

Interactive Learning for LLM Reasoning

Hehai Lin^{1*}, Shilei Cao^{2*}, Sudong Wang¹, Haotian Wu¹,
Minzhi Li^{4,5}, Linyi Yang³, Juepeng Zheng², Chengwei Qin^{1†}

¹The Hong Kong University of Science and Technology (Guangzhou),

²Sun Yat-sen University, ³Southern University of Science and Technology,

⁴National University of Singapore, ⁵Institute for Infocomm Research (I²R), A*STAR

Abstract

Existing multi-agent learning approaches explicitly foster collaboration among Large Language Models (LLMs) to build stronger multi-agent systems (MAS), yet they still rely on re-executing the MAS during inference. This contrasts with human cognition, wherein individuals can internalize insights from interactions to improve later independent reasoning. To investigate whether multi-agent interaction can enhance LLMs’ independent problem-solving ability, we propose ILR (Interactive Learning for LLM Reasoning), a co-learning framework that integrates Dynamic Interaction and Perception Calibration. Dynamic Interaction adaptively selects cooperative or competitive strategies based on question difficulty and model capability, after which LLMs exchange information via Idea3 framework (Idea Sharing, Idea Analysis, and Idea Fusion), an interaction paradigm simulating human discussion, before producing final answers. Perception Calibration employs Group Relative Policy Optimization (GRPO) while integrating one LLM’s reward characteristics into another’s to strengthen interaction cohesion. We evaluate the effectiveness of ILR across three LLMs from two model families of varying scales on five mathematical, one coding, one general question answering, and one scientific reasoning benchmarks. We further investigate the advantage of Dynamic Interaction (i.e., boosting the robustness of stronger LLMs and surpassing pure strategy), and the scalability of ILR beyond two-model interactions. Code is available at <https://github.com/linhh29/Interactive-Learning-for-LLM-Reasoning>.

1 Introduction

Efforts to enhance the reasoning capabilities of Large Language Models (LLMs) have largely relied on training paradigms such as Supervised Fine-Tuning (SFT) (Achiam et al., 2023), Preference

Learning (PL) (Rafailov et al., 2023), and Reinforcement Learning (RL) (Schulman et al., 2017; Guo et al., 2025; Zeng et al., 2025). These methods allow LLMs to iteratively interact with data and refine their behavior, essentially engaging in trial-and-error learning to acquire problem-solving skills, which can be viewed as self-learning for LLMs. However, real-world knowledge acquisition is rarely an isolated activity (Bloembergen et al., 2015; Canese et al., 2021). Humans continuously exchange knowledge through collaborative learning, as in peer discussions within classroom settings. While *single-agent learning* (self-learning) serves as the foundation of human education, *multi-agent learning* represents a more advanced and often more effective paradigm: multiple learners bring diverse perspectives, challenge each other’s reasoning, and provide mutual feedback, ultimately leading to deeper understanding and more robust solutions (Kahveci and Imamoglu, 2007; Hsiung, 2012; Zambrano et al., 2019; Mende et al., 2021). The same principle suggests that multi-agent learning can benefit LLMs¹: by exposing models to diverse reasoning strategies and peer-based feedback, it may help them overcome individual blind spots and develop stronger problem-solving abilities.

Recent studies have explored multi-agent learning. For example, MALT (Motwani et al., 2024) designs a sequential multi-agent system (MAS) consisting of Generator, Verifier, and Refiner, each independently trained to sample trajectories. ReMA (Wan et al., 2025) introduces a hierarchical framework with a high-level agent responsible for problem decomposition and a low-level agent for concrete steps, trained alternately to achieve complementary expertise. MAPoRL (Park et al., 2025) proposes a Post-Co-Training framework to enhance collaboration alignment through debate. However, during inference, these methods are required to re-

*Equal contribution

†Corresponding to chengweiqin@hkust-gz.edu.cn

¹In our paper, “Agent” and “LLM” refer to the same entity

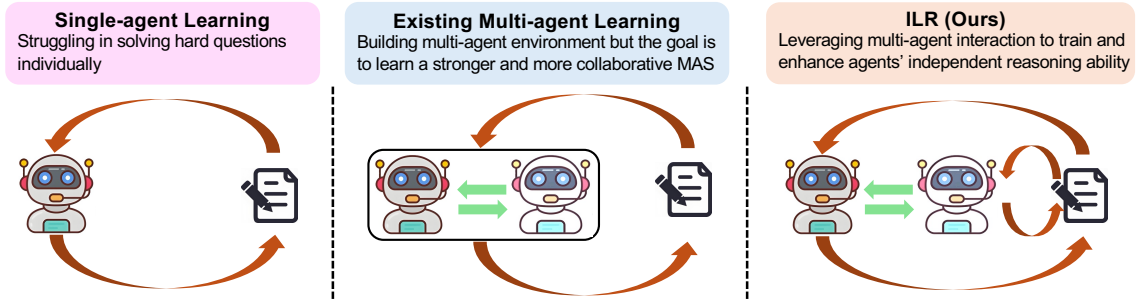


Figure 1: Conceptual comparison of ILR, existing Multi-agent Learning, and Single-agent Learning. With ILR, LLMs engage in multi-agent learning at training time and can resolve questions independently at inference time.

execute the MAS to obtain final solutions, a process misaligned with human cognition, where individuals improve reasoning through peer interactions and subsequently solve problems independently.

In this paper, we address this gap by treating each agent as an autonomous entity and investigating whether multi-agent learning can enhance an LLM’s individual problem-solving capacity (see Figure 1). We propose **ILR (Interactive Learning for LLM Reasoning)**, a co-learning framework consisting of two key components: *Dynamic Interaction* and *Perception Calibration*.

The *Dynamic Interaction* module simulates human discussion. For “Dynamic”, when confronted with complex problems, humans tend to cooperate, whereas for simpler problems, they often compete to identify the most efficient solution (Richard et al., 2002; Schneider et al., 2011). To emulate this behavior, an LLM estimates question difficulty through self-ranking and applies Item Response Theory (Cai et al., 2016; Benedetto et al., 2023) to calculate the probability of solving it independently. If the probability is low, the model engages in cooperation; otherwise, it chooses competition. For “Interaction”, we design a novel **Idea3** framework, comprising three sequential stages: Idea Sharing (each LLM proposes its own solution), Idea Analysis (each LLM analyzes and reflects on the peer’s solution), and Idea Fusion (the insights are synthesized into a refined and potentially novel solution). Following Dynamic Interaction, the *Perception Calibration* module is applied. Prior work (Ma et al., 2024; Park et al., 2025) has shown that incorporating tailored reward signals can effectively guide LLMs toward better multi-agent learning. Instead of relying on predefined signals, we propose a fully automated mechanism that integrates one LLM’s reward distribution characteristics, derived from answer group sampling on the same input, into another LLM’s reward function. We then employ the

standard GRPO (Shao et al., 2024) to update each LLM based on calibrated rewards. This calibration allows LLMs to perceive the quality of peer-generated solutions and adapt their reasoning accordingly. We further validate the generalizability of ILR on another RL algorithm in Appendix D.8.

We evaluate the effectiveness of ILR across three LLMs from two model families of varying scales on five mathematical, one coding, one general question answering, and one scientific reasoning benchmarks. Experimental results demonstrate that ILR consistently outperforms both single-agent learning and multi-agent learning. Our investigation further reveals three findings: (1) Dynamic Interaction enhances the robustness of stronger LLMs during multi-agent inference scenarios. Analyzing and reflecting on the peer’s solutions reduces the probability of being misled by weaker LLMs when exchanging information. (2) Dynamically determining interaction types can boost the efficacy of multi-agent learning and surpass pure cooperation or competition strategies. (3) ILR has a positive scaling trend with a larger number of LLMs. Our main contributions are summarized as follows:

- Unlike prior work that focuses on improving system-level collaboration in multi-agent learning, we are the first to explore its potential for enhancing an individual LLM’s independent reasoning capability.
- Inspired by human interaction, we design a novel multi-agent learning framework ILR, which comprises Dynamic Interaction and Perception Calibration.
- Experimental results validate ILR’s superiority over existing baselines. We further discover the advantage of Dynamic Interaction and the Scalability of ILR, which offer valuable insights toward future MAS design.

2 Related Work

Multi-agent learning (Busoniu et al., 2006; Han et al., 2024; Li et al., 2024b) first requires designing a multi-agent system that defines interaction paradigms among multiple agents, such as equi-level (Chan et al., 2023), hierarchical (Gronauer and Diepold, 2022), or nested structures (Zhao et al., 2025). Then, within this architectural framework, distinct agents engage in interactive sampling to acquire experience, which subsequently undergoes optimization through learning algorithms. Therefore, we systematically review prior works from the following two perspectives: Multi-Agent Communication (interactive paradigms) and Multi-Agent Training (optimization methods).

2.1 Multi-Agent Communication

Traditionally, researchers employ recurrent neural networks (RNNs) as agents and utilize attention mechanisms to facilitate communication (Yu et al., 2019; Ding et al., 2024; Sun et al., 2024). For instance, TarMAC leverages multi-head attention to enable agents to learn both message content and targeted recipient (Das et al., 2019). After the emergence of LLMs, researchers develop numerous explainable prompt-based multi-agent communication. Notable examples include Debate, where multiple agents articulate arguments culminating in a final answer through majority voting mechanisms (Liang et al., 2023), and Actor-Critic, where actor agents generate solutions subsequently evaluated by critic agents through iterative feedback processes (Shinn et al., 2023; Estornell et al., 2024; Yuan and Xie, 2025). However, existing communication paradigms collectively conceptualize individual agents as components with optimization objectives centered on MAS performance.

In contrast, our work treats each agent as an independent entity. We emulate human discussion dynamics through a novel Idea3 interaction, specifically designed to facilitate critical thinking communication among agents via its three-stage process: Idea Sharing, Idea Analysis, and Idea Fusion.

2.2 Multi-Agent Training

Conventional multi-agent training typically trains agents independently without awareness of other agents' states. Researchers employ multi-agent inference to collectively sample experiences, subsequently applying SFT or DPO to update individual agents independently, e.g., MALT (Motwani

et al., 2024), Multiagent-FT (Subramaniam et al., 2025), and DEBATUNE (Li et al., 2024a). However, this static one-time sampling fundamentally compromises the dynamic nature of multi-agent interactions. To address this limitation, recent advances in Multi-Agent Reinforcement Learning (MARL) have enabled continuous, real-time interaction sampling among agents (Ma et al., 2024; Chen et al., 2025; Liao et al., 2025). For example, MAPoRL (Park et al., 2025) implements a multi-agent proximal policy optimization algorithm, defining the agent state as the concatenation of interaction histories and incorporating manually predefined hyperparameters into rewards to incentivize collaboration.

Building on their design, we introduce a fully automated reward calibration, enabling automatic peer perception without manual intervention.

3 Methodology

As illustrated in Figure 2, for each question, each LLM first dynamically selects an appropriate interaction mode and then engages in Idea3 interactions (Idea Sharing, Idea Analysis, and Idea Fusion) to produce a group of answers. Finally, a reward calibration allows LLMs to perceive the quality of peers' solutions, fostering more effective multi-agent learning. By doing so, we aim to more accurately simulate and study the human-like behaviors in existing LLMs within a context that mirrors real-world human learning scenarios. We provide the pseudocode of ILR in Appendix E.

3.1 Question Difficulty Estimation

In real-world learning environments, such as classrooms, students often adapt their strategies to the complexity of the problem: for challenging tasks, they are more likely to collaborate, whereas for simpler ones, they tend to compete to demonstrate the efficiency and superiority of their solutions (Richard et al., 2002; Green and Rechis, 2006; Schneider et al., 2011; Fülöp, 2022). Both cooperation and competition can serve as drivers of multi-agent learning, which provides the underlying motivation for our Dynamic Interaction design. Since most problems lack explicit or continuous difficulty annotations, we adopt a self-ranking to estimate the question difficulty (Wang et al., 2025).

Given a training dataset of N questions, we divide it into B random batches to avoid the long-context lost problem (Liu et al., 2023), with each

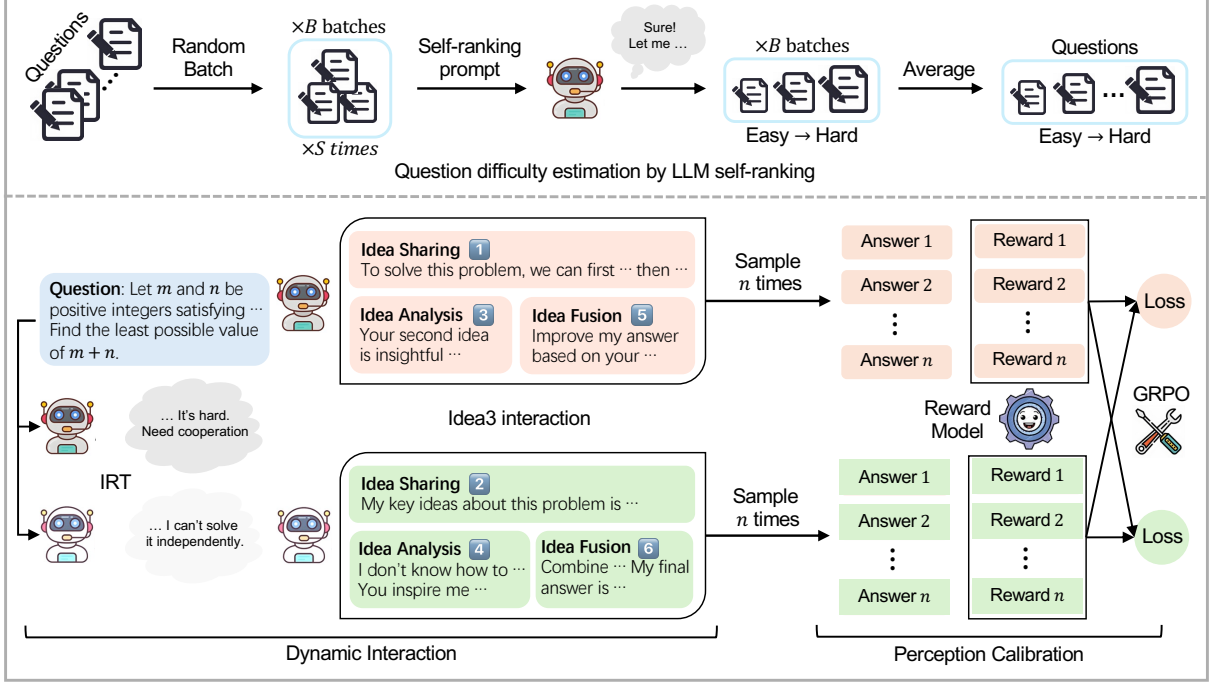


Figure 2: Illustration of proposed ILR. **Top:** LLMs first sort questions within each batch using the self-ranking prompt and then compute the average rank. **Bottom:** We depict the cooperation interaction here. If the question is too hard, LLMs will share their ideas, identify those complementary steps from other LLMs’ ideas, and improve their answers. The competition interaction also follows Idea3, which only requires a minor change in prompts.

batch containing $N' = \frac{N}{B}$ questions ($N' \ll N$). Using a self-ranking prompt (detailed in Appendix B.1), the LLM is instructed to order the questions within each batch in ascending difficulty. The ranks are then normalized into difficulty scores, where the easiest question is assigned $\frac{1}{N'}$ (rank 1) and the hardest one is assigned 1.0 (rank N'). Since a single random split provides only relative difficulty within its batch, we perform S independent splits and average their results to obtain a more stable and robust difficulty estimate for the entire dataset. We set N' and S to 10 in our work. For a set of LLMs $M = \{\mathcal{M}_i | i = 0, 1, \dots, m\}$ in a multi-agent learning scenario, the difficulty score D_q of a question q is computed as:

$$D_q = \frac{1}{m} \sum_{i=1}^m \mathcal{M}_i \left(\frac{1}{S} \sum_{j=1}^S \frac{r_{q,j}}{N'} \right) \quad (1)$$

where m is the number of LLMs, $\mathcal{M}_i(\cdot)$ denotes the estimation given by the i -th LLM, and $r_{q,j}$ is the rank of question q in the j -th random split.

3.2 Dynamic Interaction

When an LLM with reasoning ability γ_i receives a question of difficulty level D_q , it can quantify the probability ($P_{q,i}$) of correctly answering the

question using IRT (Benedetto et al., 2023):

$$P_{q,i} = \frac{1}{1 + e^{-1.7 \times (\gamma_i - D_q)}} \quad (2)$$

where the empirically derived coefficient 1.7 has been shown to yield reliable predictions across diverse conditions (Baker, 2001; De Ayala, 2013; Benedetto et al., 2023). IRT is a mature framework from educational measurement designed to simultaneously model the relationship between three key variables: the ability of an individual, the difficulty of a question, and the probability of the individual correctly solving the question. This is highly analogous to our scenario, and previous works have investigated the intersection of IRT and AI/LLM (Wang et al., 2023; Lalor et al., 2024; Zhuang et al., 2023), demonstrating its relevance and utility in our domain. Since $P_{q,i} = 0.5$ when $\gamma_i = D_q$, we adopt 0.5 as the decision boundary between different interaction modes. We average the probability of m LLMs to derive the overall probability (P_q) of independently solving question q and determine the interaction mode. We study the effect of mixed-type interaction for each question in Appendix D.5.

$$\text{Mode} = \begin{cases} \text{Cooperation} & \text{if } P_q < 0.5 \\ \text{Competition} & \text{if } P_q \geq 0.5 \end{cases} \quad (3)$$

To simulate human discussion, we design a novel and unified three-stage Idea3 interaction for multi-agent communication: Idea Sharing (each LLM proposes its own solution), Idea Analysis (each LLM analyzes and reflects on the peer’s solution), and Idea Fusion (the insights are synthesized into a refined and potentially novel solution). For different modes, we only need to slightly modify the prompt to inject the corresponding signal (see Appendix B.1). Unlike debate frameworks, which directly encourage one LLM to take another’s output as advice without further thinking, our Idea3 is designed to foster critical thought during communication. We show two cases in Appendix F.

Idea Sharing. Each LLM begins by presenting its problem-solving strategy, explaining the reasoning process and methods employed to address the given problem. For example, when solving a complex algebraic equation, one model might focus on factoring, while another may rely on graphical analysis. This stage produces the *initial answer*.

Idea Analysis. Subsequently, LLMs engage in a critical evaluation of each other’s proposed methods. In the cooperation mode, they may identify complementary strengths from different approaches, such as combining graphical insights with algebraic manipulation to generate a more comprehensive solution. In the competition mode, however, they rigorously assess the merits and limitations of the shared strategies. For example, one LLM might argue that the factoring approach, while effective, overlooks potential solutions that could be derived from the quadratic formula, thereby revealing a potential improvement.

Idea Fusion. Finally, LLMs synthesize the insights gained during previous analyses to generate a refined answer. This may involve integrating the most effective elements of both approaches, yielding a solution that not only accurately addresses the problem but also leverages complementary techniques from each LLM. For example, the final resolution to the algebraic equation might incorporate both the graphical representation for visual clarity and the algebraic methods for precision, culminating in a solution that is both robust and comprehensible. This stage’s output is the *updated answer*.

Prior research has noted that inter-agent communication may introduce noise into final outputs (Pan et al., 2025; Zhang et al., 2025). To address this, we adopt a label-based selection mechanism: the initial answer is retained only if it is correct and the updated answer is incorrect; in all other cases,

the updated answer is chosen. Finally, the rollout experience only contains this chosen solution, and the reward is calculated solely based on this “final answer”. This design choice can ensure: (1) To avoid a potential prompt shift. Since the rollout only contains a higher-quality step-by-step solution without intermediate results, the LLM is trained on the same type of data that it is expected to produce during inference. Therefore, there is no misalignment between the training and inference prompts. (2) To ensure a fair comparison with single-agent learning baselines. Including the rich information from the interaction trajectory in the reward signal could introduce a confounding variable, making it difficult to determine whether performance improvements stem from ILR or simply from the additional information in the trajectory.

3.3 Perception Calibration

Prior works rely on adding manually predefined signals to rewards (Park et al., 2025), which are discrete and coarse-grained, and in turn limit the scalability of reward shaping. In contrast, we introduce a fully automated method that integrates the distributional characteristics of one LLM’s reward data into another LLM’s reward function.

For a given input question, m LLMs each perform n sampling rounds, producing m groups of responses, with n answers in each group. A reward model is first used to assign initial rewards R to all responses. Each group is then summarized by its maximum (R_{max}), minimum (R_{min}), and average (R_{avg}) scores, collectively reflecting the model’s overall answer quality for that question. These statistics are arithmetically normalized and injected into the reward shaping process of peer models, yielding the final reward \bar{R} . For example, the k -th final reward of LLM i is computed as:

$$\bar{R}_{i,k} = R_{i,k} + \sum_{l \in M \setminus \{M_i\}} \text{clip}(\delta_{i,l}, -U, U) \quad (4)$$

where $\delta_{i,l} = \frac{R_{i,k} - R_{l,avg}}{R_{l,max} - R_{l,min}}$, $U = \frac{1}{m-1}$, $\text{clip}(\cdot)$ is a stabilization operation to prevent extreme values. We further discuss the suitability of $\text{clip}(\cdot)$ in Appendix B.4 and study the effect of z-score variant of $\delta_{i,l}$ in Appendix D.1. We then apply standard GRPO to optimize LLMs using calibrated rewards.

4 Experimental Setting

Test Models. We conduct experiments on three representative LLMs spanning two series and

two scales: Llama-3.1-8B-Instruct (Dubey et al., 2024), Qwen2.5-7B-Instruct, and Qwen2.5-14B-Instruct (Team, 2024). We organize these models into three groups and apply ILR within each. We further study a homogeneous setting (same LLM with different input prompts) in Appendix D.7.

- Group1 (different **series**): Llama-3.1-8B-Instruct and Qwen2.5-7B-Instruct
- Group2 (different **series & scale**): Llama-3.1-8B-Instruct and Qwen2.5-14B-Instruct
- Group3 (distinct **scale**): Qwen2.5-7B-Instruct and Qwen2.5-14B-Instruct

Benchmarks. For ILR *training*, we use the MATH (Hendrycks et al., 2021) following (Zeng et al., 2025) and randomly select 1,000 samples as the validation set to estimate each LLM’s γ_i offline before training, while the other 11,000 samples are used for training. We compare this one-time measurement strategy (γ_i) with a periodic-update strategy in Appendix D.4. For ILR **evaluation**, we conduct a comprehensive assessment across multiple mathematical reasoning benchmarks, which encompass both standard benchmarks, including GSM8K (Cobbe et al., 2021), MATH-500 (Hendrycks et al., 2021), Minerva Math (Lewkowycz et al., 2022), and Olympiad Bench (He et al., 2024), and a competition-level benchmark AIME24&25 (MAA-Committees, 2025), consisting of AIME2024 and AIME 2025. To further assess generalization beyond math, we evaluate ILR and baselines on a code generation benchmark MBPP (Austin et al., 2021), a general QA benchmark NQ (Kwiatkowski et al., 2019), and a scientific reasoning benchmark MMLU (Hendrycks et al., 2020). Note that except for Table 1, which is averaged across 5 seeds, all the other results are based on seed 0.

Baselines. We compare our ILR with six single-agent learning (self-learning) baselines, including SFT (Achiam et al., 2023), DPO (Rafailov et al., 2023), PPO (Schulman et al., 2017), GRPO (Shao et al., 2024), GSPO (Zheng et al., 2025), and Reinforce++ (Hu, 2025), and three multi-agent learning baselines, including DebateFT (Subramaniam et al., 2025), COPY (Ma et al., 2024), and MA-PoRL (Park et al., 2025). These baselines represent widely adopted and empirically effective approaches. Details about baselines, hyperparameters, training, and evaluation are provided in Appendix B.2 and Appendix B.3. We further compare

ILR with more self-learning baselines like self-reflection learning and advanced GRPO variants in Appendix D.3.

5 Results and Analysis

This section presents our main findings on ILR. Additional discussions on role assignment and application scenarios are provided in Appendix C. Further analyses, including z-score reward normalization, periodic update of γ_i , mixed-type interactions, cost analysis, homogeneous LLM settings, and another RL algorithm, are reported in Appendix D.

5.1 LLMs Perform Better Through ILR

Overall Comparison. Table 1 presents the quantitative comparison on mathematical benchmarks and Table 2 depicts the out-of-domain performance on MBPP, NQ, and MMLU. The results show that ILR’s multi-agent learning framework consistently enhances the independent problem-solving capability of individual LLMs. From Table 1, we highlight two further insights: (1) **ILR Promotes Complex Reasoning.** On the competition-level dataset AIME24&25 (comprising AIME24 and AIME25), ILR significantly improves LLMs’ capability in solving complex problems, which demonstrates that both weaker and stronger models can benefit from ILR’s multi-agent learning framework, enabling them to independently tackle challenging reasoning tasks. (2) **Balanced Grouping Improves Learning.** Each LLM is trained with ILR under two different grouping settings. Across comparisons, models achieve stronger results when paired with peers of more similar initial reasoning ability. For example, Llama-3.1-8B-Instruct performs better in ILR-Group1 with Qwen2.5-7B-Instruct (40.88% vs. 40.44% in Group2). We attribute this phenomenon to the fact that excessive initial performance disparities may lead to imbalanced interactions where the stronger LLM overwhelmingly dominates the entire process. While this pattern consistently emerges across three groups, the observed performance differences are modest. We leave more comprehensive empirical validation of this finding to future research.

Ablation Study. To further analyze contributions of each component, we conduct an ablation study and report the average accuracy in Table 3. Removing either component leads to consistent performance drops across all models, underscoring their joint contribution to the overall performance.

	GSM8K	MATH-500	Minerva Math	Olympiad Bench	AIME 24&25	Avg
Llama-3.1-8B-Instruct	82.38	48.60	23.31	13.30	2.33	33.99
SFT	85.17	49.68	23.90	16.68	3.00	35.69
DPO	86.06	50.12	24.63	17.36	2.67	36.17
PPO	86.17	52.44	26.03	19.26	4.00	37.58
GRPO	85.70	52.32	26.69	19.79	4.00	37.70
GSPO	85.82	53.76	27.43	18.90	5.00	38.18
Reinforce++	86.88	51.40	29.78	18.52	3.33	37.98
DebateFT-Group1	84.91	49.60	25.14	17.93	2.67	36.05
DebateFT-Group2	84.82	49.24	24.19	16.89	2.67	35.56
MAPoRL-Group1	86.81	53.00	30.88	21.93	6.67	39.86
MAPoRL-Group2	87.41	53.60	31.25	21.48	5.00	39.75
COPY-Group1	86.88	54.40	28.68	19.70	5.00	38.93
COPY-Group2	86.58	53.00	29.41	19.11	6.67	38.95
ILR-Group1	88.40	54.56	30.37	22.07	9.00	40.88
ILR-Group2	87.38	54.36	32.28	21.51	6.67	40.44
Qwen2.5-7B-Instruct	91.83	74.84	40.95	36.83	9.33	50.76
SFT	91.90	75.92	41.25	37.24	11.33	51.53
DPO	92.11	75.72	42.43	37.54	10.33	51.63
PPO	92.12	76.44	43.31	38.07	11.67	52.32
GRPO	92.48	75.88	41.91	38.19	13.00	52.29
GSPO	92.33	76.80	42.06	38.34	13.66	52.64
Reinforce++	92.57	76.56	42.13	38.22	12.67	52.43
DebateFT-Group1	92.16	75.96	41.25	37.28	10.33	51.40
DebateFT-Group3	92.17	76.44	41.11	37.36	12.33	51.88
MAPoRL-Group1	92.03	76.60	44.12	38.67	15.00	53.28
MAPoRL-Group3	92.10	77.20	42.65	39.41	13.33	52.94
COPY-Group1	92.57	77.00	42.65	39.11	16.33	53.53
COPY-Group3	92.27	77.20	43.01	38.52	15.00	53.20
ILR-Group1	93.00	77.00	44.56	39.50	16.33	54.08
ILR-Group3	92.88	77.80	44.73	39.17	17.00	54.31
Qwen2.5-14B-Instruct	94.83	80.08	46.03	40.21	12.33	54.70
SFT	94.98	80.20	46.32	40.86	13.00	55.07
DPO	95.09	80.84	47.28	40.59	14.33	55.63
PPO	95.06	80.88	48.45	41.10	16.67	56.43
GRPO	94.86	80.48	47.43	42.11	16.00	56.17
GSPO	95.30	80.12	48.82	42.96	16.34	56.71
Reinforce++	95.03	80.68	48.31	41.39	16.67	56.41
DebateFT-Group2	94.85	80.56	46.62	41.15	13.67	55.37
DebateFT-Group3	94.92	80.76	47.57	41.48	12.00	55.35
MAPoRL-Group2	95.22	80.80	49.26	42.41	20.00	57.54
MAPoRL-Group3	95.45	81.40	50.00	42.67	18.33	57.57
COPY-Group2	95.08	81.20	50.00	42.81	16.67	57.15
COPY-Group3	95.28	81.80	48.90	43.11	18.33	57.48
ILR-Group2	95.50	81.56	49.85	43.11	20.00	58.00
ILR-Group3	95.42	82.32	50.07	43.53	19.66	58.20

Table 1: The quantification comparison (accuracy %) of ILR and other baselines.

5.2 Dynamic Interaction Enhances Stronger LLMs’ Robustness

In Section 5.1, we demonstrate how multi-agent learning through ILR training strengthens the independent reasoning abilities of individual LLMs. Here, we evaluate the effectiveness of Idea3 communication during the inference stage. In pure inference scenarios where ground-truth labels are unavailable, we employ a summarization prompt to synthesize the initial and updated responses, thereby mitigating noise from multi-agent interactions. Specifically, for a given input question,

two LLMs first engage in Idea3 communication, after which each model evaluates both its own initial answer and the updated answer to produce a final prediction. For comparison, we also include Debate, which treats other agents’ outputs as additional advice to inform final answer generation. Full prompt details are provided in Appendix B.1.

Table 4 shows the multi-agent inference result on MATH-500 using untrained LLMs. We consistently observe that within different groups, Debate is more beneficial for weaker LLMs, while our Idea3 enhances the robustness of stronger LLMs

Benchmark	Model	Base	DPO	GRPO	ILR-G1	ILR-G2	ILR-G3
MBPP	Llama-3.1-8B-Instruct	54.00	56.40	56.80	57.40	57.60	-
	Qwen2.5-7B-Instruct	64.80	65.20	65.20	65.60	-	66.20
	Qwen2.5-14B-Instruct	71.40	71.60	71.80	-	72.40	71.60
NQ	Llama-3.1-8B-Instruct	51.12	52.32	52.11	54.77	55.31	-
	Qwen2.5-7B-Instruct	40.23	41.38	41.80	43.23	-	43.60
	Qwen2.5-14B-Instruct	49.13	49.60	50.00	-	52.45	53.87
MMLU	Llama-3.1-8B-Instruct	67.11	67.55	68.77	70.65	70.12	-
	Qwen2.5-7B-Instruct	52.15	52.54	52.56	54.00	-	54.14
	Qwen2.5-14B-Instruct	78.33	78.45	78.60	-	79.66	80.22

Table 2: Out-of-domain evaluation of ILR, DPO, and GRPO on MBPP (Pass@1), NQ (Accuracy), and MMLU (Accuracy). G_i means Group i . Compared with representative baselines, ILR further improves the performance.

	Llama-3.1-8B-Instruct		Qwen2.5-7B-Instruct		Qwen2.5-14B-Instruct	
	Group1	Group2	Group1	Group3	Group2	Group3
ILR	41.51	41.10	54.44	54.59	58.95	59.30
DI-only	39.12	39.25	53.95	53.23	58.04	58.57
PC-only	40.14	38.41	54.04	53.91	57.66	58.07

Table 3: Ablation Study of ILR. We report the average accuracy (%) of five mathematical evaluation benchmarks. DI, PC represent Dynamic Interaction and Perception Calibration.

by making them less susceptible to low-quality responses generated from weaker LLMs during multi-agent communication. We attribute this phenomenon to two primary reasons: **First**, for weaker LLMs, Debate directly incorporates stronger models’ answers as additional guidance, enabling them to refine their outputs, which often results in a more significant improvement. **Second**, for stronger LLMs, Debate similarly compels them to consider answers from weaker LLMs, which are usually lower in quality and can potentially degrade performance. In contrast, our Idea3 prompts stronger LLMs to critically evaluate and selectively integrate peer contributions, filtering out noise and thereby improving robustness. We further conduct multi-agent inference using ILR-trained LLMs as base models. Results show that ILR-trained LLMs achieve better performance than untrained LLMs in the same multi-agent inference setup, which indicates that ILR does not diminish an LLM’s ability to collaborate effectively at deployment.

5.3 Cooperation vs. Competition

In ILR training, we employ Item Response Theory (IRT) to dynamically determine interaction types, i.e., cooperation or competition. To further investigate the influence of cooperation, we vary the

cooperation ratio (p) from 0.0 to 1.0 in increments of 0.2. Here, $p = 0.0$ corresponds to full competition, $p = 1.0$ to full cooperation, and intermediate values designate the first p -proportion of questions (ranked by difficulty) as cooperative, with the remainder treated competitively. Due to training costs, we restrict this study to Group1.

Figure 3 shows the results, with IRT highlighted in red. Two key findings emerge for dynamic interaction design: (1) **Suboptimality of Extreme Strategies**. Relying solely on competition or cooperation is suboptimal for ILR, underscoring the necessity of adaptive interaction in multi-agent learning. This is intuitive: for challenging problems, cooperation allows LLMs to leverage complementary strengths and produce more comprehensive solutions, whereas for simpler tasks that can be effectively solved independently, excessive cooperation provides little benefit and may even introduce noise into the final outputs. (2) **Configuration of p** . The optimal cooperation ratio p requires careful design. One option is to manually partition data into subsets and tune p , but this is costly. IRT offers a practical alternative by approximating problem difficulty and aligning it with model reasoning capability. Although not always optimal (e.g., for

Inference Paradigm	Llama-3.1-8B-Instruct		Qwen2.5-7B-Instruct		Qwen2.5-14B-Instruct	
	Group1	Group2	Group1*	Group3	Group2*	Group3*
Single	49.80	49.80	75.60	75.60	81.20	81.20
Debate	64.00	66.20	74.60	80.00	79.20	81.00
Idea3 (Ours)	63.40	62.00	75.60	77.80	79.80	82.00
ILR-Single	55.80	55.20	77.60	78.00	81.80	82.60
ILR-Debate	65.60	71.20	81.20	81.00	79.40	81.20
ILR-Idea3 (Ours)	65.00	67.80	76.00	78.80	80.40	82.40

Table 4: Multi-agent inference results on MATH-500. ‘Single’ denotes the single-agent inference performance of the base models, and * indicates the stronger LLM within each group. For Debate and Idea3, the better-performing result is highlighted in bold. ILR- means we utilize LLMs after ILR-training as base models.

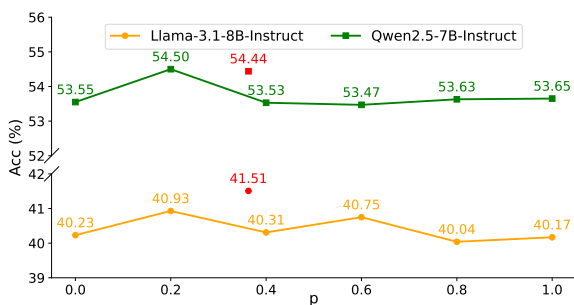


Figure 3: Average accuracy of Group1 under varying cooperation ratios. IRT is marked in red.

Qwen2.5-7B-Instruct), it achieves competitive results while eliminating manual intervention. This demonstrates the feasibility of IRT as a principled mechanism for integrating problem difficulty with LLM reasoning abilities. Future work may enhance robustness by incorporating additional conditional parameters into the IRT formulation.

5.4 Scalability of ILR

On balance of training costs and effectiveness, we primarily conduct experiments on three pairwise groups. To further evaluate the scalability of ILR beyond two-model interactions, we introduce a tri-model configuration, designated as Group 4 (ILR-G4). This experiment group integrates Llama-3.1-8B-Instruct, Qwen2.5-7B-Instruct, and Qwen2.5-14B-Instruct into a unified collaborative framework. The primary objective is to investigate the scalability of ILR with more than two LLMs.

As evidenced in Table 5, the deployment of three LLMs results in a consistent and monotonic performance improvement across all constituent models compared to their dyadic counterparts (ILR-G1, G2, and G3). All LLMs surpass their best two-model performance by nearly 1.0%, which suggests that ILR effectively leverages the heterogeneity of

Model	ILR-G1	ILR-G2	ILR-G3	ILR-G4
Llama	41.51	41.10	-	42.48
Qwen7B	54.44	-	54.59	55.12
Qwen14B	-	58.95	59.30	59.95

Table 5: Scaling trend of ILR (two vs. three LLMs per group). G_i denotes Group i . We report the average accuracy on five mathematical evaluation benchmarks.

larger model groups. The addition of a third model does not introduce noise or redundancy; rather, it provides complementary reasoning paths that further refine the consensus, confirming the robustness and positive scaling trend of our ILR. This presents an intriguing finding. Future work could focus on mitigating the training costs associated with multi-agent reinforcement learning, thereby enabling the investigation of interactive learning behaviors across broader ensembles of LLMs.

6 Conclusion

In this paper, we investigate whether interactive learning among multiple LLMs (*multi-agent learning*) can outperform traditional self-learning (*single-agent learning*). Inspired by real-world human interaction, we propose ILR, a novel framework built on two key components: Dynamic Interaction and Perception Calibration. Extensive experiments across different model series and scales demonstrate the effectiveness of ILR, showing that interactive learning consistently yields greater performance improvements than self-learning. We further investigate the advantage of Dynamic Interaction and the Scalability of ILR beyond two-model interactions. These findings align with human learning patterns and provide insights into analyzing high-level, human-like behaviors in LLMs.

Limitations

Our current reliance on intrinsic difficulty may be less designable for creative generation tasks, as even simple questions in this domain benefit significantly from the stylistic diversity of cooperation strategy. Therefore, additional conditioning variables are required to extend the robustness of Dynamic Interaction and the scenario of ILR.

Acknowledgments

This research/project is supported by the Youth S&T Talent Support Programme of GDSTA (SKXRC2025462).

References

- Josh Achiam, Steven Adler, Sandhini Agarwal, Lama Ahmad, Ilge Akkaya, Florencia Leoni Aleman, Diogo Almeida, Janko Altenschmidt, Sam Altman, Shyamal Anadkat, and 1 others. 2023. Gpt-4 technical report. *arXiv preprint arXiv:2303.08774*.
- Jacob Austin, Augustus Odena, Maxwell Nye, Maarten Bosma, Henryk Michalewski, David Dohan, Ellen Jiang, Carrie Cai, Michael Terry, Quoc Le, and 1 others. 2021. Program synthesis with large language models. *arXiv preprint arXiv:2108.07732*.
- Frank B Baker. 2001. *The basics of item response theory*. ERIC.
- Luca Benedetto, Paolo Cremonesi, Andrew Caines, Paula Buttery, Andrea Cappelli, Andrea Giussani, and Roberto Turrin. 2023. A survey on recent approaches to question difficulty estimation from text. *ACM Computing Surveys*, 55(9):1–37.
- Daan Bloembergen, Karl Tuyls, Daniel Hennes, and Michael Kaisers. 2015. Evolutionary dynamics of multi-agent learning: A survey. *Journal of Artificial Intelligence Research*, 53:659–697.
- Lucian Busoniu, Robert Babuska, and Bart De Schutter. 2006. Multi-agent reinforcement learning: A survey. In *2006 9th international conference on control, automation, robotics and vision*, pages 1–6. IEEE.
- Li Cai, Kilchan Choi, Mark Hansen, and Lauren Harrell. 2016. Item response theory. *Annual Review of Statistics and Its Application*, 3(1):297–321.
- Lorenzo Canese, Gian Carlo Cardarilli, Luca Di Nunzio, Rocco Fazzolari, Daniele Giardino, Marco Re, and Sergio Spanò. 2021. Multi-agent reinforcement learning: A review of challenges and applications. *Applied Sciences*, 11(11):4948.
- Chi-Min Chan, Weize Chen, Yusheng Su, Jianxuan Yu, Wei Xue, Shanghang Zhang, Jie Fu, and Zhiyuan Liu. 2023. Chateval: Towards better llm-based evaluators through multi-agent debate. *arXiv preprint arXiv:2308.07201*.
- Yiqun Chen, Lingyong Yan, Weiwei Sun, Xinyu Ma, Yi Zhang, Shuaiqiang Wang, Dawei Yin, Yiming Yang, and Jiaxin Mao. 2025. Improving retrieval-augmented generation through multi-agent reinforcement learning. *arXiv preprint arXiv:2501.15228*.
- Karl Cobbe, Vineet Kosaraju, Mohammad Bavarian, Mark Chen, Heewoo Jun, Lukasz Kaiser, Matthias Plappert, Jerry Tworek, Jacob Hilton, Reiichiro Nakano, and 1 others. 2021. Training verifiers to solve math word problems. *arXiv preprint arXiv:2110.14168*.
- Abhishek Das, Théophile Gervet, Joshua Romoff, Dhruv Batra, Devi Parikh, Mike Rabbat, and Joelle Pineau. 2019. Tarmac: Targeted multi-agent communication. In *International Conference on machine learning*, pages 1538–1546. PMLR.
- Rafael Jaime De Ayala. 2013. *The theory and practice of item response theory*. Guilford Publications.
- Shifei Ding, Wei Du, Ling Ding, Lili Guo, and Jian Zhang. 2024. Learning efficient and robust multi-agent communication via graph information bottleneck. In *Proceedings of the AAAI Conference on Artificial Intelligence*, volume 38, pages 17346–17353.
- Zi-Yi Dou, Cheng-Fu Yang, Xueqing Wu, Kai-Wei Chang, and Nanyun Peng. 2024. Re-rest: Reflection-reinforced self-training for language agents. *arXiv preprint arXiv:2406.01495*.
- Abhimanyu Dubey, Abhinav Jauhri, Abhinav Pandey, Abhishek Kadian, Ahmad Al-Dahle, Aiesha Letman, Akhil Mathur, Alan Schelten, Amy Yang, Angela Fan, and 1 others. 2024. The llama 3 herd of models. *arXiv e-prints*, pages arXiv–2407.
- Andrew Estornell, Jean-François Ton, Yuanshun Yao, and Yang Liu. 2024. Acc-collab: An actor-critic approach to multi-agent llm collaboration. *arXiv preprint arXiv:2411.00053*.
- Márta Fülöp. 2022. Cooperation and competition. *The Wiley-Blackwell handbook of childhood social development*, pages 555–572.
- Vanessa A Green and Ruth Rechis. 2006. Children’s cooperative and competitive interactions in limited resource situations: A literature review. *Journal of applied developmental psychology*, 27(1):42–59.
- Sven Gronauer and Klaus Diepold. 2022. Multi-agent deep reinforcement learning: a survey. *Artificial Intelligence Review*, 55(2):895–943.
- Daya Guo, Dejian Yang, Haowei Zhang, Junxiao Song, Ruoyu Zhang, Runxin Xu, Qihao Zhu, Shiron Ma, Peiyi Wang, Xiao Bi, and 1 others. 2025. Deepseek-r1: Incentivizing reasoning capability in llms via reinforcement learning. *arXiv preprint arXiv:2501.12948*.

- Shanshan Han, Qifan Zhang, Yuhang Yao, Weizhao Jin, and Zhaozhuo Xu. 2024. Llm multi-agent systems: Challenges and open problems. *arXiv preprint arXiv:2402.03578*.
- Chaoqun He, Renjie Luo, Yuzhuo Bai, Shengding Hu, Zhen Leng Thai, Junhao Shen, Jinyi Hu, Xu Han, Yujie Huang, Yuxiang Zhang, and 1 others. 2024. Olympiadbench: A challenging benchmark for promoting agi with olympiad-level bilingual multimodal scientific problems. *arXiv preprint arXiv:2402.14008*.
- Jiayi He, Hehai Lin, Qingyun Wang, Yi R Fung, and Heng Ji. 2025. Self-correction is more than refinement: A learning framework for visual and language reasoning tasks. In *Findings of the Association for Computational Linguistics: ACL 2025*, pages 6405–6421.
- Dan Hendrycks, Collin Burns, Steven Basart, Andy Zou, Mantas Mazeika, Dawn Song, and Jacob Steinhardt. 2020. Measuring massive multitask language understanding. *arXiv preprint arXiv:2009.03300*.
- Dan Hendrycks, Collin Burns, Saurav Kadavath, Akul Arora, Steven Basart, Eric Tang, Dawn Song, and Jacob Steinhardt. 2021. Measuring mathematical problem solving with the math dataset. *arXiv preprint arXiv:2103.03874*.
- Chin-min Hsiung. 2012. The effectiveness of cooperative learning. *Journal of engineering Education*, 101(1):119–137.
- Jian Hu. 2025. Reinforce++: A simple and efficient approach for aligning large language models. *arXiv preprint arXiv:2501.03262*.
- Jian Hu, Xibin Wu, Zilin Zhu, Xianyu, Weixun Wang, Dehao Zhang, and Yu Cao. 2024. Openrlhf: An easy-to-use, scalable and high-performance rlhf framework. *arXiv preprint arXiv:2405.11143*.
- Murat Kahveci and Yesim Imamoglu. 2007. Interactive learning in mathematics education: Review of recent literature. *Journal of Computers in Mathematics and Science Teaching*, 26(2):137–153.
- Tom Kwiatkowski, Jennimaria Palomaki, Olivia Redfield, Michael Collins, Ankur Parikh, Chris Alberti, Danielle Epstein, Illia Polosukhin, Jacob Devlin, Kenton Lee, and 1 others. 2019. Natural questions: a benchmark for question answering research. *Transactions of the Association for Computational Linguistics*, 7:453–466.
- John P Lalor, Pedro Rodriguez, João Sedoc, and Jose Hernandez-Orallo. 2024. Item response theory for natural language processing. In *Proceedings of the 18th Conference of the European Chapter of the Association for Computational Linguistics: Tutorial Abstracts*, pages 9–13.
- Aitor Lewkowycz, Anders Andreassen, David Dohan, Ethan Dyer, Henryk Michalewski, Vinay Ramasesh, Ambrose Slone, Cem Anil, Imanol Schlag, Theo Gutman-Solo, and 1 others. 2022. Solving quantitative reasoning problems with language models. *Advances in neural information processing systems*, 35:3843–3857.
- Ming Li, Jiuhai Chen, Lichang Chen, and Tianyi Zhou. 2024a. Can llms speak for diverse people? tuning llms via debate to generate controllable controversial statements. *arXiv preprint arXiv:2402.10614*.
- Xinyi Li, Sai Wang, Siqi Zeng, Yu Wu, and Yi Yang. 2024b. A survey on llm-based multi-agent systems: workflow, infrastructure, and challenges. *Vicinatearth*, 1(1):9.
- Tian Liang, Zhiwei He, Wenxiang Jiao, Xing Wang, Yan Wang, Rui Wang, Yujiu Yang, Shuming Shi, and Zhaopeng Tu. 2023. Encouraging divergent thinking in large language models through multi-agent debate. *arXiv preprint arXiv:2305.19118*.
- Junwei Liao, Muning Wen, Jun Wang, and Weinan Zhang. 2025. Marft: Multi-agent reinforcement fine-tuning. *arXiv preprint arXiv:2504.16129*.
- Zhihang Lin, Mingbao Lin, Yuan Xie, and Rongrong Ji. 2025. Cppo: Accelerating the training of group relative policy optimization-based reasoning models. *arXiv preprint arXiv:2503.22342*.
- Nelson F Liu, Kevin Lin, John Hewitt, Ashwin Paranjape, Michele Bevilacqua, Fabio Petroni, and Percy Liang. 2023. Lost in the middle: How language models use long contexts. *arXiv preprint arXiv:2307.03172*.
- Hao Ma, Tianyi Hu, Zhiqiang Pu, Liu Boyin, Xiaolin Ai, Yanyan Liang, and Min Chen. 2024. Coevolving with the other you: Fine-tuning llm with sequential cooperative multi-agent reinforcement learning. *Advances in Neural Information Processing Systems*, 37:15497–15525.
- MAA-Committees. 2025. [Aime problems and solutions](#).
- Stephan Mende, Antje Proske, and Susanne Narciss. 2021. Individual preparation for collaborative learning: Systematic review and synthesis. *Educational Psychologist*, 56(1):29–53.
- Sumeet Ramesh Motwani, Chandler Smith, Rocktim Jyoti Das, Rafael Rafailov, Ivan Laptev, Philip HS Torr, Fabio Pizzati, Ronald Clark, and Christian Schroeder de Witt. 2024. Malt: Improving reasoning with multi-agent llm training. *arXiv preprint arXiv:2412.01928*.
- Ansong Ni, Srini Iyer, Dragomir Radev, Ves Stoyanov, Wen-tau Yih, Sida I Wang, and Xi Victoria Lin. 2023. Lever: Learning to verify language-to-code generation with execution. In *Proceedings of the 40th International Conference on Machine Learning (ICML'23)*.

- Melissa Z Pan, Mert Cemri, Lakshya A Agrawal, Shuyi Yang, Bhavya Chopra, Rishabh Tiwari, Kurt Keutzer, Aditya Parameswaran, Kannan Ramchandran, Dan Klein, and 1 others. 2025. Why do multiagent systems fail? In *ICLR 2025 Workshop on Building Trust in Language Models and Applications*.
- Chanwoo Park, Seungju Han, Xingzhi Guo, Asuman Ozdaglar, Kaiqing Zhang, and Joo-Kyung Kim. 2025. Maporl: Multi-agent post-co-training for collaborative large language models with reinforcement learning. *arXiv preprint arXiv:2502.18439*.
- Rafael Rafailov, Archit Sharma, Eric Mitchell, Christopher D Manning, Stefano Ermon, and Chelsea Finn. 2023. Direct preference optimization: Your language model is secretly a reward model. *Advances in neural information processing systems*, 36:53728–53741.
- Matthew Renze and Erhan Guven. 2024. Self-reflection in llm agents: Effects on problem-solving performance. *arXiv preprint arXiv:2405.06682*.
- Jacques F Richard, Ada Fonzi, Franca Tani, Fulvio Tassi, Giovanna Tomada, Barry H Schneider, and 1 others. 2002. Cooperation and competition. *Blackwell handbook of childhood social development*, pages 515–532.
- Barry H Schneider, Joyce Benenson, Márta Fülöp, Mihály Berkics, and Mónika Sándor. 2011. Cooperation and competition. *The Wiley-Blackwell handbook of childhood social development*, pages 472–490.
- John Schulman, Filip Wolski, Prafulla Dhariwal, Alec Radford, and Oleg Klimov. 2017. Proximal policy optimization algorithms. *arXiv preprint arXiv:1707.06347*.
- Zhihong Shao, Peiyi Wang, Qihao Zhu, Runxin Xu, Junxiao Song, Xiao Bi, Haowei Zhang, Mingchuan Zhang, YK Li, Yang Wu, and 1 others. 2024. Deepseekmath: Pushing the limits of mathematical reasoning in open language models. *arXiv preprint arXiv:2402.03300*.
- Noah Shinn, Federico Cassano, Ashwin Gopinath, Karthik Narasimhan, and Shunyu Yao. 2023. Reflexion: Language agents with verbal reinforcement learning. *Advances in Neural Information Processing Systems*, 36:8634–8652.
- Vighnesh Subramaniam, Yilun Du, Joshua B Tenenbaum, Antonio Torralba, Shuang Li, and Igor Mordatch. 2025. Multiagent finetuning: Self improvement with diverse reasoning chains. *arXiv preprint arXiv:2501.05707*.
- Chuxiong Sun, Zehua Zang, Jiabao Li, Jiangmeng Li, Xiao Xu, Rui Wang, and Changwen Zheng. 2024. T2mac: Targeted and trusted multi-agent communication through selective engagement and evidence-driven integration. In *Proceedings of the AAAI Conference on Artificial Intelligence*, volume 38, pages 15154–15163.
- Zexu Sun, Yongcheng Zeng, Erxue Min, Heyang Gao, Bokai Ji, and Xu Chen. 2025. Cog-rethinker: Hierarchical metacognitive reinforcement learning for llm reasoning. *arXiv preprint arXiv:2510.15979*.
- Qwen Team. 2024. Qwen2 technical report. *arXiv preprint arXiv:2407.10671*.
- Ziyu Wan, Yunxiang Li, Xiaoyu Wen, Yan Song, Hanjing Wang, Linyi Yang, Mark Schmidt, Jun Wang, Weinan Zhang, Shuyue Hu, and 1 others. 2025. Rema: Learning to meta-think for llms with multi-agent reinforcement learning. *arXiv preprint arXiv:2503.09501*.
- Xinglin Wang, Shaoxiong Feng, Yiwei Li, Peiwen Yuan, Yueqi Zhang, Chuyi Tan, Boyuan Pan, Yao Hu, and Kan Li. 2025. Make every penny count: Difficulty-adaptive self-consistency for cost-efficient reasoning. In *Findings of the Association for Computational Linguistics: NAACL 2025*, pages 6904–6917.
- Xiting Wang, Liming Jiang, Jose Hernandez-Orallo, David Stillwell, Luning Sun, Fang Luo, and Xing Xie. 2023. Evaluating general-purpose ai with psychometrics. *arXiv preprint arXiv:2310.16379*.
- An Yang, Baosong Yang, Binyuan Hui, Bo Zheng, Bowen Yu, Chang Zhou, Chengpeng Li, Chengyuan Li, Dayiheng Liu, Fei Huang, and 1 others. 2024. Qwen2 technical report. *arXiv preprint arXiv:2407.10671*.
- Yong Yu, Xiaosheng Si, Changhua Hu, and Jianxun Zhang. 2019. A review of recurrent neural networks: Lstm cells and network architectures. *Neural computation*, 31(7):1235–1270.
- Yurun Yuan and Tengyang Xie. 2025. Reinforce llm reasoning through multi-agent reflection. *arXiv preprint arXiv:2506.08379*.
- Jimmy Zambrano, Femke Kirschner, John Sweller, and Paul A Kirschner. 2019. Effects of prior knowledge on collaborative and individual learning. *Learning and Instruction*, 63:101214.
- Weihao Zeng, Yuzhen Huang, Qian Liu, Wei Liu, Keqing He, Zejun Ma, and Junxian He. 2025. Simplerl-zoo: Investigating and taming zero reinforcement learning for open base models in the wild. *arXiv preprint arXiv:2503.18892*.
- Shaokun Zhang, Ming Yin, Jieyu Zhang, Jiale Liu, Zhiguang Han, Jingyang Zhang, Beibin Li, Chi Wang, Huazheng Wang, Yiran Chen, and 1 others. 2025. Which agent causes task failures and when? on automated failure attribution of llm multi-agent systems. *arXiv preprint arXiv:2505.00212*.
- Wenqi Zhang, Yongliang Shen, Linjuan Wu, Qiuying Peng, Jun Wang, Yueting Zhuang, and Weiming Lu. 2024. Self-contrast: Better reflection through inconsistent solving perspectives. In *Proceedings of the 62nd Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*, pages 3602–3622.

Wanjia Zhao, Mert Yuksekgonul, Shirley Wu, and James Zou. 2025. Sirius: Self-improving multi-agent systems via bootstrapped reasoning. *arXiv preprint arXiv:2502.04780*.

Chujie Zheng, Shixuan Liu, Mingze Li, Xiong-Hui Chen, Bowen Yu, Chang Gao, Kai Dang, Yuqiong Liu, Rui Men, An Yang, and 1 others. 2025. Group sequence policy optimization. *arXiv preprint arXiv:2507.18071*.

Yan Zhuang, Qi Liu, Yuting Ning, Weizhe Huang, Rui Lv, Zhenya Huang, Guanhao Zhao, Zheng Zhang, Qingyang Mao, Shijin Wang, and 1 others. 2023. Efficiently measuring the cognitive ability of llms: An adaptive testing perspective.

A Motivation of ILR

The rationale of the proposed ILR is intricately designed to function as a cohesive pipeline that mirrors humans’ interactive learning. In real-world peer discussions, individuals sequentially experience three cognitive stages: first, upon encountering a problem, they instinctively assess its complexity (Difficulty Estimation). Second, based on this assessment, they determine their interaction strategy, competing to find the most elegant solution for simple tasks, or cooperating to overcome challenging ones (Dynamic Interaction with the Idea3 framework). Finally, after the discussion, individuals reflect on and are incentivized by the quality of their peers’ solutions (Perception Calibration).

B Experiment Details

B.1 Prompt Details

Self-ranking Prompt for Difficulty Estimation

Ranking Prompt1:

Your task is to rank the given questions from easy to hard based on their difficulty level. Questions to be evaluated: $\{Q_1, Q_2, \dots, Q_{N'}\}$.

Ranking Prompt2:

You will be given a batch of questions. Your task is to rank them from easy to hard based on their difficulty level. You should carefully horizontally compare the given questions in order to assign a suitable ranking place to each question. Please make sure you read and understand these instructions carefully. Please keep this document open while reviewing, and refer to it as needed. Questions to be evaluated: $\{Q_1, Q_2, \dots, Q_{N'}\}$.

Ranking Prompt3:

You need to analyze and rank questions $Q_1-Q_{N'}$ by their difficulty level based on these criteria: (1) Cognitive load required. If a problem involves multiple steps, it will have a higher cognitive load than a problem with a single step. (2) Knowledge depth needed. Problems involving the deduction of complex formulas require deeper professional knowledge. (3) Typical error rates. For example, problems that tend to overlook a certain prerequisite or calculation step will have a relatively high error rate. Please first score each question (1-10 scale) on

three dimensions above, then calculate the average score, and rank these questions by final scores. Questions to be evaluated: $\{Q_1, Q_2, \dots, Q_{N'}\}$.

Format Prompt:

After analyzing all the questions, please give all the ranking places (from easy to hard) in order, following the template “Ranking: $[Q_{\{\text{number of the easiest question}\}}, \dots, Q_{\{\text{number of the hardest question}\}}]$ ”.

Prompt for Idea3 communication

Idea Sharing (Cooperation&Competition):

Question: {Input question}

Please reason step by step, and your final answer should be in the form boxed{answer} given at the end of your response.

Idea Analysis (Cooperation):

Partner’s Contribution: {Ideas from other LLMs}

Collaboratively analyze the key steps in the partner’s contribution, identify those steps that can help you improve your answer, and serve as additional advice.

Idea Analysis (Competition):

Opponent’s Solution: {Ideas from other LLMs}

Critically analyze the opponent’s ideas, identify the weaknesses and strengths of his ideas.

Idea Fusion (Cooperation&Competition):

Based on the above analysis, give an updated answer to the Original Question: {Input question}. Please reason step by step, and your final answer should be in the form boxed{answer} given at the end of your response.

Prompt for Evaluation

Single-agent Evaluation:

Please reason step by step, and your final answer should be in the form boxed{answer} given at the end of your response.

Multi-agent Evaluation:

Debate:

Here are solutions from other agents:

One agent response: {other_agent_response}

Using each response as additional advice based on the correctness of each response. Can you give an updated bullet-by-bullet answer to {Input question}. Please reason step by step, and your final answer should be in the form boxed{answer} given at the end of your response.

Summarization (After Idea3 communication):

The original question is {Input question}. There are two solutions you provided:

Solution 1: {Initial answer}

Solution 2: {Updated answer}

Please answer the original question step-by-step based on these two solutions, and your final answer should be in the form boxed{answer} given at the end of your response.

To mitigate potential prompt bias, we utilize three different self-ranking prompts of varying levels of granularity for each question and average the rankings to obtain the final estimation.

Furthermore, we add a post-validation (see Table 6) to study the effectiveness of our difficulty estimation, which measures the correlation of self-ranking difficulty and empirical “true solve rate” derived from 20 samplings per question on the val-

validation set (1,000 samples). The resulting high correlation coefficient between the self-ranking scores and the actual success rates confirms that our estimation mechanism accurately and robustly reflects the true difficulty relative to the LLM’s capabilities. Specifically, the first three columns report the correlation between each individual llm’s self-ranking scores and its respective solve probabilities, while the final “Avg” column presents the correlation derived from the averaged self-ranking scores and averaged solve probabilities across all llms.

B.2 Implementation Details

Baselines: For multi-agent learning, direct comparisons are challenging because existing methods typically train specialized LLMs with complementary roles for problem-solving. To enable fair evaluation, we introduce a minor modification to Multiagent-FT (Subramaniam et al., 2025), sampling answers through Debate and optimizing each LLM using the original training algorithm. Each LLM will solve questions independently at inference time. We rename this baseline as DebateFT.

Training: We use full-tuning to optimize the LLMs for one epoch. We use a batch size of 256 and a learning rate of 1e-6 for Llama-3.1-8B-Instruct, 1e-6 for Qwen2.5-7B-Instruct, and 9e-7 for Qwen2.5-14B-Instruct. The temperature is 0.5 for all LLMs, and the KL coefficient is 0 for Llama-3.1-8B-Instruct, 5e-7 for Qwen2.5-7B-Instruct, and 0 for Qwen2.5-14B-Instruct. The maximum output token number of the sampled answer is 2K. As for other hyperparameters, we strictly use the original parameters of GRPO. We utilize Llama-3-8b-rm-mixture (Hu et al., 2024) as the reward model to rate sampled answers. The initial reasoning ability γ_i , which can be measured on the validation set, is 0.59, 0.75, and 0.78 for Llama-3.1-8B-Instruct, Qwen2.5-7B-Instruct, and Qwen2.5-14B-Instruct. According to the IRT in Section 3.2, the proportion of cooperation/competition across the three groups is as follows: 36.30/63.70%, 32.76/67.24%, and 16.87/83.13% for Group1, Group2, and Group3. The approximate training time is 8 hours for Group1, 12 hours for Group2, and 10 hours for Group3.

Evaluation: We set the temperature as 0 and the maximum output token number of evaluation is set to 8K for AIME, while 2K for other benchmarks. All training experiments are conducted on eight H100 GPUs, and evaluation experiments are conducted on one H100 GPU.

B.3 Statistics of Benchmarks

Training Dataset

Following (Zeng et al., 2025), we only use MATH (Hendrycks et al., 2021) as our training data source. Excluding the common MATH-500 (Hendrycks et al., 2021) as the evaluation set, there are 12000 samples in the remaining dataset. We randomly select 1000 samples as the validation set to assess the LLMs’ reasoning ability for Dynamic Interaction, while the remaining 11000 samples serve as the training set to fine-tune LLMs.

Each question of the MATH training set will have a continuous difficulty measured by LLMs’ self-ranking, and we depict the question difficulty distribution in Figure 4. As illustrated in Section 4, the initial reasoning ability γ_i , which can be measured on the validation set, is 0.59, 0.75, and 0.78 for Llama-3.1-8B-Instruct, Qwen2.5-7B-Instruct, and Qwen2.5-14B-Instruct. According to the IRT in Section 3.2, the proportion of cooperation/competition across the three groups is as follows: 36.30/63.70%, 32.76/67.24%, and 16.87/83.13% for Group 1, Group 2, and Group 3, respectively.

Evaluation Dataset

GSM8K (Cobbe et al., 2021): This dataset comprises 1319 single-step mathematical questions designed for elementary school students in English. As one of the most widely adopted benchmarks in the field, it plays a basic role in assessing the arithmetic reasoning capabilities of LLMs.

MATH-500 (Hendrycks et al., 2021): The dataset encompasses 500 intermediate-level mathematical problems systematically curated across core domains, including algebra, geometry, precalculus, probability, and number theory.

Minerva Math (Lewkowycz et al., 2022): This dataset contains 272 mathematical problems across core domains, including algebra, geometry, precalculus, probability, and number theory.

Olympiad Bench (He et al., 2024): This benchmark constitutes a bilingual multimodal evaluation framework comprising 8476 Olympiad-level problems curated from prestigious mathematics and physics competitions. We utilize the subset processed by (Yang et al., 2024) as our test set, which contains 675 English text-only questions.

AIME24&25 (MAA-Committees, 2025): This benchmark collection contains 60 questions and derives from the 2024 and 2025 editions of the American Invitational Mathematics Examination (AIME), comprising two distinct problem sets. Each set con-

Meta-Llama-3.1-8B-Instruct	Qwen2.5-7B-Instruct	Qwen2.5-14B-Instruct	Avg
0.7138	0.6905	0.7254	0.7522

Table 6: The Pearson correlation of self-ranking difficulty and empirical “true solve rate”.

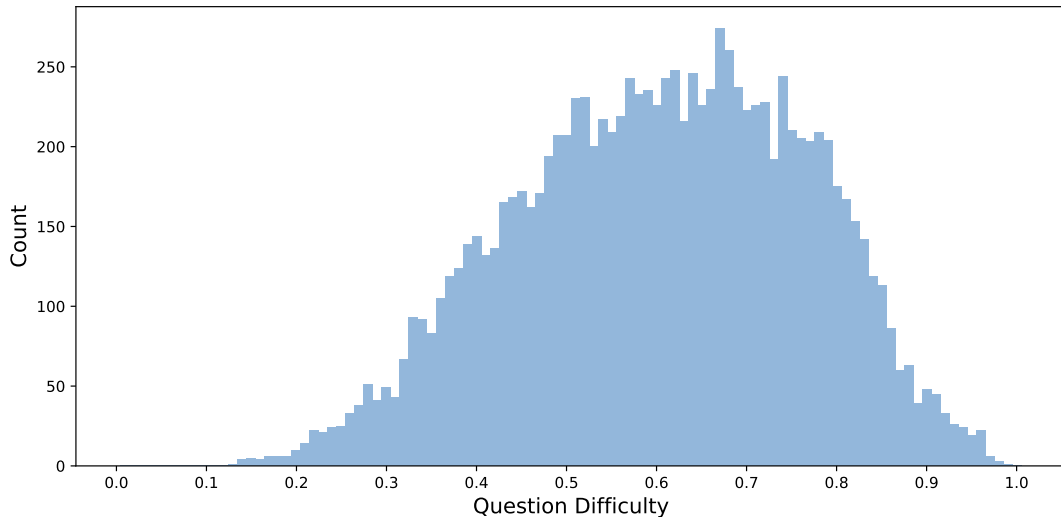


Figure 4: Question Difficulty Distribution (the interval is 0.01) of MATH training set used in ILR.

tains 30 rigorously vetted mathematical questions characterized by high cognitive demand. The evaluative focus lies in probing advanced mathematical competencies.

MBPP (Austin et al., 2021): This benchmark comprises 974 crowd-sourced Python programming problems, meticulously curated to align with the competency level of entry-level programmers. These problems encompass foundational programming concepts, standard library implementations, and essential algorithmic paradigms, ensuring a comprehensive evaluation of introductory programming proficiency. Following (Ni et al., 2023), we evaluate ILR and baselines using the test subset, which contains 500 questions.

NQ (Kwiatkowski et al., 2019): This benchmark is designed to evaluate open-domain question answering systems on natural, information-seeking questions that require reasoning over entire documents. We randomly sample 500 questions and evaluate ILR and baselines on this set.

MMLU (Hendrycks et al., 2020): This comprehensive benchmark spans subjects across STEM, the humanities, and the social sciences, designed to measure a model’s multitask accuracy and world knowledge. It includes questions ranging from elementary to professional levels, evaluating both knowledge acquisition and problem-solving abil-

ities in zero-shot and few-shot settings. We randomly sample 500 questions and evaluate ILR and baselines on this set.

B.4 Clip ratio

In Section 3.3, we utilize a $\text{clip}(\cdot)$ function to stabilize the ILR training. In Table 7, we report the proportion of rewards that are actually clipped. It shows that the percentage of instances where the normalized reward difference exceeded the boundary is small. Therefore, we believe that this clip will not overly suppress the informative signals.

C Discussion

C.1 Heterogeneous Role Assignments in ILR

Our ILR intentionally employs a homogeneous design where all agents are problem solvers. This choice is directly tied to our core research question: to investigate whether and how interaction within a multi-agent environment can enhance the independent problem-solving capabilities of each individual LLM. In this context, all agents share the same role and learn from each other through interaction to achieve mutual improvement. However, we believe the idea of incorporating heterogeneous roles is interesting, and an exciting future direction would be to design a hybrid system. For example, one could have multiple “generator”

Model	ILR-Group1	ILR-Group2	ILR-Group3
Llama-3.1-8B-Instruct	5.04%	7.15%	-
Qwen2.5-7B-Instruct	2.19%	-	1.64%
Qwen2.5-14B-Instruct	-	3.59%	2.34%

Table 7: The clip ratio of each group.

Model	Avg	z-Avg	ratio	z-ratio
Llama-3.1-8B-Instruct	41.51	40.74 (↓ 0.77)	5.04	20.45 (↑ 15.41)
Qwen2.5-7B-Instruct	54.44	53.94 (↓ 0.50)	2.19	9.10 (↑ 6.91)

Table 8: Comparison with z-score normalization (z-) based on Group1. We report the average accuracy (%) on five mathematical evaluation benchmarks and the clip ratio (%).

agents that are trained internally via ILR to enhance their respective independent generation capabilities. Concurrently, the system could include one or more “verifier” or “refiner” agents responsible for evaluating and integrating the outputs from these ILR-strengthened generators. Such a hybrid model could potentially lead to an even more complex, robust, and powerful multi-agent system.

C.2 Robustness of Dynamic Interaction

For many common reasoning tasks, such as mathematics and standard question-answering, the intrinsic difficulty of a problem is one of the most direct and dominant factors influencing whether an individual chooses to persevere independently or seek collaborative assistance. However, in scenarios like creative writing, the “problem difficulty” does not necessarily correlate directly with the quality of the final output. Even a “simple problem” might benefit more from a cooperative strategy, as this could foster a richer blend of stylistic diversity by combining the unique perspectives of different LLMs. In such cases, the goal is not just to find a correct solution but to generate a high-quality, creative artifact. Therefore, for such specific task settings, we need to incorporate additional conditioning variables to improve the robustness of Dynamic Interaction.

D More Experiment Results

D.1 Comparison with z-score normalization

To investigate the effect of different normalization methods in reward shaping (Section 3.3), we compare our min-max normalization with z-score normalization:

$$\delta_{i,l} = \frac{R_{i,k} - R_{l,avg}}{R_{l,std}} \quad (5)$$

We report the average accuracy on five mathematical evaluation benchmarks and the clip ratio discussed in Appendix B.4. Table 8 shows that implementing z-score normalization within the ILR framework yields inferior performance compared to our method. We attribute this degradation to the aggressive clipping associated with z-score normalization, which tends to suppress the magnitude of incentive signals. This observation further corroborates the superiority of our min-max normalization.

D.2 Confidence Intervals of Table 1

In Table 9, we report the 95% confidence intervals of Table 1. Note that except for the result in Table 1 and Table 9, which are based on 5 random seeds, all the others are based on the result of seed 0.

D.3 Comparison with More Self-learning Methods

In Table 1, we compare ILR with several basic self-learning methods like SFT and GRPO. Here, we discuss more self-learning baselines and compare ILR with them based on Llama-3.1-8B-Instruct.

Self-reflection Learning. Self-reflection and self-improvement constitute a pivotal research trajectory for enhancing agent capabilities through an intra-agent feedback loop (Renze and Guven, 2024; Shinn et al., 2023; Dou et al., 2024; Zhang et al., 2024; He et al., 2025). The central premise is that an agent can refine its performance by critically analyzing its own generated trajectories, answers, or reasoning processes in a closed-loop manner. To rigorously benchmark our multi-agent approach against this single-agent paradigm, we employ ReReST as a representative baseline. This comparison aims to determine whether the diverse feedback

	GSM8K	MATH-500	Minerva Math	Olympiad Bench	AIME 24&25	Avg
Llama-3.1-8B-Instruct	± 1.78	± 1.33	± 1.79	± 0.83	± 1.13	± 0.69
SFT	± 0.39	± 1.61	± 1.83	± 1.12	± 0.92	± 0.68
DPO	± 0.63	± 1.48	± 1.07	± 0.94	± 1.13	± 0.18
PPO	± 1.74	± 1.62	± 3.08	± 0.93	± 2.36	± 0.93
GRPO	± 0.72	± 3.43	± 1.82	± 2.16	± 2.36	± 1.13
GSPO	± 0.81	± 0.32	± 0.40	± 0.25	± 0.00	± 0.34
Reinforce++	± 1.16	± 2.66	± 1.65	± 1.71	± 2.27	± 0.72
DebateFT-Group1	± 1.78	± 3.01	± 1.73	± 2.40	± 1.85	± 1.85
DebateFT-Group2	± 0.73	± 2.72	± 4.32	± 3.35	± 1.85	± 2.03
ILR-Group1	± 0.70	± 0.97	± 0.24	± 0.18	± 1.14	± 0.55
ILR-Group2	± 0.21	± 1.14	± 1.18	± 0.71	± 1.46	± 0.77
Qwen2.5-7B-Instruct	± 0.57	± 1.35	± 1.84	± 0.62	± 2.78	± 1.20
SFT	± 0.86	± 2.10	± 1.63	± 1.44	± 3.70	± 1.43
DPO	± 0.62	± 2.02	± 1.50	± 0.33	± 3.40	± 1.10
PPO	± 0.83	± 1.39	± 2.03	± 0.82	± 2.07	± 1.23
GRPO	± 0.39	± 1.12	± 1.12	± 1.05	± 2.70	± 0.97
GSPO	± 0.31	± 1.65	± 0.62	± 0.93	± 0.93	± 0.82
Reinforce++	± 0.24	± 0.88	± 1.14	± 0.53	± 1.85	± 0.80
DebateFT-Group1	± 0.54	± 1.72	± 2.88	± 1.37	± 0.93	± 1.10
DebateFT-Group3	± 0.41	± 0.59	± 3.03	± 1.08	± 1.85	± 0.86
ILR-Group1	± 0.31	± 0.46	± 1.18	± 0.28	± 2.26	± 0.58
ILR-Group3	± 0.20	± 0.25	± 0.68	± 0.28	± 2.26	± 0.48
Qwen2.5-14B-Instruct	± 0.18	± 1.45	± 2.32	± 1.18	± 3.76	± 1.13
SFT	± 0.28	± 1.02	± 2.11	± 1.17	± 3.07	± 1.23
DPO	± 0.27	± 0.79	± 2.54	± 2.39	± 1.14	± 0.95
PPO	± 0.49	± 1.33	± 1.84	± 1.91	± 2.53	± 1.43
GRPO	± 0.72	± 1.50	± 2.28	± 1.13	± 3.14	± 1.47
GSPO	± 0.17	± 0.42	± 0.38	± 0.19	± 0.93	± 0.36
Reinforce++	± 0.25	± 1.22	± 0.95	± 1.79	± 4.86	± 1.42
DebateFT-Group2	± 0.18	± 0.97	± 2.05	± 1.11	± 1.73	± 0.30
DebateFT-Group3	± 0.56	± 0.75	± 2.63	± 1.25	± 3.70	± 1.41
ILR-Group2	± 0.08	± 0.21	± 0.41	± 0.43	± 2.53	± 0.72
ILR-Group3	± 0.14	± 0.33	± 0.99	± 0.51	± 2.70	± 0.86

Table 9: 95% Confidence Intervals of Table 1.

from peers offers superior gradients for improvement compared to the internalized self-verification.

Advanced GRPO Variants. Recent advancements in reinforcement learning have introduced sophisticated variants of GRPO (Lin et al., 2025; Sun et al., 2025), which significantly outperform the vanilla GRPO implementation. By comparing against CPPO (Lin et al., 2025), we seek to unveil whether the multi-agent enhancement provides a distinct advantage over purely algorithmic improvements in single-agent training stability.

GRPO with More Sampling (Compute-Matched Baseline). To rule out the possibility that performance gains are simply artifacts of increased sampling, we introduce a computation-

aligned baseline: single-agent GRPO with doubled rollout numbers (denoted as 2-GRPO). Since our multi-agent setting involves interactions that effectively increase the inference load, this variant establishes a fair comparison by equating the sampling volume. By contrasting ILR with this high-sampling baseline, we aim to demonstrate that the benefits stem from the quality of interactive feedback rather than from more sampling.

The results presented in Table 10 demonstrate that ILR consistently outperforms the extended suite of self-learning baselines, substantiating the distinct advantage of our multi-agent framework. Notably, ILR surpasses the self-reflection method Re-ReST, suggesting that diverse feedback from

	GSM8K	MATH-500	Minerva Math	Olympiad Bench	AIME 24&25	Avg
Llama-3.1-8B-Instruct	82.87	49.80	22.79	13.63	1.67	34.15
Re-ReST	84.23	49.80	29.04	15.56	5.00	36.73
CPPO	84.84	54.60	27.57	18.85	6.67	38.51
2-GRPO	86.43	53.20	28.31	18.22	6.67	38.57
ILR-Group1	89.39	55.80	30.15	22.22	10.00	41.51
ILR-Group2	87.26	55.20	33.82	22.52	6.67	41.10

Table 10: Comparison with Re-ReST, CPPO, and 2-GRPO, on Llama-3.1-8B-Instruct.

Model	Avg	periodic-Avg	proportion	periodic-proportion
Llama-3.1-8B-Instruct	41.51	42.35 (\uparrow 0.84)	36.30/63.70%	32.44/67.56%
Qwen2.5-7B-Instruct	54.44	54.82 (\uparrow 0.38)		

Table 11: Comparison with periodic updates of LLM’s reasoning ability γ_i (periodic-) based on Group1. We report the average accuracy on five mathematical evaluation benchmarks and the proportion of cooperation/competition.

Model	Avg	mixed-Avg	proportion	mixed-proportion
Llama-3.1-8B-Instruct	41.51	41.91 (\uparrow 0.40)	36.30/63.70%	40.99/59.01%
Qwen2.5-7B-Instruct	54.44	53.97 (\downarrow 0.47)		

Table 12: Comparison with mixed types interaction (mixed-) based on Group1. We report the average accuracy on five mathematical evaluation benchmarks and the proportion of cooperation/competition.

peer agents effectively mitigates the inherent blind spots of single-agent introspection. Moreover, ILR demonstrates its superiority over CPPO, underscoring that the benefits of multi-agent enhancement transcend purely algorithmic improvements in single-agent stability. Importantly, ILR also outperforms the computation-aligned 2-GRPO baseline, which confirms that the effectiveness of ILR stems from the high quality of interactive rather than the sheer quantity of samples, demonstrating greater data efficiency than force-scaling of rollouts.

D.4 Comparison with periodic updates of LLM’s reasoning ability

As a model’s reasoning ability can indeed evolve during training. To investigate this, we conduct an additional experiment to study the impact of dynamically updating γ_i . Specifically, we re-evaluated γ_i periodically every 5 training steps (we have a total of 42 steps). We report the average accuracy on five mathematical evaluation benchmarks and the proportion of cooperation/competition. Table 11 indicates that the performance of periodic-updates surpasses that of one-time measurement.

As the capability gap diminishes, the cooperation rate declines, a dynamic that appears more organic, thereby further enhancing performance.

D.5 Comparison with mixed types interaction

In real human collaboration, cooperation and competition often coexist within the same reasoning process. To investigate this scenario, we modified our sampling procedure as follows: for a given problem, even when our IRT-based mechanism determines to use a cooperation strategy, we enforce a mixture of strategies across the multiple samples. For example, if we generated five distinct rollouts for a single question, we would randomly assign the designated cooperation strategy to 3-5 of these rollouts, while assigning the competition strategy to the remaining rollouts.

As shown in Table 12, this shift yielded asymmetric outcomes. Llama-3.1-8B-Instruct achieved a performance gain of 0.40%, likely because the mixed strategy exposed it to more cooperative reasoning paths on challenging problems. In contrast, Qwen2.5-7B-Instruct experienced a regression of 0.47%. This decline suggests that the “mixed” samples acted as noise for the stronger model, where

Model	GRPO	ILR-Group1	ILR- Group2	ILR-Group3
Llama-3.1-8B-Instruct	38.39[3.75]	41.51[8.00]	41.10[12.00]	-
Qwen2.5-7B-Instruct	53.58[3.25]	54.44[8.00]	-	54.59[10.00]
Qwen2.5-14B-Instruct	57.55[5.50]	-	58.95[12.00]	59.30[10.00]

Table 13: Quantitative analysis of computational overhead versus performance gains. We report the average accuracy on five mathematical evaluation benchmarks (% , without []) and the training time (hours, within []).

	GSM8K	MATH-500	Minerva Math	Olympiad Bench	AIME 24&25	Avg
Llama-3.1-8B-Instruct	82.87	49.80	22.79	13.63	1.67	34.15
Prompt1	87.79	54.20	29.41	20.30	6.67	39.67
Prompt2	87.72	55.60	30.51	20.15	6.67	40.13
ILR-Group1	89.39	55.80	30.15	22.22	10.00	41.51
ILR-Group2	87.26	55.20	33.82	22.52	6.67	41.10

Table 14: Comparison of the same LLM with different prompts based on Llama-3.1-8B-Instruct.

	Optimization algorithms	GSM8K	MATH -500	Minerva Math	Olympiad Bench	AIME 24&25	Avg
Llama-3.1-8B-Instruct	-	82.87	49.80	22.79	13.63	1.67	34.15
ILR-Group1	GRPO	89.39	55.80	30.15	22.22	10.00	41.51
ILR-Group2	GRPO	87.26	55.20	33.82	22.52	6.67	41.10
ILR-Group1	GSPO	89.72	56.00	31.25	21.78	10.00	41.75↑
ILR-Group2	GSPO	89.58	54.80	30.15	23.15	10.00	41.54↑
Qwen2.5-7B-Instruct	-	92.34	75.60	41.54	37.63	11.67	51.76
ILR-Group1	GRPO	93.40	77.60	43.01	39.85	18.33	54.44
ILR-Group3	GRPO	92.65	78.00	45.01	38.96	18.33	54.59
ILR-Group1	GSPO	93.48	77.80	45.22	40.15	16.67	54.66↑
ILR-Group3	GSPO	92.95	78.60	44.85	40.30	20.00	55.34↑
Qwen2.5-14B-Instruct	-	94.84	81.20	47.43	41.04	13.33	55.57
ILR-Group2	GRPO	95.53	81.80	50.37	43.70	23.33	58.95
ILR-Group3	GRPO	95.30	82.60	51.10	44.15	23.33	59.30
ILR-Group2	GSPO	95.98	82.80	50.37	43.70	23.33	59.24↑
ILR-Group3	GSPO	95.75	82.20	52.37	43.96	23.33	59.52↑

Table 15: Quantitative comparison of using different reinforcement learning algorithms (GRPO vs. GSPO) as the ILR optimization algorithms. Green lines are based on GRPO, and pink lines are based on GSPO.

forcing it to collaborate (when it should compete) likely constrained its reasoning potential or introduced errors from the weaker partner.

D.6 Cost Table

We compare the training time and performance of GRPO and ILR in Table 13. Note that the training durations reported for ILR-Group i represent the aggregated time required for both LLMs within the group. For instance, in the ILR-Group1 setting, the combined training time for Llama-3.1-8B-Instruct

and Qwen2.5-7B-Instruct amounts to 8 hours. In comparison, training these models individually using standard GRPO requires a total of 7 hours (3.75 hours for Llama-3.1-8B-Instruct and 3.25 hours for Qwen2.5-7B-Instruct). Consequently, incorporating ILR introduces only a marginal computational overhead of one hour. This slight increase in cost is justified by the performance gains, yielding improvements of approximately 3% for Llama-3.1-8B-Instruct and 1% for Qwen2.5-7B-Instruct.

D.7 Same LLM with different input prompts

In our primary experiments, we conceptualized distinct LLMs receiving the same input prompt as individual agents to investigate interactive learning among heterogeneous models. In this section, we extend this inquiry by examining the efficacy of interactive learning within a homogeneous setting, where a single LLM utilizes different input prompts to simulate distinct agents.

- **Prompt1:** $\langle \text{Question} \rangle$ Please reason step by step.
- **Prompt2:** $\langle \text{Question} \rangle$ Please double-check the question, first extract the keywords and think about the knowledge about these keywords, then solve the question.

Two notable findings emerge from this analysis (see Table 14). **First**, while ILR with an identical LLM still yields performance improvements over single-agent baselines, it does not surpass the performance of ILR with different LLMs. This suggests that while diversity can be simulated through varied prompting, distinct models provide a richer and more effective source of alternative perspectives. The interaction between truly heterogeneous agents appears to generate a more powerful learning signal. **Second**, the post-training performance of the Llama model utilizing Prompt2 was superior to its counterpart. This observation underscores that the granularity of instruction prompts plays a significant role in influencing the outcomes of interactive learning.

D.8 Extendability of Optimization Algorithm

In fact, the optimization algorithm used in our ILR is not limited to GRPO and can be replaced by another method. Any reinforcement learning (RL) algorithm capable of generating a group of responses can serve as our optimization algorithm. Therefore, in this section, we investigate the impact of different optimization algorithms (GRPO vs. GSPO) on ILR. As shown in the baseline results of Table 1, GSPO outperforms GRPO across all LLMs. Meanwhile, from Table 15, it can be observed that ILR with GSPO achieves better results than ILR using GRPO. This demonstrates the high extensibility and universality of ILR with respect to RL algorithms. With the advancement of more sophisticated RL algorithms in the future, ILR is expected to achieve even stronger performance.

E Pseudocode of ILR

Algorithm 1 Interactive Learning for LLM Reasoning (ILR Training Algorithm)

- 1: **Input:** Training set $\mathcal{D} = \{q_1, \dots, q_N\}$, Validation set \mathcal{D}_{val} , Set of LLMs $\mathcal{M} = \{M_1, \dots, M_m\}$, Reward model \mathcal{M}_{reward} , Self-ranking splits S , Batch size N' , Sampling rounds n
- 2: **Output:** Optimized Models \mathcal{M}^*

Phase 1: Ability and Difficulty Estimation

- 3: Estimate model ability γ_i for each $M_i \in \mathcal{M}$ based on accuracy on \mathcal{D}_{val}
- 4: **for** each question $q \in \mathcal{D}$ **do**
- 5: Divide \mathcal{D} into random batches of size N' across S splits
- 6: Collect ranks $r_{q,j}$ from M_i via self-ranking prompt
- 7: $D_q \leftarrow \frac{1}{m} \sum_{i=1}^m M_i \left(\frac{1}{S} \sum_{j=1}^S \frac{r_{q,j}}{N'} \right)$ ▷ Eq. 1
- 8: **end for**

Phase 2: Training Loop

- 9: **for** each batch $\mathcal{B} \subset \mathcal{D}$ **do**
 - 10: $\mathcal{T} \leftarrow \emptyset, \mathcal{R} \leftarrow \emptyset, \bar{\mathcal{R}} \leftarrow \emptyset$
 - 11: **for** each question $q \in \mathcal{B}$ **do**
 - 12: // 1. Dynamic Interaction
 - 13: Calculate success prob.
 - 14: $P_{q,i} \leftarrow \frac{1}{1+e^{-1.7(\gamma_i - D_q)}}$ for all M_i ▷ Eq. 2
 - 15: $P_q \leftarrow \frac{1}{m} \sum_{i=1}^m P_{q,i}$
 - 16: **if** $P_q < 0.5$ **then** Mode \leftarrow Cooperation
 - 17: **else** Mode \leftarrow Competition
 - 18: **end if** ▷ Eq. 3
 - 19: // 2. Idea3 Interaction
 - 20: **for** $k = 1$ to $n, i = 1$ to m **do**
 - 21: $a_{init} \leftarrow M_i.\text{Gen}(q, \text{Prompt}_S)$
 - 22: $a_{ana} \leftarrow M_i.\text{Gen}(q, a_{init}^{\text{peer}}, \text{Mode}, \text{Prompt}_A)$
 - 23: $a_{fus} \leftarrow M_i.\text{Gen}(q, a_{ana}, \text{Mode}, \text{Prompt}_E)$
 - 24: **if** $\text{IsCorrect}(a_{init}) \wedge \neg \text{IsCorrect}(a_{fus})$ **then**
 - 25: $a_{final} \leftarrow a_{init}$
 - 26: **else**
 - 27: $a_{final} \leftarrow a_{fus}$
 - 28: **end if**
 - 29: $\tau_{i,k} \leftarrow (q, a_{final}); \mathcal{T}.\text{add}(\tau_{i,k})$
 - 30: $R_{i,k} \leftarrow \mathcal{M}_r(\tau_{i,k}); \mathcal{R}.\text{add}(R_{i,k})$
 - 31: **end for**
 - 32: // 3. Perception Calibration
 - 33: **for** $i = 1$ to m **do**
 - 34: **for** $k = 1$ to n **do**
 - 35: $V_{calib} \leftarrow 0$
 - 36: **for** $l \in \mathcal{M} \setminus \{M_i\}$ **do**
 - 37: Get peer stats $R_{l,max}, R_{l,min}, R_{l,avg}$
 - 38: $\Delta \leftarrow (R_{i,k} - R_{l,avg}) / (R_{l,max} - R_{l,min})$
 - 39: $V_{calib} += \text{clip}(\Delta, -\frac{1}{m-1}, +\frac{1}{m-1})$
 - 40: **end for**
 - 41: $\bar{R}_{i,k} \leftarrow R_{i,k} + V_{calib}$ ▷ Eq. 4
 - 42: $\bar{\mathcal{R}}.\text{add}(\bar{R}_{i,k})$
 - 43: **end for**
 - 44: **end for**
 - 45: **end for**
 - 46: // 4. Optimization
 - 47: **for** $i = 1$ to m **do**
 - 48: Update M_i via GRPO using trajectories \mathcal{T} and calibrated rewards $\bar{\mathcal{R}}$
 - 49: **end for**
 - 50: **end for**
-

F Case Study

We show some cases of **Dynamic Interaction (Idea3)** based on Group1 (LLM1: Llama-3.1-8B-Instruct and LLM2: Qwen2.5-7B-Instruct).

Cooperation mode (both LLMs initially fail but succeed after interaction)

Problem:

The function $f(x)$ satisfies $f(x) + f\left(\frac{x+1}{1-3x}\right) = x$ for all $x \neq \frac{1}{3}$. Find $f(2)$.

LLM1 (Idea Sharing):

To find $f(2)$, we can use the given functional equation and plug in $x = 2$. This gives us

$$f(2) + f\left(\frac{2+1}{1-3 \cdot 2}\right) = 2 \Rightarrow f(2) + f\left(-\frac{3}{5}\right) = 2$$

Now, we can plug in $x = -\frac{3}{5}$ into the original functional equation to get

$$\begin{aligned} f\left(-\frac{3}{5}\right) + f\left(\frac{-\frac{3}{5}+1}{1-3 \cdot (-\frac{3}{5})}\right) &= -\frac{3}{5} \\ \Rightarrow f\left(-\frac{3}{5}\right) + f\left(\frac{1}{7}\right) &= -\frac{3}{5} \end{aligned}$$

Now, we can plug in $x = \frac{1}{7}$ into the original functional equation to get

$$f\left(\frac{1}{7}\right) + f\left(\frac{\frac{1}{7}+1}{1-3 \cdot \frac{1}{7}}\right) = \frac{1}{7} \Rightarrow f\left(\frac{1}{7}\right) + f(2) = \frac{1}{7}$$

Now we have two equations and two unknowns. We can solve for $f(2)$ by subtracting the second equation from the first equation:

$$\begin{aligned} f(2) + f\left(-\frac{3}{5}\right) - \left(f\left(\frac{1}{7}\right) + f(2)\right) &= 2 - \frac{1}{7} \\ \Rightarrow f\left(-\frac{3}{5}\right) - f\left(\frac{1}{7}\right) &= \frac{13}{7} \end{aligned}$$

Now we can plug in $f\left(-\frac{3}{5}\right) = 2 - f\left(\frac{1}{7}\right)$ into the equation $f\left(-\frac{3}{5}\right) + f\left(\frac{1}{7}\right) = -\frac{3}{5}$ to get

$$2 - f\left(\frac{1}{7}\right) - f\left(\frac{1}{7}\right) = -\frac{3}{5} \Rightarrow f\left(\frac{1}{7}\right) = \frac{13}{10}$$

Now we can plug in $f\left(\frac{1}{7}\right) = \frac{13}{10}$ into the equation $f\left(-\frac{3}{5}\right) - f\left(\frac{1}{7}\right) = \frac{13}{7}$ to get

$$f\left(-\frac{3}{5}\right) - \frac{13}{10} = \frac{13}{7} \Rightarrow f\left(-\frac{3}{5}\right) = \frac{221}{70}$$

Now we can plug in $f\left(-\frac{3}{5}\right) = \frac{221}{70}$ into the equation $f(2) + f\left(-\frac{3}{5}\right) = 2$ to get

$$f(2) + \frac{221}{70} = 2 \Rightarrow f(2) = -\frac{81}{70}$$

The final answer is: $\boxed{-\frac{81}{70}}$. ✖

LLM2 (Idea Sharing):

To solve for $f(2)$ given the functional equation

$$f(x) + f\left(\frac{x+1}{1-3x}\right) = x$$

We need to find a way to express $f(x)$ in a form that allows us to determine $f(2)$. First, let's denote the transformation $g(x) = \frac{x+1}{1-3x}$. We need to find the inverse of $g(x)$ to understand the behavior of the function under repeated applications. To find the inverse, we solve for x in terms of y where $y = g(x)$:

$$y = \frac{x+1}{1-3x} \Rightarrow x = \frac{y-1}{1+3y}$$

Thus, the inverse function is

$$g^{-1}(x) = \frac{x-1}{1+3x}$$

Now, we apply the functional equation twice. First, we substitute x with $g(x)$:

$$f(g(x)) + f\left(\frac{g(x)+1}{1-3g(x)}\right) = g(x)$$

Since $g\left(\frac{x+1}{1-3x}\right) = x$, we have:

$$f(g(x)) + f(x) = g(x)$$

Next, we substitute x with $g^{-1}(x)$:

$$f(g^{-1}(x)) + f\left(\frac{g^{-1}(x)+1}{1-3g^{-1}(x)}\right) = g^{-1}(x)$$

Since $g\left(\frac{x-1}{1+3x}\right) = x$, we have:

$$f(g^{-1}(x)) + f(x) = g^{-1}(x)$$

We now have two equations: 1. $f(g(x)) + f(x) = g(x)$ 2. $f(g^{-1}(x)) + f(x) = g^{-1}(x)$. To find $f(2)$, we need to determine the value of x such that $g(x) = 2$:

$$2 = \frac{x+1}{1-3x} \Rightarrow x = \frac{1}{7}$$

Thus, $g\left(\frac{1}{7}\right) = 2$. Using the first equation with $x = \frac{1}{7}$:

$$f(2) + f\left(\frac{1}{7}\right) = 2$$

Next, we need to find $f\left(\frac{1}{7}\right)$. We use the second equation with $x = 2$:

$$\begin{aligned} f\left(\frac{2-1}{1+3 \cdot 2}\right) + f(2) &= \frac{2-1}{1+3 \cdot 2} \\ \Rightarrow f\left(\frac{1}{7}\right) + f(2) &= \frac{1}{7} \end{aligned}$$

We now have the system of linear equations: $f(2) + f\left(\frac{1}{7}\right) = 2$ and $f\left(\frac{1}{7}\right) + f(2) = \frac{1}{7}$. This is a contradiction,

indicating an error in the setup. Revisiting the problem, we find: $f(2) = \frac{2}{3}$. Thus, the correct value is: $\frac{2}{3}$. ✗

LLM1 (Idea Analysis):

The partner's contribution provides a clear and step-by-step approach to solving the problem. The key steps that can help improve the original answer are: 1. Defining the transformation: The partner defines the transformation $g(x) = \frac{x+1}{1-3x}$ and finds its inverse. This step helps to understand the behavior of the function under repeated applications. 2. Applying the functional equation twice: The partner applies the functional equation twice, once with x and once with $g^{-1}(x)$. This step leads to two equations that can be used to solve for $f(2)$. This step is similar to the step-by-step plug-in solution. 3. Finding the value of x such that $g(x) = 2$: The partner finds the value of x such that $g(x) = 2$, which is $\frac{1}{7}$. This step is crucial for using the first equation to solve for $f(2)$. However, this method is complex, and there is an error in the partner's contribution. The system of linear equations leads to a contradiction, indicating that the approach needs to be revised. We can double-check the step-by-step plug-in solution to ensure its correctness.

LLM2 (Idea Analysis):

Let's analyze the key steps in the partner's contribution and identify areas for improvement. Key Steps in the Partner's Contribution: 1. Substitute $x = 2$ into the functional equation. 2. Substitute $x = -\frac{3}{5}$ into the functional equation. 3. Substitute $x = \frac{1}{7}$ into the functional equation. 4. Solve the system of equations. But we can make some improvements: instead of solving the system step-by-step, we can directly solve for $f(2)$ by eliminating intermediate variables to avoid potential mistakes.

LLM1 (Idea Fusion):

To find $f(2)$, we can use the given functional equation and plug in $x = 2$. This gives us

$$f(2) + f\left(\frac{2+1}{1-3 \cdot 2}\right) = 2 \Rightarrow f(2) + f\left(-\frac{3}{5}\right) = 2$$

Now, we can plug in $x = -\frac{3}{5}$ into the original functional equation to get

$$\begin{aligned} f\left(-\frac{3}{5}\right) + f\left(\frac{-\frac{3}{5}+1}{1-3 \cdot \left(-\frac{3}{5}\right)}\right) &= -\frac{3}{5} \\ \Rightarrow f\left(-\frac{3}{5}\right) + f\left(\frac{1}{7}\right) &= -\frac{3}{5} \end{aligned}$$

Now, we can plug in $x = \frac{1}{7}$ into the original functional equation to get

$$\begin{aligned} f\left(\frac{1}{7}\right) + f\left(\frac{\frac{1}{7}+1}{1-3 \cdot \frac{1}{7}}\right) &= \frac{1}{7} \\ \Rightarrow f\left(\frac{1}{7}\right) + f(2) &= \frac{1}{7} \end{aligned}$$

Now we have two equations and two unknowns. We can solve for $f(2)$ by adding the two equations:

$$\begin{aligned} f(2) + f\left(-\frac{3}{5}\right) + f\left(\frac{1}{7}\right) + f(2) &= 2 + \frac{1}{7} \\ \Rightarrow 2f(2) + f\left(-\frac{3}{5}\right) + f\left(\frac{1}{7}\right) &= \frac{15}{7} \end{aligned}$$

Now we can plug in $f\left(-\frac{3}{5}\right) + f\left(\frac{1}{7}\right) = -\frac{3}{5}$ into the equation to get

$$2f(2) - \frac{3}{5} = \frac{15}{7} \Rightarrow f(2) = \frac{48}{35}$$

The final answer is: $\frac{48}{35}$. ✓

LLM2 (Idea Fusion):

To solve for $f(2)$ given the functional equation

$$f(x) + f\left(\frac{x+1}{1-3x}\right) = x$$

We will follow a systematic approach. 1. Substitute $x = 2$ into the functional equation:

$$f(2) + f\left(\frac{2+1}{1-3 \cdot 2}\right) = 2 \Rightarrow f(2) + f\left(-\frac{3}{5}\right) = 2$$

2. Substitute $x = -\frac{3}{5}$ into the functional equation:

$$\begin{aligned} f\left(-\frac{3}{5}\right) + f\left(\frac{-\frac{3}{5}+1}{1-3 \cdot \left(-\frac{3}{5}\right)}\right) &= -\frac{3}{5} \\ \Rightarrow f\left(-\frac{3}{5}\right) + f\left(\frac{1}{7}\right) &= -\frac{3}{5} \end{aligned}$$

3. Substitute $x = \frac{1}{7}$ into the functional equation:

$$f\left(\frac{1}{7}\right) + f\left(\frac{\frac{1}{7}+1}{1-3 \cdot \frac{1}{7}}\right) = \frac{1}{7} \Rightarrow f\left(\frac{1}{7}\right) + f(2) = \frac{1}{7}$$

4. Solve the system of equations: We have the following system of equations:

$$f(2) + f\left(-\frac{3}{5}\right) = 2 \quad (1)$$

$$f\left(-\frac{3}{5}\right) + f\left(\frac{1}{7}\right) = -\frac{3}{5} \quad (2)$$

$$f\left(\frac{1}{7}\right) + f(2) = \frac{1}{7} \quad (3)$$

5. Add equations (1) and (3):

$$(f(2) + f\left(-\frac{3}{5}\right)) + (f\left(\frac{1}{7}\right) + f(2)) = 2 + \frac{1}{7}$$

$$2f(2) + f\left(-\frac{3}{5}\right) + f\left(\frac{1}{7}\right) = \frac{15}{7}$$

6. Substitute equation (2) into the result:

$$2f(2) + \left(-\frac{3}{5}\right) = \frac{15}{7}$$

$$2f(2) - \frac{3}{5} = \frac{15}{7}$$

7. Solve for $f(2)$:

$$2f(2) = \frac{15}{7} + \frac{3}{5}$$

Find a common denominator:

$$2f(2) = \frac{75}{35} + \frac{21}{35} = \frac{96}{35}$$

$$f(2) = \frac{48}{35}$$

Thus, the value of $f(2)$ is $\frac{48}{35}$. ✓

Competition mode (LLMs raise better solutions after interaction)

Problem:

Find the minimum value of $\frac{x^2+7}{\sqrt{x^2+3}}$ over all real numbers x .

LLM1 (Idea Sharing):

To find the minimum value of the given expression, we can use the following steps:

Step 1: Rewrite the expression. We can rewrite the given expression as

$$\frac{x^2+7}{\sqrt{x^2+3}} = \frac{x^2+3+4}{\sqrt{x^2+3}} = \frac{x^2+3}{\sqrt{x^2+3}} + \frac{4}{\sqrt{x^2+3}}$$

Step 2: Simplify the expression. Simplifying the expression, we get

$$\frac{x^2+3}{\sqrt{x^2+3}} + \frac{4}{\sqrt{x^2+3}} = \sqrt{x^2+3} + \frac{4}{\sqrt{x^2+3}}$$

Step 3: Use the AM-GM inequality. Using the AM-GM inequality, we have

$$\sqrt{x^2+3} + \frac{4}{\sqrt{x^2+3}} \geq 2\sqrt{\sqrt{x^2+3} \cdot \frac{4}{\sqrt{x^2+3}}} = 4$$

Step 4: Find the equality condition. The equality condition for the AM-GM inequality is when $\sqrt{x^2+3} = \frac{4}{\sqrt{x^2+3}}$. Solving for x , we get $x^2+3 = 4$, which gives $x = \pm 1$.

Therefore, the final answer is: $\boxed{4}$. ✓

LLM2 (Idea Sharing):

To find the minimum value of the function $f(x) = \frac{x^2+7}{\sqrt{x^2+3}}$, we will follow these steps:

1. Define the function and its domain:

$$f(x) = \frac{x^2+7}{\sqrt{x^2+3}}$$

The domain of $f(x)$ is all real numbers x because the expression under the square root, x^2+3 , is always positive for all real x .

2. Find the first derivative of $f(x)$: To find the critical points, we need to compute the first derivative $f'(x)$ and set it to zero.

$$f(x) = \frac{x^2+7}{\sqrt{x^2+3}}$$

Let $u = x^2+7$ and $v = \sqrt{x^2+3}$. Then $f(x) = \frac{u}{v}$. Using the quotient rule:

$$f'(x) = \frac{u'v - uv'}{v^2}$$

where $u' = 2x$ and $v' = \frac{x}{\sqrt{x^2+3}}$. Substituting u , u' , v , and v' into the quotient rule:

$$f'(x) = \frac{(2x)(\sqrt{x^2+3}) - (x^2+7)\left(\frac{x}{\sqrt{x^2+3}}\right)}{(x^2+3)}$$

Simplify the numerator:

$$f'(x) = \frac{2x\sqrt{x^2+3} - \frac{x(x^2+7)}{\sqrt{x^2+3}}}{x^2+3} = \frac{x(x^2-1)}{(x^2+3)\sqrt{x^2+3}}$$

Set the first derivative equal to zero to find critical points:

$$\frac{x(x^2-1)}{(x^2+3)\sqrt{x^2+3}} = 0 \Rightarrow x = 0, \quad x = 1, \quad x = -1$$

3. Evaluate $f(x)$ at the critical points and endpoints:

$$f(0) = \frac{0^2+7}{\sqrt{0^2+3}} = \frac{7}{\sqrt{3}} = \frac{7\sqrt{3}}{3}$$

$$f(1) = \frac{1^2+7}{\sqrt{1^2+3}} = \frac{8}{\sqrt{4}} = \frac{8}{2} = 4$$

$$f(-1) = \frac{(-1)^2+7}{\sqrt{(-1)^2+3}} = \frac{8}{\sqrt{4}} = \frac{8}{2} = 4$$

4. Determine the minimum value: Comparing the values:

$$f(0) = \frac{7\sqrt{3}}{3} \approx 4.041$$

$$f(1) = 4$$

$$f(-1) = 4$$

The minimum value is 4. Therefore, the minimum value of the function is $\boxed{4}$. ✓

LLM1 (Idea Analysis):

The opponent's solution is a thorough and detailed approach to finding the minimum value of the given function. Here's a critical analysis of his ideas:

Strengths: 1. Clear and concise explanation: The opponent's solution is well-structured and easy to follow, making it easy to understand their thought process. 2. Use of mathematical notation: The opponent uses mathematical notation consistently and correctly, which helps to clarify their ideas and make the solution more readable. 3. Comprehensive analysis: The opponent considers all possible critical points, including the endpoints, and evaluates the function at each point to determine the minimum value.

Weaknesses: 1. Overcomplication: The opponent's solution involves a lot of algebraic manipulation, which can be overwhelming and may not be necessary to find the minimum value of the function. 2. Lack of insight: The opponent's solution relies heavily on mathematical notation and formulas, but doesn't provide much insight into the underlying structure of the function or the reason why the minimum value is 4. 3. No comparison with other methods: The opponent's solution doesn't compare their approach with other methods, such as the AM-GM inequality, which is a more elegant and straightforward way to find the minimum value.

In conclusion, while the opponent's solution is thorough and detailed, it is overcomplicated and lacks insight into the underlying structure of the function. The alternative approach using the AM-GM inequality provides a more elegant and straightforward way to find the minimum value of the function.

LLM2 (Idea Analysis):

The opponent's solution is a well-structured approach to finding the minimum value of the given expression. Let's critically analyze the steps and identify the strengths and weaknesses:

Strengths: 1. Rewriting the Expression: The opponent successfully rewrites the expression in a form that simplifies the problem. This step is crucial as it breaks down the original expression into more manageable

parts. 2. Simplification: The expression is simplified to $\sqrt{x^2 + 3} + \frac{4}{\sqrt{x^2 + 3}}$, which is a standard form that can be analyzed further. 3. Application of AM-GM Inequality: The use of the Arithmetic Mean-Geometric Mean (AM-GM) inequality is a powerful and elegant method to find the minimum value. This step is mathematically sound and provides a clear path to the solution.

Weaknesses: 1. Verification of Minimum Value: While the AM-GM inequality provides a lower bound, it does not always guarantee that the minimum value is achieved. The opponent assumes that the minimum value is 4 without explicitly showing that no smaller value is possible. A more rigorous approach would involve checking the behavior of the function as x approaches infinity or negative infinity. 2. Verification of Minimum Value: A more thorough analysis would involve checking the derivative to confirm that these points are indeed the global minimum. Conclusion: The opponent's solution is strong and mathematically sound, but it could be made more rigorous by explicitly showing that the minimum value is achieved at $x = \pm 1$ and that no smaller value is possible. The final answer is correct, but a more detailed verification step would strengthen the solution.

LLM1 (Idea Fusion):

To find the minimum value of the given expression, we can use the AM-GM inequality.

Step 1: Apply the AM-GM inequality. The AM-GM inequality states that for any non-negative real numbers a and b , we have

$$\frac{a + b}{2} \geq \sqrt{ab}.$$

We can rewrite the given expression as

$$\frac{x^2 + 7}{\sqrt{x^2 + 3}} = \frac{x^2 + 3 + 4}{\sqrt{x^2 + 3}} = \sqrt{x^2 + 3} + \frac{4}{\sqrt{x^2 + 3}}.$$

Step 2: Use the AM-GM inequality on the expression. Applying the AM-GM inequality to the expression, we get

$$\sqrt{x^2 + 3} + \frac{4}{\sqrt{x^2 + 3}} \geq 2\sqrt{\sqrt{x^2 + 3} \cdot \frac{4}{\sqrt{x^2 + 3}}} = 4.$$

Step 3: Find the equality condition. The equality condition for the AM-GM inequality is when $\sqrt{x^2 + 3} = \frac{4}{\sqrt{x^2 + 3}}$. Solving for x , we get $x^2 + 3 = 4$, which gives $x = \pm 1$.

Step 4: Check the equality condition. We need to check if the equality condition is satisfied for $x = \pm 1$. Substituting $x = \pm 1$ into the original expression, we get

$$\frac{(\pm 1)^2 + 7}{\sqrt{(\pm 1)^2 + 3}} = \frac{8}{2} = 4.$$

Step 5: Conclusion. Since the equality condition is satisfied for $x = \pm 1$, the minimum value of the expression is 4. The final answer is: $\boxed{4}$ ✓

LLM2 (Idea Fusion):

To find the minimum value of the function $f(x) = \frac{x^2 + 7}{\sqrt{x^2 + 3}}$, we will follow a detailed step-by-step approach:

1. Rewrite the Expression:

$$f(x) = \frac{x^2 + 7}{\sqrt{x^2 + 3}} = \sqrt{x^2 + 3} + \frac{4}{\sqrt{x^2 + 3}}$$

2. Simplify the Expression: We have simplified the expression to:

$$f(x) = \sqrt{x^2 + 3} + \frac{4}{\sqrt{x^2 + 3}}$$

3. Use the AM-GM Inequality: The Arithmetic Mean-Geometric Mean (AM-GM) inequality states that for any non-negative real numbers a and b ,

$$\frac{a + b}{2} \geq \sqrt{ab}$$

Applying this to our expression:

$$\sqrt{x^2 + 3} + \frac{4}{\sqrt{x^2 + 3}} \geq 2\sqrt{\sqrt{x^2 + 3} \cdot \frac{4}{\sqrt{x^2 + 3}}} = 4$$

Therefore,

$$\sqrt{x^2 + 3} + \frac{4}{\sqrt{x^2 + 3}} \geq 4$$

4. Find the Equality Condition: The equality in the AM-GM inequality holds when:

$$\sqrt{x^2 + 3} = \frac{4}{\sqrt{x^2 + 3}}$$

Solving for x :

$$(\sqrt{x^2 + 3})^2 = \left(\frac{4}{\sqrt{x^2 + 3}}\right)^2$$

$$x^2 + 3 = \frac{16}{x^2 + 3}$$

Let $y = \sqrt{x^2 + 3}$. Then:

$$y^2 = \frac{16}{y^2}$$

$$y^4 = 16$$

$$y^2 = 4 \quad (\text{since } y \geq 0)$$

$$y = 2$$

Therefore:

$$\sqrt{x^2 + 3} = 2$$

$$x^2 + 3 = 4$$

$$x^2 = 1$$

$$x = \pm 1$$

5. Verify the Minimum Value: Substitute $x = \pm 1$ into the original function to verify:

$$f(1) = \frac{1^2 + 7}{\sqrt{1^2 + 3}} = \frac{8}{2} = 4$$

$$f(-1) = \frac{(-1)^2 + 7}{\sqrt{(-1)^2 + 3}} = \frac{8}{2} = 4$$

6. Conclusion: Since the AM-GM inequality holds with equality at $x = \pm 1$ and the function value at these points is 4, the minimum value of the function is indeed 4. The final answer is: $\boxed{4}$ ✓