we consider π_1^* to be a more socially robust policy if $V_1(\pi_1^*, \arg \max_{\pi_2} V_2(\pi_1^*, \pi_2)) > V_1(\pi_1, \arg \max_{\pi_2} V_2(\pi_1, \pi_2))$. That is, a socially robust policy takes into consideration the best response of the opponent.

2.3 OPPONENT SHAPING

Prior work shows that small neural networks trained with standard RL tend to converge to the *Always Defect* strategy in IPD (Sandholm & Crites, 1996). More recently, Foerster et al. (2018) demonstrated that this undesirable outcome also arises with policy gradient methods. This outcome is expected, as policy gradients are obtained via Monte-Carlo trajectory samples with fixed policies. Retaliatory behaviour is rare to randomly sample, while greedy behaviour is easier to discover. This results in agents that are drawn into a greedy equilibrium and achieve suboptimal payoffs. LOLA (Foerster et al., 2018) removed the assumption of a static environment in Markov Games and included a model of a learning agent in its update. By explicitly modelling how opponent learning is affected by an agent's action, LOLA was able to learn the *tit-for-tat* strategy in IPD. Unfortunately, LOLA's computational complexity is quadratic in the number of parameters of the agent, making it prohibitively computationally expensive for modern Large Language Models (LLMs).

Advantage Alignment (Duque et al., 2025) is an opponent-shaping algorithm that addresses scalability by focusing on the Q-values of both agent and opponent. Assuming that agents act proportionally to the exponent of their Q-value, Advantage Alignment aims to align an opponent's Q-value with your own. This leads to a simple modification to the advantages used in the policy gradient term of

¹which is a special case of Markov Game with n=1

a REINFORCE estimator. Advantage Alignment is capable of solving social dilemmas in scenarios with high dimensional state representations (e.g. pixel spaces), partial observability, continuous action spaces, and delayed action outcomes. Due to its performance in complex scenarios, we chose Advantage Alignment as a prime candidate to train LLMs to find robust social strategies.

3 JIT-ADVANTAGE ALIGNMENT

3.1 RE-DERIVING THE ADVANTAGE ALIGNMENT FORMULATION

The original Advantage Alignment algorithm (Duque et al., 2025) made two assumptions about the opponents with which the agentinteracts: (1) each agent learns to maximize their value function; and (2) each agent acts proportionally to the *softmax* of the Q-values. Under these conditions, they derive the following opponent shaping formula:

$$\beta \cdot \mathbb{E}_{\tau \sim \Pr_{\mu}^{\pi^1, \pi^2}} \left[\sum_{t=0}^{\infty} \gamma^{t+1} \left(\sum_{k < t} \gamma^{t-k} A^1(s_k, a_k, b_k) \right) A^2(s_t, a_t, b_t) \nabla_{\theta^1} \log \pi^1(a_t | s_t) \right]. \tag{1}$$

The authors implicitly make one additional assumption: (3) agents are able to observe the actions of other players at the current time t. We relax this third assumption to derive a modified Advantage Alignment term (see Appendix A). Due to the difficulty of estimating the term under the relaxed assumption, we approximate it as follows:

$$\beta \cdot \mathbb{E}_{\tau \sim \Pr_{\mu}^{\pi^1, \pi^2}} \left[\sum_{t=0}^{\infty} \gamma^t \left(\sum_{k \le t} \gamma^{t-k} A^1(s_k, a_k, b_k) \right) A^2(s_t, a_t, b_t) \nabla_{\theta^1} \log \pi^1(a_t | s_t) \right], \tag{2}$$

Crucially, the main difference lies in incorporating the advantage of the agent at the current timestep. This change is necessary for the method to work in our LLM benchmarks. Given that our modification considers the *just-in-time* advantage, we call it *jit*-Advantage-Alignment or *jit*-AA. Implementing this formulation requires scalable advantage estimation, which is a non-trivial task in the context of LLM training. For this reason, we introduce group-relative baseline that simplifies advantage computation.

3.2 A GROUP-RELATIVE BASELINE FOR ADVANTAGE COMPUTATION

Estimating advantages with value networks has proven challenging in the context of LLM training, often leading to unstable or ineffective results (Kazemnejad et al., 2024). Recent work, such as RLOO (Ahmadian et al., 2024) and GRPO Shao et al. (2024) has shown that baseline-based approaches provide more stable and efficient advantage estimates. We build on this idea and extend it to multiagent training of LLMs over multiple rounds. In our setting, each action a_t^i corresponds to a response generated by LLM agent i given the context. We define the state at time t as $s_t = \{x_0, \mathbf{a}_0, x_1, \mathbf{a}_1, \dots, x_t\}$, where $x_0 \sim \mu$ is the *initial game context* and $x_{>0}$ represents the *intermediate game context*. The initial game context x_0 provides the prompt that explains the game setup, while each intermediate game context summarizes the outcome of previous round along with new information for the next one. To compute the advantage of an action for each player, we fix the initial game context across k parallel games. Let $A^i(s_t, a_t)$ denote the advantage for agent i, we estimate it using a leave-one-out group baseline computed over the k games at each time step t:

$$\frac{1}{k} \sum_{i=1}^{k} \left[G(a_t^{(i)}, s_t) - \frac{1}{k-1} \sum_{j \neq i} G(a_t^{(j)}, s_t) \right]$$

where $G(a_t^{(i)}, s_t)$ is the discounted return for action $a_t^{(i)}$ taken in state s_t . This group-relative baseline avoids the need for a learned value function, simplifies advantage computation, and enables multi-turn RL training with LLMs in our experiments.

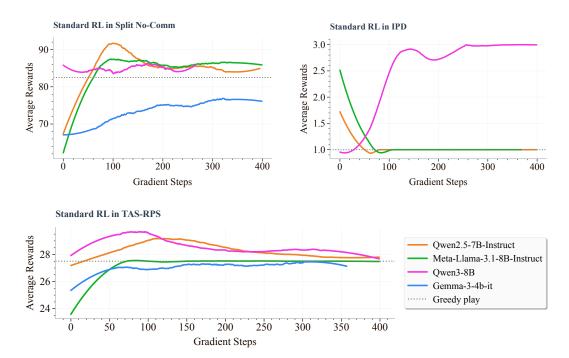


Figure 1: Training curves of standard RL on a set of open source LLMs on multiple environments. In all but one case, the models converge towards the payoffs of greedy strategies, or worst. The only exception is Qwen3-8B in IPD, which has a reciprocal component that makes it converge to cooperation.

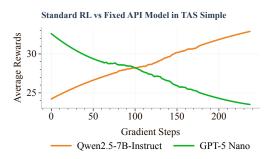
4 Social Dilemma Testbed

Our goal in this work is to study the behavior of LLM Agents trained with RL in general sum game environments with a focus on social dilemmas. To support this, we develop a novel testbed to evaluate the effects of RL training on cooperation and resistance to exploitation.

Iterated Prisoner's Dilemma The Iterated Prisoner's Dilemma is a two-player game where agents repeatedly and simultaneously choose to either *Cooperate* (C) or *Defect* (D). The per-round pay-off matrix used in our experiments is provided in Appendix C. We include IPD in our testbed because it is one of the most widely studied social dilemmas. However, since it is also likely presented in the training data of LLMs, we obsfucate the nature of the game by removing any mention of "Prisoner's Dilemma" and replace the action labels from *Cooperate* and *Defect* to A and B, respectively. This allows us to test how well LLMs generalize beyond memorization and to examine how RL interacts with any prior knowlege the model may have about this social dilemma.

Split Games While IPD captures fundamental dilemma, it lacks the richness of real-world strategic interactions. To address this, we propose Split Games, a novel class of social dilemmas that combine the natural language interaction framework of Deal or No Deal (Lewis et al., 2017) with the proposal mechanicsm from the Negotiation Game (Cao et al., 2018; Duque et al., 2025), providing a better learning signal for training agents in this dilemma. In each game, two agents are assigned private values $v_i \sim \mathcal{U}[1,10]$ over a set of items, and must negotiate through message exchanges to decide how to split the items. The setup allows for variety of behaviors, including bluffing, exaggeration, and cooperative negotiation. Formally, let $p_{k,a}$ be the proposal for the k'th item category from agent a and q_k be the quantity available. The allocation received by agent a is $q_{k,a} = q_k p_{k,a} / \max(q_k, p_{k,a} + p_{k,b})$ and similarly for agent b. The resulting payoffs are $v_a \times q_{k,a}$ and $v_b \times q_{k,b}$, respectively. This particular mechanism-design choice removes the need for explicit agreement and provides a training signal in each round.

We design two main variants of the game:



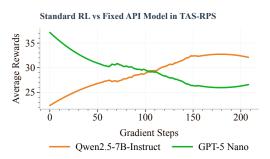


Figure 2: In TAS-Simple, a Qwen2.5-7B-Instruct model trained under standard RL learns to exploit GPT-5 Nano, a new closed-source thinking model. On TAS-RPS, GPT-5 Nano appears more robust.

Split No-Comm Here, agents observe each other's values. Proposals and payoffs are revealed at the end of each round. This variant supports reciprocity but avoids the challenges of communication.

Trust and Split In this setting, agents have private values and must communicate to negotiate. At the end of each round, values and proposals are revelead to both the players, allowing reciprocity. However, communication to negotiate multiple items introduces long contexts that can be challenging for opensource LLMs. Liao et al. (2024) find that LLMs upto the scale of 70B struggle to follow instruction in multiple item setting. To address this limitation, we design simplified sub-variants with a single item type, coins:

- TAS-Simple: Each agent privately receives a value in [1, 10] for the coin. After exchanging a single message, agents submit proposals. However, since maximum value is fixed, agents can sometimes infer if their value is higher or lower without communication.
- TAS-RPS: To remove such cases, each agent in this setting is instead given a hand of rock, paper, or scissors. The player with winning hand values coins at 10, while the losing hand player values coins at 1. Successful negotiation therefore requires communication to uncover coin values as seen in Figures 6 and 3.

An in-depth description of each variant with sample conversations can be found at Appendix D. Across all Split Games, starting agent alternates every round, and item values are assigned so that, in expectation, both agents receive the same total value. The strategy that maximizes payoffs for both the agents in each of the Split Games is to truthfully communicate private values or hands and allocate all items to the agent who values them more highly in each round, while remaining non-exploitable.

5 EXPERIMENTS

5.1 STANDARD RL LEADS TO GREEDY BEHAVIOR WITH LLMS

With the testbed in place, we now study how RL shapes LLM behavior in these setting. Each agent's policy π^i is parameterized by θ_i and implemented via LoRA finetuning (Hu et al., 2022). Throughout our experiments, we refer to the first player as Alice and the second as Bob. We use self-play unless stated otherwise, i.e, the same set of parameters for both agents, conditioned on different game contexts based on their roles. This ensures that memory usage doesn't scale with the number of agents and the model size we used is sufficient enough to handle the complexity of the different roles. Moreover, our games are symmetric, so the agent's behavior should ideally remain consistent across roles. Details of game-specific hyperparameters are provided in Appendix B. For games played over infinite rounds with a discount factor δ , we found no empirical difference between training with fixed-length versus stochastic-length trajectories. For computational efficiency, we therefore use fixed-length trajectories throughout. For models that cannot reliably follow instructions in our testbed, we use the constrained decoding setup of vLLM to enforce a regex pattern during generation.

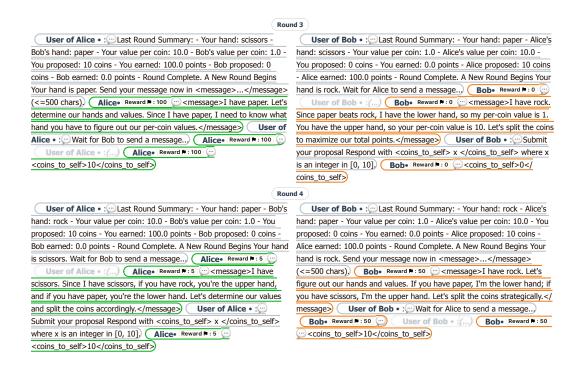


Figure 3: Two rounds of untrained Qwen3-8B playing TAS-RPS.

We evaluate multiple model families in our testbed, as shown in Figure 1. We find that standard RL consistently leads to greedy behavior across model families and testbeds. Interestingly, if we treat the game as fully cooperative, as seen in Figure 11 in the Appendix, some of these same models are capable of learning cooperative strategies. This reinforces that while models can learn cooperation, standard RL training drives them toward greedy behavior. In Split No Comm, we find that agents learned to give the highest bid for every item even when they value it less. From qualitative analysis of conversation in TAS-RPS game, we find that agents are communicating honestly, but still bid high for all coins. This greedy behavior is strictly worse for all players, compared to cooperating by bidding only when they value the item highly.

Next, we scale our experiments by training a Qwen-2.5-7B-Instruct agent against a frozen GPT-5-nano opponent using standard RL. Here, one agent (Alice) is the learning agent finetuned via LoRA, while the other (Bob) is fixed GPT-5-nano. We focus on two games in the testbed: TAS-Simple and TAS-RPS, both of which involve splitting 10 coins. In TAS-Simple, players privately receive numerical per coin value and can communicate before making the proposal. In TAS-RPS, the coin's value of each player depends on the players' hands (rock, paper, or scissors), which makes trust and communication even more critical. We chose these games because they are tractable with LLMs at the 8B scale. Figure 2 shows that the learning agent trained with standard RL continues to exploit the fixed GPT-5-nano in TAS-Simple, where the average reward of the learning agent steadily increases, while that of GPT-5-nano consistently decreases. In TAS-RPS, the learning agent also achieves higher average rewards compared to GPT-5-nano, although convergence has not yet been achieved so it is difficult to draw strong conclusions in this setting. These results underscore that RL-trained agents become increasingly greedy in social dilemma contexts, and that even advanced closed-source models remain vulnerable to such exploitation.

5.2 JIT-ADVANTAGE ALIGNMENT LEARNS ROBUST STRATEGIES

In this section, we evaluate jit-Advantage Alignment on two games from our testbed: the Iterated Prisoner's Dilemma (IPD) and Split No-Comm game. We selected these games because LLM agents at the 8B scale can reliably follow instructions.and be effectively trained. In the IPD, agents learned with jit-Advantage Alignment cooperate with themselves and with fully cooperative agents, while remaining robust against defector (Figure 4). On empirical evaluation, we find that our agent learned

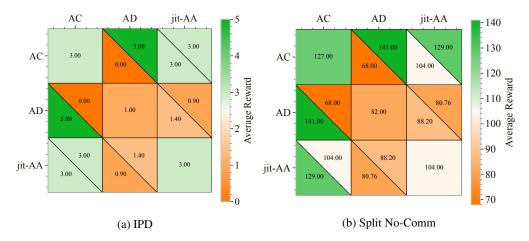


Figure 4: Average rewards obtained under different permutations of agents. "AC" stands for "always cooperate" and "AD" for "always defect". In IPD, our model has a comparable performance to a *tit-for-tat* agent. In Split No-Comm, our model performs better in self-play than the "AD" agent and is not exploitable like the "AC" agent. The "AC" and "AD" agents here are hard coded.

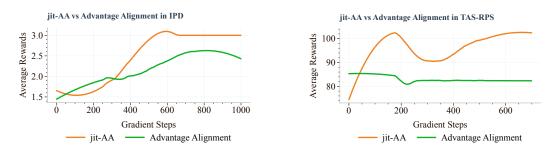


Figure 5: Comparison of jit-AA and Advantage Alignment.

tit-for-tat strategy. The slight drop in performance against defectors is due to losing the first round, where the agent initially cooperates and receives a lower payoff. In the Split No-Comm game, jit-Advantage Alignment agents achieve about 80% of the efficiency of full cooperation while maintaining robustness. When paired with defectors, their performance decreases only slightly, indicating they are not easily exploitable as seen in Figure 4. Moreover, jit-Advantage Alignment achieves higher average reward than Advantage Alignment in both settings as seen in Figure 5. Finally, when extending jit-Advantage Alignment to communication variants of our testbed, we observed training instabilities as the models struggled to reliably follow game instructions.

6 RELATED WORK

Negotiation, especially in games like DoND (Lewis et al., 2017), inherently involves coordination and adaptation to another agent's behavior, making it a natural testbed for broader questions in multiagent cooperation. More recently, Liao et al. (2024) used the game as a benchmark to test behaviour cloning training to train closed source Large Language Models. Coordination and negotiation pose significant challenges in multi-agent reinforcement learning (MARL). Dafoe et al. (2020) highlight key open problems in MARL such as communication and cooperation in mixed-motive settings. Unlike competitive settings, cooperative settings demand that agents develop shared norms and robust coordination protocols. Agashe et al. (2025) propose the LLM-Coordination Benchmark to evaluate LLMs in multi-agent pure coordination games through two tasks: Agentic Coordination and CoordQA. Their results reveal key limitations in LLMs' ability to reason about partners' beliefs and intentions, an essential component for effective coordination. (Li et al., 2023) evaluate LLM-based agents in a multi-agent cooperative text game involving Theory of Mind inference tasks and observe evidence of emergent collaborative behavior. In contrast, our work leverages RL fine-tuning to di-

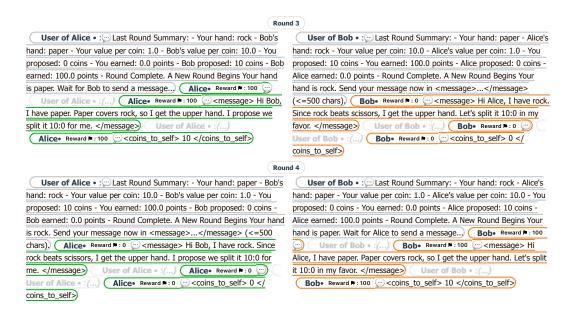


Figure 6: Example of perfect cooperative play TAS-RPS for two rounds.

rectly optimize agents on the outcomes of their own proposals, demonstrating that such fine-tuning can strip away the cooperative behavior and instead drive more outcome-oriented behavior.

Opponent shaping was introduced in Foerster et al. (2018) as a paradigm that assumes opponents are naive REINFORCE-based learners and attempts to shape their learning trajectories. Other opponent shaping methods treat the learning process as a meta-game in the space of policy parameters, where inter-episode returns constitute rewards and policy updates constitute actions (Lu et al., 2022). Alternatively, opponent shaping can be done by differentiating through a best response opponent (Aghajohari et al., 2024a) or by influencing the joint probability distribution over trajectories to control the Q-values (Aghajohari et al., 2024b). Advantage Alignment (Duque et al., 2025) reduces opponent shaping to a functional modification of the advantage that is used in standard policy gradient, greatly improving its scalability. We propose jit-Advantage Alignment, which builds upon this approach by relaxing the assumption that the opponent's current action is known, leading to better performance in our testbed.

7 CONCLUSION

In this work, we investigated the shortcomings of training large language models (LLMs) with standard rei end, we introduced a testbed of social dilemmas, Split Games and trained across model families. Our results demonstrate that standard RL training leads to greedy behavior even with cooperative prior of LLMs. Furthermore, we found closed source advanced LLMs to be exploitable when interacting with trained agents. These findings highlight critical limitations in the current approach to training LLM agents, namely standard RL, when evaluated in realistic test beds such as social dilemmas. To address these challenges, we proposed jit-Advantage Alignment which learns cooperative behavior while remaining robust against exploitation. We show that jit-Advantage Alignment improves upon recent work by relaxing the assumption of assuming the opponent's action at the current time step and performs better. In particular, jit-Advantage Alignment learns tit-for-tat strategy in IPD bechmark with LLMs, and achieves higher payoffs while being less exploitable against greedy agents in Splits No Comm. game. In future work, we are interested in exploring jit-Advantage Alignment applied to communication variants of our testbed, with the aim of learning cooperative and non-exploitable strategies with LLMs.

8 ETHICS STATEMENT

We are not aware of any either negative or positive societal implications of our work. Our work is primarily focused on diagnosing issues related with RL with LLMs in academic benchmarks. Our work does not involve any large scale training, restricting itself to training small scale models.

9 REPRODUCIBILITY STATEMENT

We include detailed prompts, game specifications, and payoff rules in the appendix C and D We also include training/eval hyperparameters used in our experiments in the appendix B. We will release code, configs, prompts, and evaluation logs to replicate figures and tables and to rerun all baselines.

REFERENCES

- Saaket Agashe, Yue Fan, Anthony Reyna, and Xin Eric Wang. Llm-coordination: Evaluating and analyzing multi-agent coordination abilities in large language models, 2025. URL https://arxiv.org/abs/2310.03903.
- Agent2AgentProtocol. Agent2agentprotocol. https://github.com/google/A2A, 2024.
- Milad Aghajohari, Tim Cooijmans, Juan Agustin Duque, Shunichi Akatsuka, and Aaron Courville. Best response shaping, 2024a. URL https://arxiv.org/abs/2404.06519.
- Milad Aghajohari, Juan Agustin Duque, Tim Cooijmans, and Aaron Courville. Loqa: Learning with opponent q-learning awareness, 2024b. URL https://arxiv.org/abs/2405.01035.
- Arash Ahmadian, Chris Cremer, Matthias Gallé, Marzieh Fadaee, Julia Kreutzer, Olivier Pietquin, Ahmet Üstün, and Sara Hooker. Back to basics: Revisiting reinforce style optimization for learning from human feedback in llms. *arXiv preprint arXiv:2402.14740*, 2024.
- Anthropic. Model context protocol (mcp). https://docs.anthropic.com/en/docs/agents-and-tools/mcp, 2024.
- Kris Cao, Angeliki Lazaridou, Marc Lanctot, Joel Z. Leibo, Karl Tuyls, and Stephen Clark. Emergent communication through negotiation. In *ICLR*, 2018.
- Allan Dafoe, Edward Hughes, Yoram Bachrach, Tantum Collins, Kevin R. McKee, Joel Z. Leibo, Kate Larson, and Thore Graepel. Open problems in cooperative ai, 2020. URL https://arxiv.org/abs/2012.08630.
- Juan Agustin Duque, Milad Aghajohari, Tim Cooijmans, Razvan Ciuca, Tianyu Zhang, Gauthier Gidel, and Aaron Courville. Advantage alignment algorithms, 2025. URL https://arxiv.org/abs/2406.14662.
- FAIR, Anton Bakhtin, Noam Brown, Emily Dinan, Gabriele Farina, Colin Flaherty, Daniel Fried, Andrew Goff, Jonathan Gray, Hengyuan Hu, et al. Human-level play in the game of diplomacy by combining language models with strategic reasoning. *Science*, 378(6624):1067–1074, 2022.
- Jakob N. Foerster, Richard Y. Chen, Maruan Al-Shedivat, Shimon Whiteson, Pieter Abbeel, and Igor Mordatch. Learning with opponent-learning awareness, 2018. URL https://arxiv.org/abs/1709.04326.
- Daya Guo, Dejian Yang, Haowei Zhang, Junxiao Song, Ruoyu Zhang, Runxin Xu, Qihao Zhu, Shirong Ma, Peiyi Wang, Xiao Bi, et al. Deepseek-r1: Incentivizing reasoning capability in llms via reinforcement learning. *arXiv* preprint arXiv:2501.12948, 2025.
- Edward J Hu, Yelong Shen, Phillip Wallis, Zeyuan Allen-Zhu, Yuanzhi Li, Shean Wang, Lu Wang, Weizhu Chen, et al. Lora: Low-rank adaptation of large language models. *ICLR*, 1(2):3, 2022.
- Amirhossein Kazemnejad, Milad Aghajohari, Eva Portelance, Alessandro Sordoni, Siva Reddy, Aaron Courville, and Nicolas Le Roux. Vineppo: Unlocking rl potential for llm reasoning through refined credit assignment. 2024.
- Mike Lewis, Denis Yarats, Yann N. Dauphin, Devi Parikh, and Dhruv Batra. Deal or no deal? end-to-end learning for negotiation dialogues, 2017. URL https://arxiv.org/abs/1706.05125.
- Huao Li, Yu Chong, Simon Stepputtis, Joseph Campbell, Dana Hughes, Charles Lewis, and Katia Sycara. Theory of mind for multi-agent collaboration via large language models. In *Proceedings of the 2023 Conference on Empirical Methods in Natural Language Processing*. Association for Computational Linguistics, 2023. doi: 10.18653/v1/2023.emnlp-main.13. URL http://dx.doi.org/10.18653/v1/2023.emnlp-main.13.
- Austen Liao, Nicholas Tomlin, and Dan Klein. Efficacy of language model self-play in non-zero-sum games, 2024. URL https://arxiv.org/abs/2406.18872.

- Chris Lu, Timon Willi, Christian Schroeder de Witt, and Jakob Foerster. Model-free opponent shaping, 2022. URL https://arxiv.org/abs/2205.01447.
- Anatol Rapoport and Albert Chammah. *Prisoner's Dilemma: A Study in Conflict and Cooperation*. University of Michigan Press, 1965.
- Jillian Ross, Yoon Kim, and Andrew W. Lo. Llm economicus? mapping the behavioral biases of llms via utility theory. In *COLM*, 2024.
- Tuomas Sandholm and Robert Crites. Multiagent reinforcement learning in the iterated prisoner's dilemma. *Bio Systems*, 37(1-2):147–166, 1996.
- Zhihong Shao, Peiyi Wang, Qihao Zhu, Runxin Xu, Junxiao Song, Xiao Bi, Haowei Zhang, Mingchuan Zhang, YK Li, Y Wu, et al. Deepseekmath: Pushing the limits of mathematical reasoning in open language models. *arXiv preprint arXiv:2402.03300*, 2024.
- SourceHut. Sourcehut. https://status.sr.ht/issues/2025-03-17-git.sr. ht-llms/, 2025.
- Guanzhi Wang, Yuqi Xie, Yunfan Jiang, Ajay Mandlekar, Chaowei Xiao, Yuke Zhu, Linxi Fan, and Anima Anandkumar. Voyager: An open-ended embodied agent with large language models. arXiv preprint arXiv:2305.16291, 2023.
- Junzi Zhang, Jongho Kim, Brendan O'Donoghue, and Stephen Boyd. Sample efficient reinforcement learning with reinforce, 2020. URL https://arxiv.org/abs/2010.11364.

A JIT-ADALIGN DERIVATION

We relax the observability of the actions of other players at the current time-step by considering instead the expectation over their policies. Therefore, the opponent modeling assumption is now:

$$\hat{\pi}^i(b|s) = \frac{\exp \beta \mathbb{E}_{a \sim \pi^{-i}(\cdot|s)}[Q^i(s, a, b)]}{\sum_b \exp \beta \mathbb{E}_{a \sim \pi^{-i}(\cdot|s)}[Q^i(s, a, b)]}$$

Recall the opponent shaping policy gradient expression:

$$\nabla_{\theta_1} V^1(\mu) = \mathbb{E}_{\tau \sim \Pr_{\mu}^{\pi^1, \pi^2}} \left[\sum_{t=0}^{\infty} \gamma^t A^1(s_t, a_t, b_t) \left(\underbrace{\nabla_{\theta_1} \log \pi^1(a_t | s_t)}_{(A)} + \underbrace{\nabla_{\theta_1} \log \hat{\pi}^2(b_t | s_t)}_{(B)} \right) \right]$$

We expand the term (B) as:

$$\mathbb{E}_{\tau \sim \Pr_{\mu}^{\pi^1, \pi^2}} \left[\sum_{t=0}^{\infty} \gamma^t A_t^1 \nabla_{\theta_1} \log \pi^2(b_t | s_t) \right]$$
(3)

$$= \mathbb{E}_{\tau \sim \Pr_{\mu}^{\pi^{1}, \pi^{2}}} \left[\sum_{t=0}^{\infty} \gamma^{t} A_{t}^{1} \nabla_{\theta_{1}} \log \frac{\exp \beta \mathbb{E}_{a'_{t} \sim \pi^{1}(\cdot|s_{t})}[Q^{2}(s_{t}, a'_{t}, b_{t})]}{\sum_{b} \exp \beta \mathbb{E}_{a'_{t} \sim \pi^{1}(\cdot|s_{t})}[Q^{2}(s_{t}, a'_{t}, b_{t})]} \right]$$
(4)

$$= \mathbb{E}_{\tau \sim \Pr_{\mu}^{\pi^{1}, \pi^{2}}} \left[\sum_{t=0}^{\infty} \gamma^{t} A_{t}^{1} \left(\nabla_{\theta_{1}} \log \exp \beta \mathbb{E}_{a'_{t} \sim \pi^{1}(|s_{t})} [Q^{2}(s_{t}, a'_{t}, b_{t})] - \nabla_{\theta_{1}} \log \sum_{t \in \mathcal{S}_{t}} (\dots) \right) \right]$$
(5)

$$= \beta \cdot \mathbb{E}_{\tau \sim \Pr_{\mu}^{\pi^{1}, \pi^{2}}} \left[\sum_{t=0}^{\infty} \gamma^{t} A_{t}^{1} \nabla_{\theta_{1}} \mathbb{E}_{a'_{t} \sim \pi^{1}(|s_{t})} [Q^{2}(s_{t}, a'_{t}, b_{t})] \right]. \tag{6}$$

where in line 6 line we used the fact that any term of the form $A^1(s_t, a_t, b_t) f(s_t)$ in the expectation will vanish². For convenience of notation, we define:

$$r_t^i := r^i(s_t, a_t', b_t), \ A_t^i := A^i(s_t, a_t, b_t), \ Q_t^i := Q^i(s_t, a_t', b_t)$$

These are the reward and advantage of agent i at time step t after taking action a_t and opponent taking action b_t .

$$\beta \cdot \mathbb{E}_{\tau \sim \Pr_{\mu}^{\pi^{1}, \pi^{2}}} \left[\sum_{t=0}^{\infty} \gamma^{t} A^{1}(s_{t}, a_{t}, b_{t}) \nabla_{\theta^{1}} \mathbb{E}_{a'_{t} \sim \pi^{1}(|s_{t})} [Q^{2}(s_{t}, a'_{t}, b_{t})] \right]$$
(7)

$$= \beta \cdot \mathbb{E}_{\tau \sim \Pr_{\mu}^{\pi^{1}, \pi^{2}}} \left[\sum_{t=0}^{\infty} \gamma^{t} A^{1}(s_{t}, a_{t}, b_{t}) \nabla_{\theta^{1}} \sum_{a'_{t}} \pi_{1}(a'_{t}|s_{t}) Q^{2}(s_{t}, a'_{t}, b_{t}) \right]$$
(8)

$$= \beta \cdot \mathbb{E}_{\tau \sim \Pr_{\mu}^{\pi^{1}, \pi^{2}}} \left[\sum_{t=0}^{\infty} \gamma^{t} A^{1}(s_{t}, a_{t}, b_{t}) \sum_{a'_{t}} \nabla_{\theta^{1}} \pi_{1}(a'_{t}|s_{t}) Q^{2}(s_{t}, a'_{t}, b_{t}) \right]$$
(9)

$$= \beta \cdot \mathbb{E}_{\tau} \left[\sum_{t=0}^{\infty} \gamma^{t} A_{t}^{1} \sum_{a'_{t}} \left[\pi_{1}(a'_{t}|s_{t}) Q_{t}^{2} \nabla_{\theta_{1}} \log \pi_{1}(a'_{t}|s_{t}) + \pi_{1}(a'_{t}|s_{t}) \nabla_{\theta_{1}} Q_{t}^{2} \right] \right]$$
(10)

$$= \beta \cdot \mathbb{E}_{\tau \sim \Pr_{\mu}^{\pi^{1}, \pi^{2}}} \left[\sum_{t=0}^{\infty} \gamma^{t} A_{t}^{1} \mathbb{E}_{a'_{t} \sim \pi^{1}(a'_{t}|s_{t})} \left[Q_{t}^{2} \nabla_{\theta_{1}} \log \pi_{1}(a'_{t}|s_{t}) + \nabla_{\theta_{1}} Q_{t}^{2} \right] \right]$$
(11)

$${}^{2}\mathbb{E}_{a_{t},b_{t}}[A^{1}(s_{t},a_{t},b_{t})f(s_{t})] = f(s_{t})\mathbb{E}_{a_{t},b_{t}}[Q^{1}(s_{t},a_{t},b_{t}) - V^{1}(s_{t})] = f(s_{t})(V^{1}(s_{t}) - V^{1}(s_{t})) = 0$$

In first term of line 11, we can use the fact that $\mathbb{E}_{a'_t \sim \pi^1(a'_t|s_t)} \left[V^2(s_t) \nabla_{\theta_1} \log \pi_1(a'_t|s_t) \right] = 0$, which allows us to update Q^2 to A^2 :

$$\beta \cdot \mathbb{E}_{\tau \sim \Pr_{\mu}^{\pi^1, \pi^2}} \left[\sum_{t=0}^{\infty} \gamma^t A^1(s_t, a_t, b_t) \mathbb{E}_{a_t' \sim \pi^1(a_t'|s_t)} \left[Q^2(s_t, a_t', b_t) \nabla_{\theta_1} \log \pi_1(a_t'|s_t) \right] \right]$$
(12)

$$= \beta \cdot \mathbb{E}_{\tau \sim \Pr_{\mu}^{\pi^{1}, \pi^{2}}} \left[\sum_{t=0}^{\infty} \gamma^{t} A^{1}(s_{t}, a_{t}, b_{t}) \mathbb{E}_{a'_{t} \sim \pi^{1}(a'_{t}|s_{t})} \left[A^{2}(s_{t}, a'_{t}, b_{t}) \nabla_{\theta_{1}} \log \pi_{1}(a'_{t}|s_{t}) \right] \right]$$
(13)

Second term of line 11:

$$\beta \mathbb{E}_{\tau \sim \Pr_{\mu}^{\pi^{1}, \pi^{2}}} \left[\sum_{t=0}^{\infty} \gamma^{t} A^{1}(s_{t}, a_{t}, b_{t}) \mathbb{E}_{a'_{t} \sim \pi^{1}(a'_{t}|s_{t})} \left[\nabla_{\theta_{1}} Q^{2}(s_{t}, a'_{t}, b_{t}) \right] \right]$$
(14)

$$= \beta \mathbb{E}_{\tau \sim \Pr_{\mu}^{\pi^{1}, \pi^{2}}} \left[\sum_{t=0}^{\infty} \gamma^{t} A_{t}^{1}(s_{t}, a_{t}, b_{t}) \mathbb{E}_{a'_{t} \sim \pi^{1}(a'_{t}|s_{t})} \left[\nabla_{\theta^{1}} \left[r_{t}^{2} + \gamma \cdot \mathbb{E}_{s'} \left[V^{2}(s') \right] \middle| s_{t}, b_{t} \right] \right] \right]$$
(15)

$$= \beta \mathbb{E}_{\tau \sim \Pr_{\mu}^{\pi^{1}, \pi^{2}}} \left[\sum_{t=0}^{\infty} \gamma^{t} A_{t}^{1} \mathbb{E}_{\tau'} \left[\sum_{k=t+1}^{\infty} \gamma^{k-t} A_{k}^{2} \left(\nabla_{\theta^{1}} [\log(\pi^{1}(a'_{k}|s'_{k}) \hat{\underline{\pi}^{2}}(b'_{k}|s'_{k}))] \right) \middle| s_{t}, b_{t} \right] \right]$$
(16)

$$= \beta \mathbb{E}_{\tau \sim \Pr_{\mu}^{\pi^{1}, \pi^{2}}} \left[\sum_{t=0}^{\infty} \gamma^{t} A_{t}^{1} \mathbb{E}_{\tau' \sim \Pr_{\mu}^{\pi^{1}, \pi^{2}}} \left[\sum_{k=t+1}^{\infty} \gamma^{k-t} A_{k}^{2} \nabla_{\theta^{1}} \log \pi^{1}(a'_{k} | s'_{k}) \middle| s_{t}, b_{t} \right] \right]$$
(17)

Summing up lines 13 and 17, we get

$$\beta \cdot \mathbb{E}_{\substack{\tau \sim \Pr_{\mu}^{\pi^1, \pi^2}}} \left[\sum_{t=0}^{\infty} \gamma^t A_t^1 \mathbb{E}_{\tau' \sim \Pr_{\mu}^{\pi^1, \pi^2}} \left[\sum_{k=t}^{\infty} \gamma^{k-t} A_k^2 \nabla_{\theta^1} \log \pi^1(a_k' | s_k') \middle| s_t' = s_t, b_t' = b_t \right] \right]$$
(18)

$$= \beta \cdot \mathbb{E}_{\tau \sim \Pr_{\mu}^{\pi^1, \pi^2}} \left[\sum_{t=0}^{\infty} A_t^1 \mathbb{E}_{\tau' \sim \Pr_{\mu}^{\pi^1, \pi^2}} \left[\sum_{k=t}^{\infty} \gamma^k A_k^2 \nabla_{\theta^1} \log \pi^1(a_k' | s_k') \middle| s_t, b_t \right] \right]$$
(19)

(20)

We use Bellman equation in line (15), and the policy gradient theorem in line (16). Denoting the inner policy gradient³ by $\nabla_t(\tau')$, we rewrite 20 as

$$\beta \sum_{\tau} \Pr(\tau) \sum_{t=0}^{\infty} A^{1}(s_{t}, a_{t}, b_{t}) \sum_{\tau'} \nabla_{t}(\tau') \Pr(\tau' \mid s'_{t} = s_{t}, b'_{t} = b_{t})$$
(21)

$$= \beta \sum_{t=0}^{\infty} \sum_{\tau} \Pr(\tau) A^{1}(s_{t}, a_{t}, b_{t}) \sum_{\tau'} \mathbf{\nabla}_{t}(\tau') \Pr(\tau' \mid \tau)$$
(22)

$$= \beta \sum_{t=0}^{\infty} \sum_{\tau} \Pr(\tau) \sum_{\tau'} A^{1}(s'_{t}, a_{t}, b'_{t}) \boldsymbol{\nabla}_{t}(\tau') \Pr(\tau' \mid \tau)$$
(23)

$$= \beta \sum_{t=0}^{\infty} \sum_{\tau} \Pr(\tau) \sum_{\tau'} \left(\sum_{a_t^*} \Pr(a_t^* \mid \tau) A^1(s_t', a_t^*, b_t') \right) \boldsymbol{\nabla}_t(\tau') \Pr(\tau' \mid \tau)$$
(24)

$$= \beta \sum_{t=0}^{\infty} \sum_{\tau} \Pr(\tau) \sum_{\tau'} \left(\sum_{a_t^*} \Pr(a_t^* \mid \tau) \Pr(\tau' \mid \tau) A^1(s_t', a_t^*, b_t') \right) \boldsymbol{\nabla}_t(\tau')$$
(25)

$$= \beta \sum_{t=0}^{\infty} \sum_{\tau} \Pr(\tau) \sum_{\tau'} \sum_{a_t^*} \Pr(\tau', a_t^* \mid \tau) A^1(s_t', a_t^*, b_t') \nabla_t(\tau')$$
 (26)

$$= \beta \sum_{t=0}^{\infty} \sum_{\tau'} \sum_{a_t^*} \sum_{\tau} \Pr(\tau) \Pr(\tau', a_t^* \mid \tau) A^1(s_t', a_t^*, b_t') \nabla_t(\tau')$$
(27)

$$= \beta \sum_{t=0}^{\infty} \sum_{\tau'} \sum_{a_t^*} \Pr(\tau', a_t^*) A^1(s_t', a_t^*, b_t') \nabla_t(\tau')$$
(28)

$$= \beta \mathbb{E}_{\tau \sim \Pr_{\mu}^{\pi^1, \pi^2}} \left[\sum_{t=0}^{\infty} \mathbb{E}_{a'_t \sim \pi^1(\cdot|s_t)} \left[A^1(s_t, a'_t, b_t) \right] \boldsymbol{\nabla}_t(\tau) \right]$$
(29)

$$= \beta \mathbb{E}_{\tau \sim \Pr_{\mu}^{\pi^{1}, \pi^{2}}} \left[\sum_{t=0}^{\infty} \mathbb{E}_{a'_{t}} \left[A^{1}(s_{t}, a'_{t}, b_{t}) \right] \sum_{k=t}^{\infty} \gamma^{k} A^{2}(s_{k}, a_{k}, b_{k}) \nabla_{\theta^{1}} \log \pi^{1}(a_{k}|s_{k}) \right]$$
(30)

$$= \beta \mathbb{E}_{\tau \sim \Pr_{\mu}^{\pi^{1}, \pi^{2}}} \left[\sum_{t=0}^{\infty} \gamma^{t} \left(\sum_{k \leq t} \mathbb{E}_{a'_{k}} \left[A^{1}(s_{k}, a'_{k}, b_{k}) \right] \right) A^{2}(s_{t}, a_{t}, b_{t}) \nabla_{\theta^{1}} \log \pi^{1}(a_{t}|s_{t}) \right]$$
(31)

In line 22, we use the fact that since $\nabla_t(\tau')$, by definition, only depends on b_t, s_t from τ , $\Pr(\tau' \mid \tau) = \Pr(\tau' | s_t' = s_t, b_t' = b_t)$. In line 23, we switch from $A^1(s_t, a_t, b_t)$ to $A^1(s_t', a_t, b_t')$ in the summation. This is allowed by the fact that for each τ' with $s_t' \neq s_t$ or $b_t' \neq b_t$, $\Pr(\tau' \mid \tau) = \Pr(\tau' \mid b_t, s_t) = 0$. In line 24, we use the fact that:

$$\sum_{a_t^*} \Pr(a_t^* \mid \tau) A^1(s_t', a_t^*, b_t') = 1 \cdot A^1(s_t', a_t, b_t') + \sum_{a_t^* \neq a_t} 0 \cdot A^1(s_t', a_t^*, b_t') = A^1(s_t', a_t, b_t').$$

In line 26, we use the fact that τ' and a_t^* are independent knowing τ .

Let the trajectory length be T. An unbiased estimator of (31) would require $O(T^2)$ steps to be generated. To avoid this computationnal overhead, jit-AdAlign approximates this term as

$$\beta \mathbb{E}_{\tau \sim \Pr_{\mu}^{\pi^1, \pi^2}} \left[\sum_{t=0}^{\infty} \gamma^t \left(\sum_{k \leq t} A^1(s_k, a_k', b_k) \right) A^2(s_t, a_t, b_t) \nabla_{\theta^1} \log \pi^1(a_t | s_t) \right]$$

In the future, one could use TD-learning to approximate the term with less bias.

$$^{3}\sum_{k=t}^{\infty}\gamma^{k}A_{k}^{2}\nabla_{\theta^{1}}\log\pi^{1}(a_{k}|s_{k})$$

Hyperparameter	Value
Optimizer	Adam
Learning Rate	1e-6
Batch Size	64
Number of Rounds	10
Self-play Used	Yes
LoRA Rank	32
Lora α	64
LoRA Dropout	None
Data Type	bfloat16
Sampling Temperature	1.0
Reward Normalization Constant	5.0
Entropy Coefficient	0.01
KL Coefficient	0
Discount Factor	1
jit-AA eta	0.5
jit-AA γ	0.9
Buffer Prob.	0.5

Table 1: Hyperparameters for IPD experiments.

Hyperparameter	Value
Optimizer	Adam
Learning Rate	3e-6
Batch Size	64
Number of Rounds	10
Self-play Used	Yes
LoRA Rank	32
Lora α	64
LoRA Dropout	None
Data Type	bfloat16
Sampling Temperature	1.0
Reward Normalization Constant	100.0
Entropy Coefficient	0
KL Coefficient	0
Discount Factor	1
jit-AA eta	1
jit-AA γ	0.9
Buffer Prob.	0.5

Table 2: Hyperparameters for Split Games.

B EXPERIMENTAL DETAILS

C IPD

	Cooperate (C)	Defect (D)
Cooperate (C)	(3, 3)	(0,5)
Defect (D)	(5,0)	(1,1)

D SPLIT GAMES

This family of games feature two agents who, in each round, may briefly communicate and then simultaneously propose how to split a fixed resource (most commonly 10 coins). Rewards are the amount kept multiplied by an agent's per-unit value. The starting speaker alternates deterministically

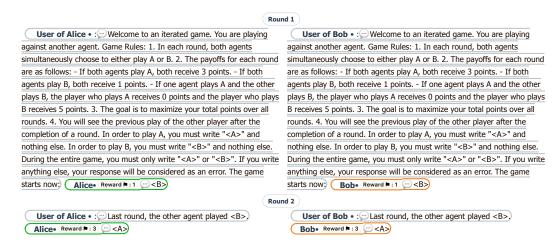


Figure 7: IPD introduction prompt example.

across rounds. Importantly, actions and private values from the previous round are revealed to both agents. This ensures that retaliatory strategies can take place.

Communication is optional and variant-dependent: some settings encourage rich messaging to share private information, while others remove messaging entirely to focus on allocation behavior.

Proportional splitting is used when the two proposals exceed the available total: allocations are scaled proportionally rather than discarded. This preserves a useful learning signal even when agents over-claim.

We now introduce the variants roughly in terms of complexity and credit alignment difficulty.

SPLIT NO-COMM.

- Single item type (coins).
- Values are public.
- No communication; agents go straight to making split proposals, with the starting player alternating deterministically.
- **Motivation:** mirrors no-communication setups (e.g., Advantage Alignment) while keeping the split decision nontrivial.
- No-Press Split: 10-1-Exclusive: values are either 1 or 10 and mutually exclusive. If one agent gets 10, the other gets 1 (and vice versa).
- No-Press Split: 10-1-Ties: values are either 1 or 10 and uncorrelated.
- No-Press Split: 1-20-Stochastic: values range from 1 to 20 (inclusive), are random and uncorrelated.

TRUST-AND-SPLIT RPS (TAS-RPS)

- Single item type (coins).
- Each round, a rock-paper-scissors hand draw creates a strong asymmetry: the winner's per-coin value is 10, the loser's is 1.
- Each agent initially sees only their own hand and must communicate to coordinate an optimal split.
- **Motivation:** enforce large value disparity so one's own value reveals little about the other's (avoiding ceiling effects) and incentivize meaningful communication.

932

933

934

935

936

937

938

939

940

941

942

943

944

945

946

947

948 949

950

951

952

953

954

955

956

957 958

Round 1 User of Alice • : Welcome to an iterated game. You are Alice. The User of Bob • : Welcome to an iterated game. You are Bob. The other agent is Bob. Setup: 1. The game consists of multiple independent other agent is Alice. Setup: 1. The game consists of multiple independent rounds. 2. In each round, there are multiple items to split between the rounds. 2. In each round, there are multiple items to split between the two agents. 3. Both agents are assigned a per-item value between 1 and two agents. 3. Both agents are assigned a per-item value between 1 and 20 (inclusive) in each round. 4. You can observe per-item values of both 20 (inclusive) in each round. 4. You can observe per-item values of both agents, 5. Because assignments are random, both agents are equally agents, 5. Because assignments are random, both agents are equally likely to have same expected per-item value. Protocol: 1. Both agents likely to have same expected per-item value. Protocol: 1. Both agents simultaneously propose the amount of each item they will keep. 2. If the simultaneously propose the amount of each item they will keep. 2. If the total sum of proposals is less than or equal to the item quantity, both total sum of proposals is less than or equal to the item quantity, both agents receive their proposed amounts. 3. If the total sum of proposals agents receive their proposed amounts. 3. If the total sum of proposals exceeds the item quantity, they are allocated proportionally. 4. Your points exceeds the item quantity, they are allocated proportionally. 4. Your points for the round = (amount you receive per item) x (your per-item value for for the round = (amount you receive per item) x (your per-item value for that round), added across all items. 5. Points are accumulated across that round), added across all items. 5. Points are accumulated across rounds. Your goal: Maximize your total points over the whole game. A rounds. Your goal: Maximize your total points over the whole game. A New Round Begins The items to split are 10 hats, 10 books, 10 balls. Your New Round Begins The items to split are 10 hats, 10 books, 10 balls. Your per-item values are hats=10, books=1, balls=1 and Bob's per-item values per-item values are hats=1, books=10, balls=1 and Alice's per-item values are hats=1, books=10, balls=1. Submit Your Proposal Respond as are hats=10, books=1, balls=1. Submit Your Proposal Respond as Proposal: x hats, y books, z balls where x: 0-10 (integer), y: 0-10 Proposal: x hats, y books, z balls where x: 0-10 (integer), y: 0-10 (integer), z: 0-10 (integer). Alice Reward №:55 Proposal: 10 hats, 0 (integer), z: 0-10 (integer). Bob• Reward №: 10 Proposal: 10 hats, 0 books, 10 balls 10 hats, 0 books, 10 balls books, 10 balls User of Alice • : Last Round Summary: - Items to split: 10 hats, User of Bob • : Last Round Summary: - Items to split: 10 hats, 10 10 books, 10 balls - Your per-item values: hats=10, books=1, balls=1 books, 10 balls - Your per-item values: hats=1, books=10, balls=1

Bob's per-item values: hats=1, books=10, balls=1 - You proposed: 10

0 books, 10 balls - Bob earned: 10.0 points - Round Complete. A New

hats, 0 books, 10 balls - You earned: 55.0 points - Bob proposed: 10 hats,

Round Begins The items to split are 10 hats, 10 books, 10 balls. Your per-

item values are hats=1, books=10, balls=1 and Bob's per-item values are

hats=1, books=1, balls=10. Submit Your Proposal Respond as Proposal: x

hats, y books, z balls where x: 0-10 (integer), y: 0-10 (integer), z: 0-10

(integer). Alice Reward : 55 Proposal: 0 hats, 10 books, 10 balls

Figure 8: Split No-Comm introduction prompt.

10 books, 10 balls

Alice's per-item values: hats=10, books=1, balls=1 - You proposed: 10

hats, 0 books, 10 balls - You earned: 10.0 points - Alice proposed: 10

hats, 0 books, 10 balls - Alice earned: 55.0 points - Round Complete. A

are hats=1, books=10, balls=1. Submit Your Proposal Respond as

Proposal: x hats, y books, z balls where x: 0-10 (integer), y: 0-10

New Round Begins The items to split are 10 hats, 10 books, 10 balls. Your

per-item values are hats=1, books=1, balls=10 and Alice's per-item values

(integer), z: 0-10 (integer). Bob• Reward ■: 65 Proposal: 10 hats,

981

982

983

984

985

986

987

988

989

990

991

992

993

994

995

996

997

998

999

1000

1001

1002

1003

1004

1005 1006

1007

1008

1009

1010

1011

1012

1013

1014

1015

1016

1017 1018

User of Alice • : Welcome to an iterated game. You are Alice. The User of Bob • : Welcome to an iterated game. You are Bob. The other agent is Bob. Setup: 1. The game has multiple independent rounds. other agent is Alice. Setup: 1. The game has multiple independent rounds. 2. In each round, there are 10 coins to split between the two agents, 3. 2. In each round, there are 10 coins to split between the two agents, 3. Each agent's per-coin value for that round is determined as follows: - Both Each agent's per-coin value for that round is determined as follows: - Both agents are randomly assigned a rock, paper or scissors hands - Rock has agents are randomly assigned a rock, paper or scissors hands - Rock has the upper hand over scissors, scissors has the upper hand over paper and the upper hand over scissors, scissors has the upper hand over paper and paper has the upper hand over rock. - The agent with the upper hand has paper has the upper hand over rock. - The agent with the upper hand has a per-coin value of 10. - The agent with the lower hand has a per-coin a per-coin value of 10. - The agent with the lower hand has a per-coin value of 1. 4. You only see your own hand, but you may communicate it in value of 1. 4. You only see your own hand, but you may communicate it in messages and infer your value based on the other agent's hand. 5. Over messages and infer your value based on the other agent's hand. 5. Over many rounds both agents are equally likely to have the upper and lower many rounds both agents are equally likely to have the upper and lower hand. Protocol: 1. At the start of the round, one agent begins the hand, Protocol: 1. At the start of the round, one agent begins the conversation. The starting role alternates each round. 2. Agents exchange conversation. The starting role alternates each round. 2. Agents exchange a short chat (1 messages per round per agent) to negotiate how to split a short chat (1 messages per round per agent) to negotiate how to split the 10 coins. - Use this chat to communicate your hand so that both the 10 coins. - Use this chat to communicate your hand so that both agents can determine their per-coin values. 3. After the chat, both agents agents can determine their per-coin values. 3. After the chat, both agents simultaneously propose how many coins they keep. 4. If the total sum of simultaneously propose how many coins they keep. 4. If the total sum of proposals is less than or equal to 10, both agents receive their proposals. proposals is less than or equal to 10, both agents receive their proposals. 5. If the total sum of proposals exceeds 10, the coins are allocated 5. If the total sum of proposals exceeds 10, the coins are allocated proportionally. 6. Your points for the round = (coins you receive) x (your proportionally. 6. Your points for the round = (coins you receive) x (your per-coin value for that round). 7. The points are accumulated across per-coin value for that round). 7. The points are accumulated across rounds. Your goal: Maximize your total points over the whole game. A rounds. Your goal: Maximize your total points over the whole game. A New Round Begins Your hand is rock. Wait for Bob to send a message... New Round Begins Your hand is paper. Send your message now in Alice • Reward № : 5 User of Alice • : Bob said: <message>...</message> (<=500 chars), Bob• Reward ■:50 <message>My hand is paper. Let's communicate our hands so we can <message>My hand is paper. Let's communicate our hands so we can determine our per-coin values.</message> Send your message now in determine our per-coin values.</message> User of Bob .: (... Bob • Reward №: 50 User of Bob • : Alice said: <message>...</message> (<=500 chars). Alice • Reward №:5 <message>My hand is rock. Since rock loses to paper, I have a per-coin <message>My hand is rock. Since rock loses to paper, I have a per-coin value of 1. Let's split the coins evenly to maximize our total points.</ value of 1. Let's split the coins evenly to maximize our total points.</ message> User of Alice • : Submit your proposal Respond with message> Submit your proposal Respond with <coins_to_self> x </ <coins_to_self> x </coins_to_self> where x is an integer in [0, 10]. coins_to_self> where x is an integer in [0, 10], Bob• Reward №:50 ... <coins_to_self>5</coins_to_self> User of Alice • : Last Round Summary: - Your hand: rock - Bob's User of Bob • : Last Round Summary: - Your hand: paper - Alice's hand: paper - Your value per coin: 1.0 - Bob's value per coin: 10.0 - You hand: rock - Your value per coin: 10.0 - Alice's value per coin: 1.0 - You proposed: 5 coins - You earned: 5.0 points - Bob proposed: 5 coins - Bob proposed: 5 coins - You earned: 50.0 points - Alice proposed: 5 coins earned: 50.0 points - Round Complete. A New Round Begins Your hand is Alice earned: 5.0 points - Round Complete. A New Round Begins Your hand is paper. Wait for Alice to send a message... Bob • Reward ₱:10 scissors. Send your message now in <message>...</message> (<=500 chars). Alice Reward ►: 0 — < message > My hand is scissors. Let's User of Bob • : Alice said: < message > My hand is scissors. share hands to determine values. Since I have scissors, I need to know Let's share hands to determine values. Since I have scissors, I need to know Bob's hand to decide the split.</message> Send your message now for Bob to send a message... Alice Reward ■:0 User of Alice in <message>...</message> (<=500 chars). Bob• Reward ■:10 : Bob said: <message>My hand is paper. Since paper beats scissors, I <message>My hand is paper. Since paper beats scissors, I have a perhave a per-coin value of 10. Let's split the coins based on our values to coin value of 10. Let's split the coins based on our values to maximize our maximize our total points.</message> Submit your proposal Respond total points.</message> User of Bob • : Submit your proposal with <coins_to_self> x </coins_to_self> where x is an integer in [0, 10], Respond with <coins_to_self> x </coins_to_self> where x is an integer in [0, 10]. Bob• Reward : 10 - < coins_to_self > 10 < / coins_to_self > Alice • Reward ■: 0 : Coins_to_self>0</coins_to_self>

Figure 9: TAS-RPS introduction prompt.

1028

1029

1030

1031

1032

1033

1034

1035

1036

1037

1038

1039

1040

1041

1042

1043

1044

1045

1046

1047

1048

1049

1050

1051

1052

1053

1054

1055

1056

1057

1059

1061

1062 1063

1064

1067

1068 1069

1070

1077 1078 1079

User of Alice • : Welcome to an iterated game. You are Alice. The User of Bob • : Welcome to an iterated game. You are Bob. The other agent is Bob. Setup: 1. The game has multiple independent rounds. other agent is Alice. Setup: 1. The game has multiple independent rounds. 2. In each round, there are 10 coins to split between the two agents. 3. 2. In each round, there are 10 coins to split between the two agents. 3. Both agents are assigned a per-coin value between 1 and 10 (inclusive) in Both agents are assigned a per-coin value between 1 and 10 (inclusive) in each round. 4. You can only observe your own per-coin value. 5. Because each round, 4. You can only observe your own per-coin value, 5. Because assignments are random, both agents are equally likely to have same assignments are random, both agents are equally likely to have same expected per-coin value. Protocol: 1. At the start of the round, one agent expected per-coin value. Protocol: 1. At the start of the round, one agent begins the conversation. The starting role alternates each round. 2. begins the conversation. The starting role alternates each round. 2. Agents exchange a short chat (1 messages per round per agent) to Agents exchange a short chat (1 messages per round per agent) to negotiate how to split the coins. - Use this chat to communicate your negotiate how to split the coins. - Use this chat to communicate your private per-coin value to make informed proposals, 3. After the chat, both private per-coin value to make informed proposals, 3, After the chat, both agents simultaneously propose how many coins they keep. 4. If the total agents simultaneously propose how many coins they keep. 4. If the total sum of proposals is less than or equal to 10, both agents receive their sum of proposals is less than or equal to 10, both agents receive their proposals. 5. If the total sum of proposals exceeds 10, the coins are proposals. 5. If the total sum of proposals exceeds 10, the coins are allocated proportionally. 6. Your points for the round = (coins you receive) allocated proportionally. 6. Your points for the round = (coins you receive) x (your per-coin value for that round). 7. Points are accumulated across x (your per-coin value for that round). 7. Points are accumulated across rounds. Your goal: Maximize your total points over the whole game. A rounds. Your goal: Maximize your total points over the whole game. A New Round Begins Your per-coin value is 9. You don't know Bob's value New Round Begins Your per-coin value is 2. You don't know Alice's value yet) Bob• Reward ■: 8 User of Bob • : Alice said: yet. Send your message now in <message>...</message> (<=200 chars). Alice • Reward №: 45 — <message>Hi Bob, I have a per-coin <message>Hi Bob, I have a per-coin value of 9. To maximize our points, value of 9. To maximize our points, let's keep the total proposals at 10 or let's keep the total proposals at 10 or less. How about you propose 4 coins less. How about you propose 4 coins and I take 6?</message> User and I take 6?</message> Send your message now in <message>...</ message> (<=200 chars), Bob• Reward №: 8 💭 (...) — <message>I of Alice • : → Alice • Reward • : 45 → User of Alice • : → Bob said: <message>I have a per-coin value of 2. To be efficient, I propose 5 have a per-coin value of 2. To be efficient, I propose 5 coins for me and 5 for you (total 10).</message> User of Bob • : Submit your coins for me and 5 for you (total 10).</message> Submit your proposal Respond with <coins_to_self> x </coins_to_self> where x is an integer in proposal Respond with <coins_to_self> x </coins_to_self> where x is an [0, 10], Alice Reward : 45 <-> < coins_to_self > 5 </ coins_to_self > integer in [0, 10], Bob• Reward ■: 8 (...) (...) <coins_to_self> 4 </ coins_to_self> Round 2 User of Alice • : In the previous round, Bob had a lower value and User of Bob • : In the previous round, Alice had a higher value proposed 4 coins. A New Round Begins Your per-coin value is 6. You don't and proposed 5 coins. A New Round Begins Your per-coin value is 3, You know Bob's value yet, Alice • Reward ■: 30 User of Alice • : ... don't know Alice's value yet. Send your message now in <message>...< Bob said: <message>My per-coin value is 3. I propose I take 4 coins and per-coin value is 3. I propose I take 4 coins and you take 6 (total 10) to you take 6 (total 10) to be efficient.</message> Send your message now in <message>...</message> (<=200 chars). Alice• (Reward №: 30) be efficient.</message> User of Bob • : Bob• Reward № : 15 <message>Hi Bob, I have a per-coin value of 6. My proposal is to split User of Bob • : Alice said: < message > Hi Bob, I have a per-coin evenly at 5 coins each. This way, both of us maximize our points value of 6. My proposal is to split evenly at 5 coins each. This way, both of efficiently.</message> After the chat, I propose 5 coins for myself and </ us maximize our points efficiently.</message> After the chat, I propose 5 message> User of Alice • : Submit your proposal Respond with coins for myself and </message> Submit your proposal Respond with <coins_to_self> x </coins_to_self> where x is an integer in [0, 10]. <coins_to_self> x </coins_to_self> where x is an integer in [0, 10],

Figure 10: TAS introduction prompt.

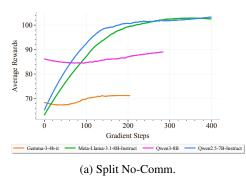
TRUST-AND-SPLIT (TAS)

Alice (Reward №: 30) — <coins_to_self> 5 </coins_to_self>

 Single item type (coins); each round, each agent's per-coin value is independently sampled in a broad range (e.g., 1–20).

- Each agent observes only their own value; they may use short messages to share and negotiate.
- **Motivation:** a simple blend that tests whether agents learn to exchange private information and coordinate proportional, value-aware splits.

E SUM OF REWARDS



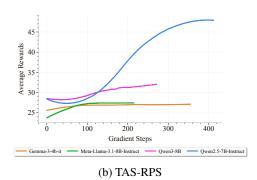


Figure 11: By changing the objective of the agents to be the total welfare (thereby removing the social dilemma aspect), we show that Qwen2.5-7B-Instruct is able to learn to play both games well.