



# OPTIMA: OPTIMIZING EFFECTIVENESS AND EFFICIENCY FOR LLM-BASED MULTI-AGENT SYSTEM

**Anonymous authors**

Paper under double-blind review

## ABSTRACT

Large Language Model (LLM) based multi-agent systems (MAS) show remarkable potential in collaborative problem-solving, yet they still face critical challenges: low communication efficiency, poor scalability, and a lack of effective parameter-updating optimization methods for multi-agent collaboration. We present **OPTIMA**, a novel framework that addresses these issues by significantly enhancing *both* communication efficiency and task effectiveness in LLM-based MAS through LLM training. At its core, OPTIMA employs an *iterative generate, rank, select, and train* paradigm, incorporating a reward function that balances task performance, token efficiency, and communication readability. We explore various RL algorithms, including Supervised Fine-Tuning, Direct Preference Optimization, and their hybrid approaches, providing insights into their effectiveness-efficiency trade-offs for iterative LLM-based MAS training. Additionally, we integrate Monte Carlo Tree Search-inspired techniques for DPO data generation, conceptualizing conversation turns as tree nodes to explore diverse interaction trajectories. We evaluate OPTIMA on common multi-agent tasks, including information-asymmetric question answering and complex reasoning. Our method demonstrates consistent and substantial improvements over single-agent baselines and vanilla MAS based on Llama 3 8B, achieving up to  $2.8x$  performance gain with less than 10% tokens on tasks requiring heavy multi-agent information exchange. Moreover, OPTIMA's efficiency gains open new possibilities for leveraging inference-compute more effectively, potentially leading to improved inference-time scaling laws. By addressing fundamental challenges in multi-agent collaboration and providing a novel optimization framework, OPTIMA shows the potential towards scalable, efficient, and effective LLM-based MAS.

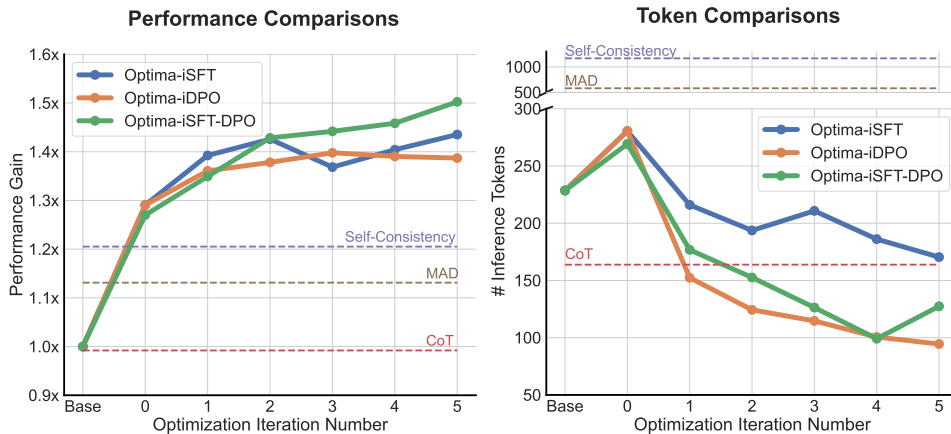


Figure 1: **Performance and efficiency of OPTIMA variants across optimization iterations. Left:** Average performance gain over iterations. OPTIMA variants consistently outperform CoT, Multi-Agent Debate (MAD), and Self-Consistency. **Right:** Average inference token numbers over iterations. All OPTIMA variants achieve better performance with substantially fewer tokens.

## 1 INTRODUCTION

Large Language Models (LLMs) have emerged as powerful tools for a wide range of tasks, from natural language processing to complex reasoning (OpenAI, 2023; Reid et al., 2024; Anthropic, 2024). A promising direction in leveraging these models is the development of autonomous multi-agent systems (MAS), which aim to harness the collective intelligence of multiple LLM-based agents for collaborative problem-solving and decision-making (Liang et al., 2023; Wang et al., 2024b; Du et al., 2024; Zhuge et al., 2024). However, for LLM-based MAS to be truly effective, they must overcome two critical challenges: **(a)** achieving efficient inter-agent communication to minimize computational costs, and **(b)** optimizing the collective performance of the system as a cohesive unit.

Current LLM-based MAS face significant difficulties in meeting these challenges. The coordination and communication between agents often lack efficiency, resulting in verbose exchanges that lead to increased token usage, longer inference times, and higher computational costs (Li et al., 2024b). This inefficiency is exacerbated by the *length bias* inherent in LLMs due to alignment training (Saito et al., 2023; Dubois et al., 2024), which favors longer responses even when concise communication would suffice (Chen et al., 2024d). Moreover, while recent work has explored training LLMs for single-agent tasks (Song et al., 2024; Xiong et al., 2024) and MAS training is well-studied in reinforcement learning (Johnson et al., 2000; Lanctot et al., 2017; Baker et al., 2020), there remains a lack of parameter-updating methods specifically designed to optimize LLM-based MAS as a unified system. Existing approaches primarily rely on simple agent profile evolution (Chen et al., 2024b) or memory evolution (Qian et al., 2024a;b; Gao et al., 2024), which fail to address the core issues of communication efficiency and collective optimization.

**Can we develop a training framework that simultaneously enhances the communication efficiency and task effectiveness of LLM-based MAS?** To address this question, we introduce **OPTIMA**, an effective framework designed to optimize LLM-based MAS. At the heart of OPTIMA is an iterative *generate, rank, select, and train* paradigm, incorporating a reward function that balances task performance, token efficiency, and communication interpretability. This approach enables the development of MAS that are not only effective and efficient but also maintain interpretable communication patterns. Based on the reward function, OPTIMA leverages a combination of techniques to induce efficient and effective communication behaviors in LLM-based agents, including Supervised Fine-Tuning (SFT) (Zelikman et al., 2022; Gülçehre et al., 2023; Aksitov et al., 2023) and Direct Preference Optimization (DPO) (Rafailov et al., 2023; Pang et al., 2024), along with their hybrid variants. Furthermore, OPTIMA introduces an integration of Monte Carlo Tree Search (MCTS)-inspired techniques for DPO data generation, conceptualizing conversation turns as tree nodes to explore diverse interaction trajectories efficiently.

Importantly, by substantially reducing the number of tokens required for inference, OPTIMA not only improves computational efficiency but also opens new possibilities for leveraging inference-compute more effectively. This reduction in token usage allows for more samples within the same computational constraints, potentially leading to *better inference-time scaling laws*. As recent work has shown the importance of inference-time compute in improving model performance (Wu et al., 2024; Brown et al., 2024; Chen et al., 2024a), OPTIMA’s efficiency gains could be combined with techniques like majority voting (Wang et al., 2023), leading to more effective LLM systems.

We evaluate OPTIMA on a diverse set of tasks spanning two multi-agent settings: **(a)** information exchange, including information-asymmetric question answering (Chen et al., 2024d; Liu et al., 2024), and **(b)** debate, encompassing mathematical and reasoning tasks (Du et al., 2024; Chen et al., 2024b; Wu et al., 2023). Using Llama 3 8B (Meta, 2024) as our base model, we demonstrate that OPTIMA consistently outperforms both single-agent MAS baselines, achieving up to 90% reduction in token usage and 2.8x increase in task performance.

To summarize, our main contribution is OPTIMA, a novel training framework that simultaneously optimizes *communication efficiency* and *task effectiveness*. To enhance high-quality training data generation in *multi-agent settings* for DPO, we introduce an integration of MCTS-like techniques. Our comprehensive empirical evaluation across diverse tasks demonstrates notable advancements in *both* token efficiency and task performance, while also providing insights into the learned communication patterns. Additionally, we examine the implications of OPTIMA’s efficiency gains for inference-time scaling laws, underscoring its potential to improve the overall capabilities of LLM systems by enabling more effective utilization of inference-compute. By addressing the dual chal-

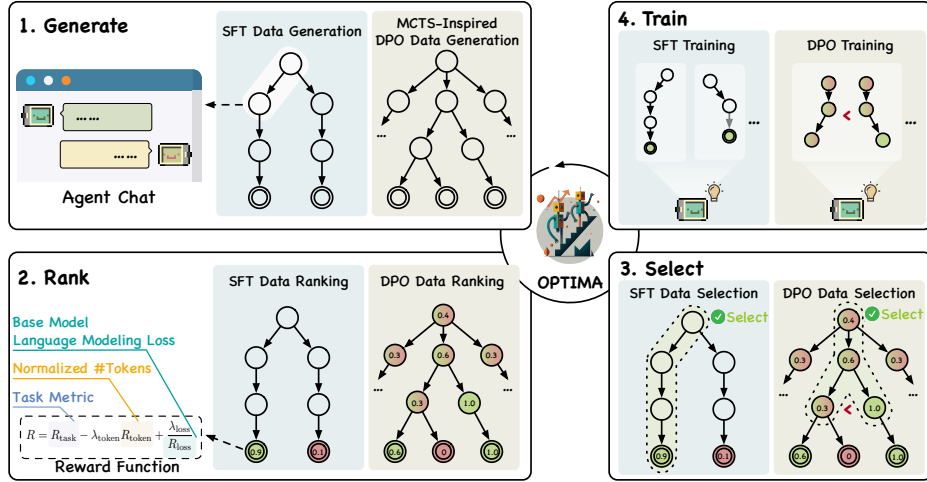


Figure 2: **Overview of the OPTIMA framework for training LLM-based MAS.** The iterative process includes four stages: *Generate*, *Rank*, *Select*, and *Train*. Note that the ranking process, while also involved in DPO data generation, is not shown in the Generate stage for simplicity.

lenges of communication efficiency and collective optimization, our work underscores the importance of developing advanced training frameworks for LLM-based MAS and highlights efficiency as a crucial metric to consider. We believe OPTIMA provides a solid foundation for future investigations into scaling and improving MAS and even general LLM systems.

## 2 OPTIMA: OPTIMIZING MULTI-AGENT LLMs VIA ITERATIVE TRAINING

### 2.1 OVERVIEW

OPTIMA is built upon an iterative *generate*, *rank*, *select*, and *train* paradigm. This approach allows for the progressive improvement of LLM-based agents in multi-agent settings, focusing on enhancing both the efficiency of inter-agent communication and the effectiveness of task completion.

Let  $\mathcal{M}_{\text{base}}$  denote the base LLM,  $\mathcal{D}$  the task dataset, and  $f$  the iterative training function. The iterative process can be formalized as  $\mathcal{M}_{t+1} = f(\mathcal{M}_t, \mathcal{D})$ , where  $\mathcal{M}_t$  represents the model at iteration  $t$ . The function  $f$  encapsulates the entire process of data generation, ranking, selection and model training. For each task instance  $d_i \in \mathcal{D}$ , we sample a set of  $N$  conversation trajectories  $\{\tau_i^j\}_{j=1}^N \subset \mathcal{T}$  using the agents powered by current model  $\mathcal{M}_t$ . Each trajectory  $\tau_i^j$  is then evaluated using a reward function  $R : \mathcal{T} \rightarrow \mathbb{R}$ , defined as:

$$R(\tau_i^j) = R_{\text{task}}(\tau_i^j) - \lambda_{\text{token}} R_{\text{token}}(\tau_i^j) + \lambda_{\text{loss}} \frac{1}{R_{\text{loss}}(\tau_i^j)}. \quad (1)$$

Here,  $R_{\text{task}} : \mathcal{T} \rightarrow \mathbb{R}$  is the task-specific performance metric,  $R_{\text{token}}(\tau_i^j) = \frac{\#\text{Tokens}(\tau_i^j)}{\max_k(\{\#\text{Tokens}(\tau_i^k)\}_k)}$  is the normalized token count, and  $R_{\text{loss}}(\tau_i^j) = g(\mathcal{L}(\mathcal{M}_{\text{base}}, d_i, \tau_i^j))$  is based on the language modeling loss of the base model  $\mathcal{M}_{\text{base}}$ , which we detail in Appendix E.2. The positive coefficients  $\lambda_{\text{token}}$  and  $\lambda_{\text{loss}}$  are hyper-parameters. This reward function is designed to balance multiple objectives simultaneously:  $R_{\text{task}}$  ensures that the model improves on the intended task,  $R_{\text{token}}$  encourages communication efficiency by penalizing verbose exchanges, and  $R_{\text{loss}}$  regularizes language naturalness and readability by favoring trajectories that are probable under the base model. By incorporating these components, we aim to develop LLM-based MAS that are not only effective in their designated tasks but also efficient in their communication, while maintaining interpretability in their outputs, unlike the often incomprehensible communication in prior RL research (Lazaridou et al., 2017; Evtimova et al., 2018; Chaabouni et al., 2022).

Based on these rewards, we apply several data selection criteria to select a subset of high-quality sampled trajectories  $\{\tau_i^*\}$  for each task instance. These selected trajectories form the training data  $\mathcal{D}_i^*$  at iteration  $i$ . The model is then updated:  $\mathcal{M}_{t+1} = \text{Train}(\mathcal{M}_t, \mathcal{D}_i^*)$ . The Train function can be

**Algorithm 1** Iterative Supervised Fine-Tuning

---

**Input:** Initialized model  $\mathcal{M}_{\text{init}}$ , dataset  $\mathcal{D}$ , sample size  $N$ , reward threshold  $\theta_{\text{sft}}$ , max iterations  $T$   
**Output:** Optimized model  $\mathcal{M}_T$

```

1:  $\mathcal{M}_0 \leftarrow \text{Initialize}(\mathcal{M}_{\text{init}}, \mathcal{D})$  ▷ Algorithm 3
2: for  $t = 0$  to  $T - 1$  do
3:    $\mathcal{D}_t^* \leftarrow \emptyset$ 
4:   for each  $d_i \in \mathcal{D}$  do
5:      $\{\tau_i^j\}_{j=1}^N \leftarrow \text{AgentChat}(\mathcal{M}_t, d_i)$  ▷ Generate N trajectories
6:      $\tau_i^* \leftarrow \arg \max_j R(\tau_i^j)$  ▷ Select best trajectory
7:     if  $R(\tau_i^*) > \theta_{\text{sft}}$  then
8:        $\mathcal{D}_t^* \leftarrow \mathcal{D}_t^* \cup \{(d_i, \tau_i^*)\}$ 
9:     end if
10:  end for
11:   $\mathcal{D}_t^* \leftarrow \text{TopK}(\mathcal{D}_t^*, 0.7|\mathcal{D}_t^*|)$  ▷ Retain top 70% trajectories
12:   $\mathcal{M}_{t+1} \leftarrow \text{SFT}(\mathcal{M}_t, \mathcal{D}_t^*)$ 
13: end for
14: return  $\mathcal{M}_T$ 

```

---

instantiated with various training algorithms, such as SFT or DPO, which we will discuss in detail in the following subsections.

Fig. 2 provides a high-level overview of OPTIMA. The specific instantiations of the generation and training processes will be detailed in the following subsections. The ranking process, consistent across all instantiations, is defined by the reward function presented in Eq. (1).

## 2.2 INITIALIZATION: DIVERSIFYING AGENT COMMUNICATION

Before starting the iterative training process, we address a critical challenge in LLM-based MAS: agents often produce responses in a similar style across conversation trajectories, even with high-temperature sampling. This homogeneity limits the exploration of diverse communication strategies, potentially hindering the optimization toward more efficient and effective interactions. Following the observation from AutoForm (Chen et al., 2024d), where LLMs can be explicitly prompted to leverage different more concise formats to communicate or reason without much compromise in performance, we introduce an initialization step that promotes diversity in agent communication.

Our approach leverages a pool of format specification prompts,  $\mathcal{P} = \{p_1, p_2, \dots, p_K\}$ , where each  $p_k$  is a string specifying a particular response format (e.g., JSON, list, see Appendix F for concrete examples and creation process). For each task instance  $d_i \in \mathcal{D}$ , we generate  $N$  conversation trajectories, each with a randomly selected format specification appended to the input task:

$$\tau_i^j = \mathcal{M}_{\text{base}}(d_i \oplus p_{k_j}), \quad k_j \sim \text{Uniform}(1, K), \quad j = 1, \dots, N, \quad (2)$$

where  $\oplus$  denotes string concatenation. This process yields a diverse set of trajectories  $\{\tau_i^j\}_{j=1}^N$  for each  $d_i$ , varying in both content and structure.

We then evaluate these trajectories using the reward function defined in Eq. (1), for each  $d_i$ , we select the trajectory with the highest reward:  $\tau_i^* = \arg \max_j R(\tau_i^j)$ . Finally, we select top  $k$  trajectories that exceed a predefined performance threshold  $\theta_{\text{init}}$ , resulting in a high-quality dataset:

$$\mathcal{D}_0^* = \text{TopK}(\{(d_i, \tau_i^*) | R_{\text{task}}(\tau_i^*) > \theta_{\text{init}}, \forall d_i \in \mathcal{D}\}, 0.7|D|). \quad (3)$$

Crucially, we remove the format specification prompts from the selected trajectories, resulting in a dataset of diverse, high-quality conversations without explicit format instructions. Using this dataset, we fine-tune the base model and obtain  $\mathcal{M}_{\text{base}}$  to obtain  $\mathcal{M}_0 = \text{SFT}(\mathcal{M}_{\text{base}}, \mathcal{D}_0^*)$ , which serves as the starting point for OPTIMA, able to generate diverse communication patterns without explicit format prompting. We provide pseudo-code in Appendix B for better understanding. This initialization sets the stage for more effective exploration and optimization in the subsequent iterative training process.

## 2.3 FRAMEWORK INSTANTIATION 1: ITERATIVE SUPERVISED FINE-TUNING

We introduce iterative Supervised Fine-Tuning (iSFT) as our first instantiation of OPTIMA. At each iteration  $t$ , iSFT follows the same general procedure outlined in Algorithm 3, generating a

216 set of  $N$  conversation trajectories for each task training instance  $d_i \in \mathcal{D}$  using the current model  
 217  $\mathcal{M}_t^{\text{iSFT}}$ . However, unlike initialization, iSFT omits the format specification pool, as  $\mathcal{M}_0$  has already  
 218 internalized diverse communication strategies. Unlike recent research on iterative training (Gülçehre  
 219 et al., 2023; Aksitov et al., 2023), iSFT maintains a fixed reward threshold  $\theta_{\text{SFT}}$  across iterations for  
 220 data selection. After data generation, the model undergoes standard SFT. This process continues  
 221 until a maximum number of iterations is reached. For clarity, the pseudo-code for iSFT is provided  
 222 in Algorithm 1.

223 iSFT provides a straightforward yet effective approach to optimize LLM-based MAS, leveraging the  
 224 diverse communication patterns established during initialization while consistently improving task  
 225 performance and communication efficiency.  
 226

## 227 2.4 FRAMEWORK INSTANTIATION 2: ITERATIVE DIRECT PREFERENCE OPTIMIZATION

228  
 229 While iSFT provides a straightforward approach to optimizing LLM-based MAS, it may be lim-  
 230 ited by its reliance on a single *best* trajectory for each task instance. To address this, we explore  
 231 iterative Direct Preference Optimization (iDPO) (Rafailov et al., 2023; Pang et al., 2024), which  
 232 optimizes models using comparative preferences and has demonstrated success in LLM alignment.  
 233 Applying DPO in multi-agent settings, however, poses distinct challenges, particularly in generating  
 234 meaningful paired data that capture the complexities of agent interactions.

235 **Data Generation:** To overcome these challenges, we integrate MCTS with DPO data collection  
 236 for high-quality paired data generation in multi-agent settings. Our MCTS-based approach con-  
 237 ceptualizes the multi-agent conversation as a tree, where nodes represent conversational turns, and  
 238 edges represent continuations. This structure allows us to explore diverse interaction trajectories  
 239 systematically and select high-quality paired data for DPO training. The MCTS process begins at  
 240 the root node (initial task prompt) and proceeds as follows: **(1) Expansion:** We select a node to  
 241 expand based on the following criteria. We first exclude leaf nodes and the second-to-last level  
 242 nodes to avoid wasting computation on low-variance expansions, then exclude nodes with content  
 243 similar to previously expanded nodes, measured based on edit distance (see Appendix E.1). From  
 244 the remaining nodes, we select 10 nodes with the highest rewards and sample one using the softmax  
 245 distribution over their rewards. **(2) Simulation:** For each selected node, we expand 3 trajectories,  
 246 simulating the conversation to completion. **(3) Backpropagation:** Once a trajectory is completed  
 247 and rewarded with Eq. (1), we update the estimated rewards of all nodes in the trajectory with the  
 248 average rewards from their children. **(4) Iteration:** We repeat the above process 8 times, resulting  
 in 24 trajectories. More iterations could potentially lead to more diverse and better-quality data.

249 **Paired Data Construction:** To generate high-quality paired data for DPO training, we traverse each  
 250 MCTS tree and identify node pairs  $(n_i, n_j)$  that satisfy three conditions: (1) shared ancestry, (2) the  
 251 higher estimated reward of  $n_i$  and  $n_j$  exceeds the threshold  $\theta_{\text{dpo-filter}}$ , and (3) their reward difference  
 252 exceeds the threshold  $\theta_{\text{dpo-diff}}$ . We sort these pairs by the higher estimated reward, and select the  
 253 top 50% pairs as part of the final training set. We construct DPO training instances by using the  
 254 common conversation history as the prompt, with  $n_i$  and  $n_j$  serving as the chosen and rejected  
 255 responses according to their estimated rewards.

256 The iDPO process then proceeds iteratively, alternating between MCTS-based data generation and  
 257 model updates using DPO. The pseudo-code for our iDPO process is presented in Algorithm 2.  
 258

## 259 2.5 FRAMEWORK INSTANTIATION 3: HYBRID ITERATIVE TRAINING

260  
 261 Building upon the strengths of both iSFT and iDPO, we investigate a hybrid approach that interleaves  
 262 SFT and DPO in the iterative training process, termed as iSFT-DPO. This hybrid method aims  
 263 to leverage the simplicity and directness of SFT in capturing high-quality trajectories, while also  
 264 benefiting from the nuanced comparative learning facilitated by DPO. By alternating between these  
 265 two training paradigms, we hypothesize that the model can more effectively balance the exploration  
 266 of diverse communication strategies with the exploitation of known effective patterns.

267 In practice, we implement this hybrid approach by performing one iteration of iSFT followed by  
 268 one iteration of iDPO, and repeating this cycle throughout the training process. This interleaving  
 269 allows the model to first consolidate learning from the best observed trajectories through SFT, and  
 then refine its understanding through the comparative preferences provided by DPO.

**Algorithm 2** Iterative Direct Preference Optimization

---

**Input:** Initial model  $\mathcal{M}_{\text{init}}$ , dataset  $\mathcal{D}$ , max iterations  $T$   
**Output:** Optimized model  $\mathcal{M}_T$

- 1:  $\mathcal{M}_0 \leftarrow \text{Initialize}(\mathcal{M}_{\text{init}}, \mathcal{D})$  ▷ Algorithm 3
- 2: **for**  $t = 0$  to  $T - 1$  **do**
- 3:    $\mathcal{D}_t^{\text{DPO}} \leftarrow \emptyset$
- 4:   **for each**  $d_i \in \mathcal{D}$  **do**
- 5:      $\mathcal{D}_i^{\text{DPO}} \leftarrow \text{MCTSDataGeneration}(\mathcal{M}_t, d_i)$  ▷ Algorithm 5
- 6:      $\mathcal{D}_t^{\text{DPO}} \leftarrow \mathcal{D}_t^{\text{DPO}} \cup \mathcal{D}_i^{\text{DPO}}$
- 7:   **end for**
- 8:    $\mathcal{M}_{t+1} \leftarrow \text{DPO}(\mathcal{M}_t, \mathcal{D}_t^{\text{DPO}})$
- 9: **end for**
- 10: **return**  $\mathcal{M}_T$

---

### 3 EXPERIMENTS

**Datasets.** We evaluate OPTIMA on two multi-agent settings: information exchange (IE) and debate. For IE, we use HotpotQA (Yang et al., 2018), 2WikiMultiHopQA (2WMHQA) (Ho et al., 2020), TriviaQA (Joshi et al., 2017), and CBT (Hill et al., 2016). For multi-hop datasets (HotpotQA, 2WikiMultiHopQA), we split relevant contexts between two agents, ensuring the answer can only be deduced from information exchange. For TriviaQA and CBT, contexts are randomly assigned, challenging agents to identify and communicate the relevant information effectively. The debate setting employs GSM8K (Cobbe et al., 2021), MATH (Hendrycks et al., 2021b), ARC’s challenge set (ARC-C) (Bhakhavatsalam et al., 2021) and MMLU (Hendrycks et al., 2021a), with one agent as solver and another as critic (Chen et al., 2024b). We use 0-shot for all benchmarks.

**Metrics.** We report F1 score between generated answers and labels for IE tasks. For debate tasks, we employ exact match accuracy (GSM8k, ARC-C, MMLU) or Sympy-based (Meurer et al., 2017) equivalence checking (MATH), following Lewkowycz et al. (2022). Conversations conclude when agents both mark the same answer with specified special tokens or reach a turn limit.

**Baselines.** We compare against single-agent approaches: Chain-of-Thought (CoT) (Wei et al., 2022) and Self-Consistency (SC) with majority voting (Wang et al., 2023) on  $n = 8$  samples. Given that the generated responses for IE tasks are in free form, direct adaptation to majority voting is impractical. Therefore, we first compute the pairwise F1 score among the sampled answers, grouping those with a pairwise F1 score exceeding 0.9, and report the average F1 score against the label for all the answers in the largest grouping. In the multi-agent context, we compare against Multi-Agent Debate (MAD) from Du et al. (2024) and AutoForm (Chen et al., 2024d). MAD utilizes natural language for inter-agent communication, providing a baseline for common multi-agent dialogue, while AutoForm encourages agents to leverage concise, non-natural-language formats to achieve a better performance-cost ratio, offering a comparison point for efficiency-oriented MAS.

**Training Setups.** We use Llama 3 8B (Meta, 2024) as our base model across all benchmarks. Our experiments focus on two-agent scenarios without external tools, a design choice that allows us to isolate and analyze the core aspects of multi-agent communication and collaboration. By constraining our initial investigation to these fundamental settings, we can more clearly demonstrate the efficacy of OPTIMA in optimizing inter-agent communication and task performance. This approach also provides a strong baseline for future research exploring more complex scenarios with multiple agents and tool use. Besides, we train a single model for both agents, although training separate models might yield improved performance, we leave it for future exploration. Detailed training configurations and prompts are provided in Appendices E and F.

#### 3.1 BENCHMARK RESULTS

Table 1 showcases OPTIMA’s performance across a diverse set of tasks, revealing consistent improvements over baseline methods in both effectiveness and efficiency. In IE tasks, OPTIMA variants demonstrate substantial gains, particularly in multi-hop reasoning scenarios like HotpotQA and 2WMHQA. Here, iSFT-DPO achieves peak performance while significantly reducing token usage

Table 1: **Performance and inference token number comparison across information exchange and debate tasks.** Best results are indicated in **bold**, and second-best results are underlined for all rows except the last three. The last three rows display self-consistency results for OPTIMA variants, with the best results highlighted in green. OPTIMA variants consistently outperform baselines in task performance and/or token efficiency.

| Method             | Information Exchange |             |             |             |             |             |             |             | Debate      |              |             |              |             |             |             |             |
|--------------------|----------------------|-------------|-------------|-------------|-------------|-------------|-------------|-------------|-------------|--------------|-------------|--------------|-------------|-------------|-------------|-------------|
|                    | HotpotQA             |             | 2WMHQA      |             | TriviaQA    |             | CBT         |             | MATH        |              | GSM8k       |              | ARC-C       |             | MMLU        |             |
|                    | F1                   | #Tok        | F1          | #Tok        | F1          | #Tok        | F1          | #Tok        | Acc         | #Tok         | Acc         | #Tok         | Acc         | #Tok        | Acc         | #Tok        |
| CoT                | 25.6                 | 123.7       | 20.5        | 139.8       | 59.8        | 110.3       | 43.4        | 135.3       | 23.9        | 329.8        | 71.5        | 230.9        | 65.2        | 138.9       | 46.0        | 132.2       |
| SC ( $n = 8$ )     | 33.8                 | 996.3       | 28.7        | 1052.8      | 70.0        | 891.4       | 52.9        | 1067.7      | <b>35.7</b> | 2600.9       | <u>80.3</u> | 1828.7       | <u>75.6</u> | 1116.7      | 54.0        | 1056.1      |
| MAD                | 28.4                 | 570.9       | 25.9        | 543.7       | 71.0        | 408.6       | 53.8        | 493.0       | 29.8        | 1517.6       | 72.5        | 514.7        | 71.4        | 478.0       | 51.5        | 516.7       |
| AutoForm           | 28.2                 | 97.7        | 24.7        | 117.7       | 60.9        | 74.0        | 35.0        | 64.8        | 26.1        | 644.3        | 71.0        | 410.5        | 60.2        | 221.2       | 43.8        | 198.5       |
| OPTIMA-iSFT        | <u>54.5</u>          | 67.6        | <u>72.4</u> | 61.2        | <u>71.9</u> | <u>51.5</u> | <b>71.8</b> | <u>38.5</u> | 30.1        | 830.3        | 79.5        | 311.5        | 74.1        | <u>92.2</u> | 56.8        | 123.8       |
| OPTIMA-iDPO        | 52.5                 | <b>45.7</b> | 66.1        | <b>35.9</b> | 69.3        | 69.2        | 66.7        | <b>37.2</b> | <u>30.4</u> | <b>272.8</b> | 78.5        | <u>270.1</u> | 74.5        | 97.8        | <u>59.6</u> | <u>61.6</u> |
| OPTIMA-iSFT-DPO    | <b>55.6</b>          | <u>63.3</u> | <b>74.2</b> | <u>54.9</u> | <b>77.1</b> | <b>32.5</b> | <u>70.1</u> | 38.9        | 29.3        | <u>488.1</u> | <b>80.4</b> | <b>246.5</b> | <b>77.1</b> | <b>88.0</b> | <b>60.2</b> | <b>56.7</b> |
| OPTIMA-iSFT SC     | 54.8                 | 806.2       | 72.6        | 245.6       | 73.7        | 413.8       | <u>72.2</u> | 847.4       | 32.4        | 2432.9       | 83.1        | 1750.7       | 77.2        | 1148.7      | 60.2        | 874.5       |
| OPTIMA-iDPO SC     | 52.8                 | 412.8       | 67.2        | 1056.2      | 71.8        | 702.8       | 66.8        | 520.6       | <u>36.9</u> | 2743.1       | <u>84.4</u> | 1750.8       | 77.0        | 1091.2      | 59.9        | 1050.4      |
| OPTIMA-iSFT-DPO SC | <u>57.4</u>          | 957.9       | <u>76.7</u> | 1096.0      | <u>77.5</u> | 494.1       | 71.8        | 417.8       | 34.8        | 2788.5       | 84.0        | 1748.7       | <u>78.8</u> | 1036.1      | <u>61.2</u> | 1026.7      |

compared to the strongest baseline SC. Notably, on 2WMHQA, iSFT-DPO improves F1 score by **38.3%** (2.8x improvement) while using only **10%** of the tokens required by MAD. This trend extends to other information exchange tasks, where OPTIMA variants maintain high performance with drastically lower token counts. The debate tasks present a more nuanced picture, yet OPTIMA’s benefits remain evident. Better task performance and token efficiency are still observed in ARC-C and MMLU, but for the MATH and GSM8k tasks, OPTIMA variants show comparable or slightly lower performance than SC, but still with much higher token efficiency. We conjecture this is due to the task’s difficulty and the small size of their training set. However, as we will demonstrate in Section 3.2, OPTIMA models trained on MATH transfer effectively to GSM8k, achieving performance nearly equivalent to models trained directly on GSM8k, with high token efficiency. More interestingly, Section 3.3 will show that applying SC to OPTIMA variants trained on MATH or GSM8k leads to better inference scaling laws on GSM8k compared to CoT SC.

A closer look at OPTIMA variants reveals interesting trade-offs. OPTIMA-iSFT often prioritizes performance at the expense of token efficiency, demonstrating the poorest efficiency in 5 of 8 tasks. In contrast, OPTIMA-iDPO often achieves remarkable reductions in token usage, occasionally with performance trade-offs. OPTIMA-iSFT-DPO emerges as a robust compromise, frequently delivering top-tier performance with satisfying token efficiency.

### 3.2 HOW WELL DOES OPTIMA GENERALIZE TO OOD TASKS?

To assess OPTIMA’s ability to generalize, we conducted transfer learning experiments across different task domains. We transferred models trained on HotpotQA to TriviaQA and 2WMHQA, as well as transferring from MATH to GSM8k. While these datasets share broad categories (question-answering and mathematical reasoning, respectively), they present different challenges in terms of complexity and required skills. The results, presented in Table 2, demonstrate OPTIMA’s robust transferability across these diverse tasks. In the question-answering domain, all OPTIMA variants significantly outperform baseline multi-agent methods on both OOD datasets. On 2WMHQA, the transferred iSFT more than doubles MAD’s F1 score while using only 14.6% of the tokens. Similar trends are observed in TriviaQA. When transferring from MATH to GSM8k, OPTIMA variants, particular iDPO, not only outperform the baselines on GSM8k but also achieve results comparable to models directly trained on GSM8k with even higher token efficiency (refer to Table 1 for comparison).

These results underscore OPTIMA’s potential for developing adaptable MAS, demonstrating that OPTIMA-trained models learn transferable skills for efficient information exchange and collabora-

Table 2: **Transfer performance of OPTIMA.** We transfer OPTIMA from Hotpot QA to 2WMH QA and Trivia QA, and from MATH to GSM8k, with MAD and AutoForm on each target task as baselines.

| Method   | 2WMH QA     |             | Trivia QA   |             | GSM8k       |              |
|----------|-------------|-------------|-------------|-------------|-------------|--------------|
|          | F1          | #Tok        | F1          | #Tok        | Acc         | #Tok         |
| MAD      | 25.9        | 543.7       | 71.0        | 408.9       | 72.5        | 514.7        |
| AutoForm | 24.7        | 117.7       | 60.9        | 74.0        | 71.0        | 410.5        |
| iSFT     | <b>56.5</b> | 79.6        | 70.0        | 90.2        | 74.6        | 293.7        |
| iDPO     | 51.6        | 84.3        | 68.0        | <b>41.1</b> | <b>77.9</b> | <b>185.7</b> |
| iSFT-DPO | 54.5        | <b>70.4</b> | <b>72.0</b> | 67.8        | 74.2        | 363.1        |

378  
379  
380  
381  
382  
383  
384  
385  
386  
387  
388  
389  
390  
391  
392  
393  
394  
395  
396  
397  
398  
399  
400  
401  
402  
403  
404  
405  
406  
407  
408  
409  
410  
411  
412  
413  
414  
415  
416  
417  
418  
419  
420  
421  
422  
423  
424  
425  
426  
427  
428  
429  
430  
431

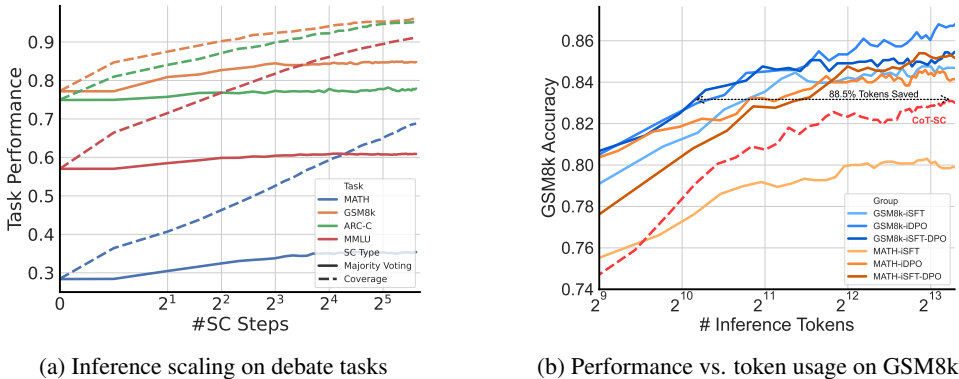


Figure 3: **OPTIMA’s impact on inference scaling laws.** (a) Relationship between OPTIMA variants’ self-consistency steps and performance on debate tasks. Solid lines represent majority voting accuracy, while dashed lines show coverage. (b) Performance of various models on GSM8k as a function of token usage, demonstrating OPTIMA’s efficiency gains.

tive reasoning. However, transferring to more distant domains remains challenging, e.g., we find it hard to transfer from HotpotQA to CBT, or from MATH to ARC-C. We believe it is a promising area for future research to explore if scaling OPTIMA to more generalized multi-task training could enhance the generalization of communication strategies in LLMs.

### 3.3 CAN OPTIMA LEAD TO BETTER INFERENCE SCALING LAW?

Recent research has highlighted the importance of inference scaling laws, which describe how model performance improves with increased compute during inference, typically by generating multiple samples per problem (Brown et al., 2024; Wu et al., 2024). While training scaling laws focus on the relationship between model size, dataset size, and performance, inference scaling laws explore the trade-off between inference compute budget and task accuracy. This paradigm offers a promising avenue for enhancing model capabilities without the need for further training models.

Fig. 3 illustrates OPTIMA’s impact on inference scaling laws. The left panel shows the relationship between the number of SC steps and performance on multi-agent debate tasks. We observe that while majority voting accuracy tends to plateau after a certain number of steps, the coverage, defined as the percentage of problems answered correctly at least once, continues to improve logarithmically with increased sampling. This trend aligns with findings in recent inference scaling law studies (Wu et al., 2024; Chen et al., 2024a) and suggests that more sophisticated answer selection techniques could further boost OPTIMA’s performance. We provide additional scaling law figures for all OPTIMA variants and on both IE and debate tasks in Appendix A, where similar trends can be observed.

The right panel of Fig. 3 demonstrates OPTIMA’s efficiency in improving inference scaling laws on the GSM8k task. OPTIMA variants, both those trained directly on GSM8k and those transferred from MATH, consistently outperform the CoT SC baseline except the iSFT variant transferred from MATH. Notably, iDPO trained on GSM8k achieves the performance of CoT-SC at around 10,000 tokens with 88.5% fewer tokens, effectively “*shifting the curve left*”. This significant reduction in token usage translates to substantial computational savings without sacrificing accuracy. Moreover, the MATH-trained OPTIMA variants, except iSFT, also deliver better inference scaling laws on GSM8k compared with CoT SC, underscoring the framework’s ability to generalize effectively across related tasks.

These results highlight OPTIMA’s potential to reshape inference scaling laws for LLM-based MAS and even general LLM systems. By enabling more efficient use of the inference compute budget, OPTIMA allows for better performance at lower computational costs or higher performance at the same cost. This efficiency gain opens new possibilities for leveraging advanced inference techniques like weighted voting or best-of-N selection (Wu et al., 2024), potentially leading to even greater performance improvements.

### 3.4 HOW DOES OPTIMA EVOLVE AGENT COMMUNICATION AND PERFORMANCE?



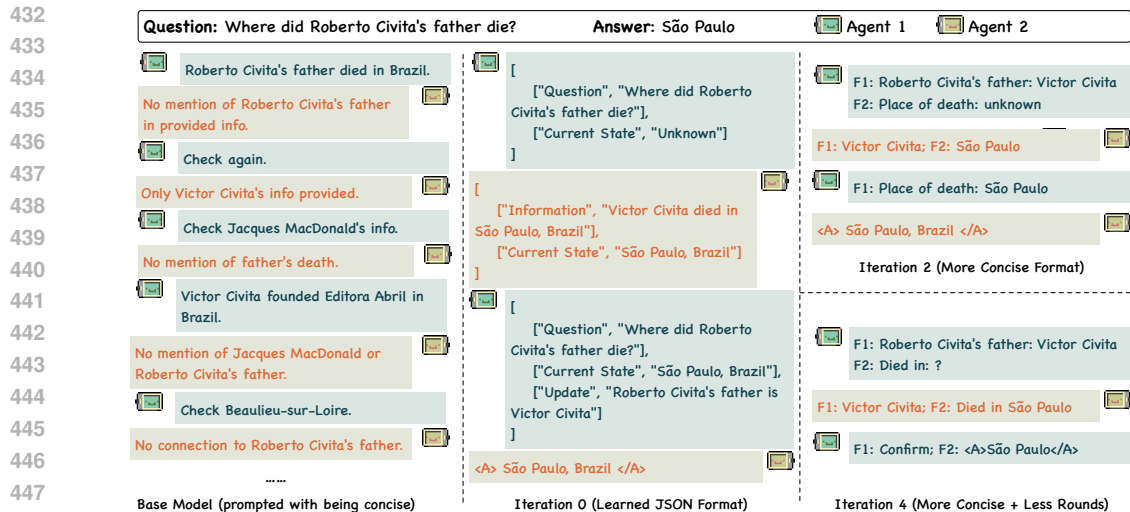


Figure 4: **Case study: Evolution of agent communication in OPTIMA-iSFT across iterations on 2WMH QA.** The different contexts given to the two agents are omitted for brevity. The progression demonstrates increasing efficiency and task-oriented communication.

To understand the impact of different components in our reward function, we conducted an ablation study on two representative tasks: 2WMHQA for IE and ARC-C for debate. We examined the performance of OPTIMA variants by removing either the token count regularization (#Tokens) or the LM loss (Loss) from the reward function. The results aim to answer two key questions: (1) *How does token count regularization affect the efficiency-performance trade-off?* (2) *What is the role of language modeling loss in maintaining communication quality?* Our findings consistently demonstrate the crucial role of each reward component in balancing task performance, communication efficiency, and language quality.

Table 3 presents the results of our ablation study. Removing the token count led to a substantial increase in the number of generated tokens across settings, with a particularly pronounced effect in the debate task. While this increased verbosity occasionally resulted in marginal performance improvements, it came at a significant computational cost. Conversely, eliminating the LM loss resulted in a decrease in token usage, often producing the most concise outputs among all variants. Examples comparing communication with and without LM loss can be found in Appendix C. Without LM loss, the model often generated overly concise messages containing insufficient information and was prone to hallucination, potentially explaining the inferior performance under this condition. These results underscore that effective LLM-based MAS should optimize not only for task performance but also for the efficiency and quality of inter-agent dialogue. The design of OPTIMA’s reward function enables this holistic optimization, leading to more effective and efficient multi-agent collaboration while highlighting the delicate balance required in optimizing such systems.

### 3.5 HOW AGENT COMMUNICATION EVOLVES OVER OPTIMIZATION ITERATIONS?

Fig. 1 illustrates the performance gains and token efficiency of OPTIMA variants across the optimization iterations, revealing a distinctive two-phase optimization pattern. In the initial phase (iterations 0-1), we observe a substantial improvement in task performance for all OPTIMA variants, accompanied by a clear increase in token usage. This suggests that OPTIMA initially prioritizes effectiveness, allowing agents to develop sophisticated problem-solving strategies through expanded communication.

Table 3: Ablation study on reward components for OPTIMA variants on two representative tasks.

| Setting     | 2WMH QA                       |                               | ARC-C                         |                               |
|-------------|-------------------------------|-------------------------------|-------------------------------|-------------------------------|
|             | F1                            | #Tok                          | Acc                           | #Tok                          |
| iSFT        | <b>72.4</b>                   | 61.2                          | 74.1                          | 92.2                          |
| w/o #Tokens | <b>72.4</b> <sub>(0.0)</sub>  | 290.3 <sub>(4.8x)</sub>       | <b>74.2</b> <sub>(+0.1)</sub> | 579.6 <sub>(6.3x)</sub>       |
| w/o Loss    | 69.7 <sub>(-2.7)</sub>        | <b>45.4</b> <sub>(0.7x)</sub> | 72.6 <sub>(-1.5)</sub>        | <b>69.7</b> <sub>(0.8x)</sub> |
| iDPO        | 66.1                          | <b>35.9</b>                   | 74.5                          | 97.8                          |
| w/o #Tokens | <b>72.9</b> <sub>(+6.8)</sub> | 183.3 <sub>(5.1x)</sub>       | <b>75.5</b> <sub>(+1.0)</sub> | 266.0 <sub>(2.7x)</sub>       |
| w/o Loss    | 63.0 <sub>(-3.1)</sub>        | 54.6 <sub>(1.5x)</sub>        | 74.4 <sub>(-0.1)</sub>        | <b>81.2</b> <sub>(0.8x)</sub> |
| iSFT-DPO    | <b>74.2</b>                   | 54.9                          | <b>77.1</b>                   | 88.0                          |
| w/o #Tokens | 63.5 <sub>(-10.7)</sub>       | 219.7 <sub>(4.0x)</sub>       | 76.9 <sub>(-0.2)</sub>        | 354.8 <sub>(4.0x)</sub>       |
| w/o Loss    | 66.7 <sub>(-7.5)</sub>        | <b>38.1</b> <sub>(0.7x)</sub> | 76.3 <sub>(-0.8)</sub>        | <b>63.4</b> <sub>(0.7x)</sub> |

486 tion. The subsequent iterations demonstrate OPTIMA’s ability to refine these strategies for efficiency  
487 without compromising performance. We observe a gradual but consistent decrease in token usage  
488 across all variants, coupled with continued performance improvements.

489 To provide concrete examples of how OPTIMA shapes agent communication, we present a case from  
490 iSFT on an information exchange task in Fig. 4. The base model exhibits unfocused and repetitive  
491 exchanges, failing to efficiently address the task at hand. At iteration 0, while more structured,  
492 the exchange is verbose and includes unnecessary metadata. By iteration 2, we observe a marked  
493 shift towards concise, task-oriented communication, with agents adopting a streamlined format that  
494 efficiently conveys key information. The final iteration demonstrates further refinement, with agents  
495 maintaining the efficient structure while eliminating any residual verbosity. This progression aligns  
496 with our quantitative findings, showcasing OPTIMA’s ability to form communication patterns that  
497 are both highly effective and remarkably efficient.

## 498 499 500 4 RELATED WORK

501  
502 **LLM-Based MAS.** LLM-based MAS have emerged as a powerful paradigm for addressing complex  
503 tasks across various domains. Seminal works by Liang et al. (2023) and Du et al. (2024) demon-  
504 strated the potential of LLM-powered agents in collaborative problem-solving through multi-agent  
505 debate. This foundation has sparked diverse research directions, including role-playing for com-  
506 plex reasoning (Wang et al., 2024b; Chen et al., 2024b), collaborative software development (Qian  
507 et al., 2024c; Hong et al., 2024; Ishibashi & Nishimura, 2024), and embodied agent interactions  
508 (Zhang et al., 2024; Mandi et al., 2024; Guo et al., 2024). Recent studies have shown that increasing  
509 the number and diversity of agents can lead to performance gains in MAS (Wang et al., 2024a; Li  
510 et al., 2024a; Chen et al., 2024c). However, as LLM-based MAS grow in scale and complexity,  
511 challenges related to computational costs and communication efficiency become more pronounced  
512 (Chen et al., 2024d; Li et al., 2024b). Notably, there is a lack of systematic training algorithms  
513 specifically designed to optimize both the effectiveness and efficiency of LLM-based multi-agent  
514 systems, with most existing approaches relying on updating agent memory (Qian et al., 2024a; Gao  
515 et al., 2024). Our work addresses this gap by introducing a training framework that simultaneously  
516 enhances communication efficiency and task effectiveness in LLM-based MAS.

517 **Iterative Refinement of LLMs.** The pursuit of continual improvement in LLMs has led to the  
518 development of various iterative refinement paradigms. While self-reflection mechanisms like Re-  
519 flexion (Shinn et al., 2023) and self-refine (Madaan et al., 2023) show promise, they heavily rely  
520 on LLMs’ limited self-correction abilities, which is relatively weak for most of the current LLMs  
521 (Huang et al., 2024; Olausson et al., 2024; Kamoi et al., 2024). More robust approaches focus on  
522 iterative parameter updates, for example, ReST (Gülçehre et al., 2023), ReST<sup>EM</sup> (Singh et al., 2024)  
523 and STaR (Zelikman et al., 2022) train models on self-generated high-quality reasoning paths. Pang  
524 et al. (2024) further integrate the incorrect self-generated paths and train models with DPO. The  
525 extension to complex, multi-step tasks (Aksitov et al., 2023) further demonstrates the versatility of  
526 these methods. However, iterative refinement remains largely unexplored in the context of LLM-  
527 based MAS. Our work addresses this gap by presenting the first effective training framework for  
528 iteratively optimizing LLMs in MAS contexts. By simultaneously enhancing communication effi-  
529 ciency and task effectiveness, our approach shows the potential of iterative training in MAS.

## 530 5 CONCLUSION

531  
532 We present OPTIMA, a novel framework for training LLM-based MAS that significantly improves  
533 communication efficiency and task performance. Extensive experiments across a range of tasks  
534 demonstrate OPTIMA’s consistent superiority over both single-agent and multi-agent baselines. The  
535 framework introduces key innovations such as iterative training techniques, a balanced reward func-  
536 tion, and an MCTS-inspired approach for data generation. OPTIMA also shows promise in enhanc-  
537 ing inference scaling laws and transferring knowledge to OOD tasks. These findings highlight the  
538 critical role of efficient communication in MAS and LLM systems. While OPTIMA marks a major  
539 step forward in multi-agent LLM training, further exploration into its scalability to larger models  
and more complex scenarios is a promising direction for future research.

## REFERENCES

- 540  
541  
542 Renat Aksitov, Sobhan Miryoosefi, Zonglin Li, Daliang Li, Sheila Babayan, Kavya Kopparapu,  
543 Zachary Fisher, Ruiqi Guo, Sushant Prakash, Pranesh Srinivasan, Manzil Zaheer, Felix X. Yu,  
544 and Sanjiv Kumar. Rest meets react: Self-improvement for multi-step reasoning LLM agent.  
545 *CoRR*, abs/2312.10003, 2023. doi: 10.48550/ARXIV.2312.10003. URL [https://doi.org/  
546 10.48550/arXiv.2312.10003](https://doi.org/10.48550/arXiv.2312.10003).
- 547 Anthropic. Claude 3.5 sonnet, 2024. URL [https://www.anthropic.com/news/  
548 claude-3-5-sonnet](https://www.anthropic.com/news/claude-3-5-sonnet).
- 549 Bowen Baker, Ingmar Kanitscheider, Todor M. Markov, Yi Wu, Glenn Powell, Bob McGrew, and  
550 Igor Mordatch. Emergent tool use from multi-agent autocurricula. In *8th International Con-  
551 ference on Learning Representations, ICLR 2020, Addis Ababa, Ethiopia, April 26-30, 2020*.  
552 OpenReview.net, 2020. URL <https://openreview.net/forum?id=SkxpxJBKwS>.
- 553  
554 Sumithra Bhakthavatsalam, Daniel Khashabi, Tushar Khot, Bhavana Dalvi Mishra, Kyle Richard-  
555 son, Ashish Sabharwal, Carissa Schoenick, Oyvind Tafjord, and Peter Clark. Think you have  
556 solved direct-answer question answering? try arc-da, the direct-answer AI2 reasoning challenge.  
557 *CoRR*, abs/2102.03315, 2021. URL <https://arxiv.org/abs/2102.03315>.
- 558 Bradley C. A. Brown, Jordan Juravsky, Ryan Saul Ehrlich, Ronald Clark, Quoc V. Le, Christopher  
559 Ré, and Azalia Mirhoseini. Large language monkeys: Scaling inference compute with repeated  
560 sampling. *CoRR*, abs/2407.21787, 2024. doi: 10.48550/ARXIV.2407.21787. URL <https://doi.org/10.48550/arXiv.2407.21787>.
- 561  
562 Rahma Chaabouni, Florian Strub, Florent Althé, Eugene Tarassov, Corentin Tallec, Elnaz Davoodi,  
563 Kory Wallace Mathewson, Olivier Tieleman, Angeliki Lazaridou, and Bilal Piot. Emergent com-  
564 munication at scale. In *The Tenth International Conference on Learning Representations, ICLR  
565 2022, Virtual Event, April 25-29, 2022*. OpenReview.net, 2022. URL [https://openreview.  
566 net/forum?id=AUGBfDIV9rL](https://openreview.net/forum?id=AUGBfDIV9rL).
- 567  
568 Lingjiao Chen, Jared Quincy Davis, Boris Hanin, Peter Bailis, Ion Stoica, Matei Zaharia, and James  
569 Zou. Are more LLM calls all you need? towards scaling laws of compound inference systems.  
570 *CoRR*, abs/2403.02419, 2024a. doi: 10.48550/ARXIV.2403.02419. URL [https://doi.org/  
571 10.48550/arXiv.2403.02419](https://doi.org/10.48550/arXiv.2403.02419).
- 572 Weize Chen, Yusheng Su, Jingwei Zuo, Cheng Yang, Chenfei Yuan, Chi-Min Chan, Heyang Yu,  
573 Yaxi Lu, Yi-Hsin Hung, Chen Qian, Yujia Qin, Xin Cong, Ruobing Xie, Zhiyuan Liu, Maosong  
574 Sun, and Jie Zhou. Agentverse: Facilitating multi-agent collaboration and exploring emergent  
575 behaviors. In *The Twelfth International Conference on Learning Representations, ICLR 2024,  
576 Vienna, Austria, May 7-11, 2024*. OpenReview.net, 2024b. URL [https://openreview.  
577 net/forum?id=EHg5GDnyq1](https://openreview.net/forum?id=EHg5GDnyq1).
- 578 Weize Chen, Ziming You, Ran Li, Yitong Guan, Chen Qian, Chenyang Zhao, Cheng Yang, Ruobing  
579 Xie, Zhiyuan Liu, and Maosong Sun. Internet of agents: Weaving a web of heterogeneous agents  
580 for collaborative intelligence. *CoRR*, abs/2407.07061, 2024c. doi: 10.48550/ARXIV.2407.07061.  
581 URL <https://doi.org/10.48550/arXiv.2407.07061>.
- 582 Weize Chen, Chenfei Yuan, Jiarui Yuan, Yusheng Su, Chen Qian, Cheng Yang, Ruobing Xie,  
583 Zhiyuan Liu, and Maosong Sun. Beyond natural language: Llms leveraging alternative formats  
584 for enhanced reasoning and communication. *CoRR*, abs/2402.18439, 2024d. doi: 10.48550/  
585 ARXIV.2402.18439. URL <https://doi.org/10.48550/arXiv.2402.18439>.
- 586  
587 Karl Cobbe, Vineet Kosaraju, Mohammad Bavarian, Mark Chen, Heewoo Jun, Lukasz Kaiser,  
588 Matthias Plappert, Jerry Tworek, Jacob Hilton, Reiichiro Nakano, Christopher Hesse, and John  
589 Schulman. Training verifiers to solve math word problems. *CoRR*, abs/2110.14168, 2021. URL  
590 <https://arxiv.org/abs/2110.14168>.
- 591 Yilun Du, Shuang Li, Antonio Torralba, Joshua B. Tenenbaum, and Igor Mordatch. Improving  
592 factuality and reasoning in language models through multiagent debate. In *Forty-first Interna-  
593 tional Conference on Machine Learning, ICML 2024, Vienna, Austria, July 21-27, 2024*. Open-  
Review.net, 2024. URL <https://openreview.net/forum?id=zj7YuTE4t8>.

- 594 Yann Dubois, Balázs Galambosi, Percy Liang, and Tatsunori B. Hashimoto. Length-controlled  
595 alpacaeval: A simple way to debias automatic evaluators. *CoRR*, abs/2404.04475, 2024. doi: 10.  
596 48550/ARXIV.2404.04475. URL <https://doi.org/10.48550/arXiv.2404.04475>.  
597
- 598 Katrina Evtimova, Andrew Drozdov, Douwe Kiela, and Kyunghyun Cho. Emergent communica-  
599 tion in a multi-modal, multi-step referential game. In *6th International Conference on Learn-  
600 ing Representations, ICLR 2018, Vancouver, BC, Canada, April 30 - May 3, 2018, Conference  
601 Track Proceedings*. OpenReview.net, 2018. URL <https://openreview.net/forum?id=rjGZq6g0->.  
602
- 603 Shen Gao, Hao Li, Zhengliang Shi, Chengrui Huang, Quan Tu, Zhiliang Tian, Minlie Huang, and  
604 Shuo Shang. 360{\deg}rea: Towards A reusable experience accumulation with 360{\deg} as-  
605 sessment for multi-agent system. *CoRR*, abs/2404.05569, 2024. doi: 10.48550/ARXIV.2404.  
606 05569. URL <https://doi.org/10.48550/arXiv.2404.05569>.  
607
- 608 Çağlar Gülçehre, Tom Le Paine, Srivatsan Srinivasan, Ksenia Konyushkova, Lotte Weerts, Abhishek  
609 Sharma, Aditya Siddhant, Alex Ahern, Miaosen Wang, Chenjie Gu, Wolfgang Macherey, Arnaud  
610 Doucet, Orhan Firat, and Nando de Freitas. Reinforced self-training (rest) for language modeling.  
611 *CoRR*, abs/2308.08998, 2023. doi: 10.48550/ARXIV.2308.08998. URL [https://doi.org/  
612 10.48550/arXiv.2308.08998](https://doi.org/10.48550/arXiv.2308.08998).  
613
- 614 Xudong Guo, Kaixuan Huang, Jiale Liu, Wenhui Fan, Natalia Vélez, Qingyun Wu, Huazheng Wang,  
615 Thomas L. Griffiths, and Mengdi Wang. Embodied LLM agents learn to cooperate in organized  
616 teams. *CoRR*, abs/2403.12482, 2024. doi: 10.48550/ARXIV.2403.12482. URL [https://  
617 doi.org/10.48550/arXiv.2403.12482](https://doi.org/10.48550/arXiv.2403.12482).  
618
- 619 Dan Hendrycks, Collin Burns, Steven Basart, Andy Zou, Mantas Mazeika, Dawn Song, and Jacob  
620 Steinhardt. Measuring massive multitask language understanding. In *9th International Confer-  
621 ence on Learning Representations, ICLR 2021, Virtual Event, Austria, May 3-7, 2021*. OpenRe-  
622 view.net, 2021a. URL <https://openreview.net/forum?id=d7KBjmI3GmQ>.  
623
- 624 Dan Hendrycks, Collin Burns, Saurav Kadavath, Akul Arora, Steven Basart, Eric Tang,  
625 Dawn Song, and Jacob Steinhardt. Measuring mathematical problem solving with  
626 the MATH dataset. In Joaquin Vanschoren and Sai-Kit Yeung (eds.), *Proceedings  
627 of the Neural Information Processing Systems Track on Datasets and Benchmarks  
628 I, NeurIPS Datasets and Benchmarks 2021, December 2021, virtual, 2021b*. URL  
629 [https://datasets-benchmarks-proceedings.neurips.cc/paper/2021/  
hash/be83ab3ecd0db773eb2dc1b0a17836a1-Abstract-round2.html](https://datasets-benchmarks-proceedings.neurips.cc/paper/2021/hash/be83ab3ecd0db773eb2dc1b0a17836a1-Abstract-round2.html).
- 630 Felix Hill, Antoine Bordes, Sumit Chopra, and Jason Weston. The goldilocks principle: Read-  
631 ing children’s books with explicit memory representations. In Yoshua Bengio and Yann LeCun  
632 (eds.), *4th International Conference on Learning Representations, ICLR 2016, San Juan, Puerto  
633 Rico, May 2-4, 2016, Conference Track Proceedings*, 2016. URL [http://arxiv.org/abs/  
634 1511.02301](http://arxiv.org/abs/1511.02301).  
635
- 636 Xanh Ho, Anh-Khoa Duong Nguyen, Saku Sugawara, and Akiko Aizawa. Constructing A multi-  
637 hop QA dataset for comprehensive evaluation of reasoning steps. In Donia Scott, Núria Bel,  
638 and Chengqing Zong (eds.), *Proceedings of the 28th International Conference on Computa-  
639 tional Linguistics, COLING 2020, Barcelona, Spain (Online), December 8-13, 2020*, pp. 6609–  
640 6625. International Committee on Computational Linguistics, 2020. doi: 10.18653/V1/2020.  
641 COLING-MAIN.580. URL [https://doi.org/10.18653/v1/2020.coling-main.  
642 580](https://doi.org/10.18653/v1/2020.coling-main.580).  
643
- 644 Sirui Hong, Mingchen Zhuge, Jonathan Chen, Xiawu Zheng, Yuheng Cheng, Jinlin Wang, Ceyao  
645 Zhang, Zili Wang, Steven Ka Shing Yau, Zijuan Lin, Liyang Zhou, Chenyu Ran, Lingfeng  
646 Xiao, Chenglin Wu, and Jürgen Schmidhuber. Metagpt: Meta programming for A multi-agent  
647 collaborative framework. In *The Twelfth International Conference on Learning Representa-  
tions, ICLR 2024, Vienna, Austria, May 7-11, 2024*. OpenReview.net, 2024. URL <https://openreview.net/forum?id=VtmBAGCN7o>.

- 648 Jie Huang, Xinyun Chen, Swaroop Mishra, Huaixiu Steven Zheng, Adams Wei Yu, Xinying Song,  
649 and Denny Zhou. Large language models cannot self-correct reasoning yet. In *The Twelfth Inter-*  
650 *national Conference on Learning Representations, ICLR 2024, Vienna, Austria, May 7-11, 2024.*  
651 OpenReview.net, 2024. URL <https://openreview.net/forum?id=Ikmd3fKBPQ>.
- 652 Yoichi Ishibashi and Yoshimasa Nishimura. Self-organized agents: A LLM multi-agent framework  
653 toward ultra large-scale code generation and optimization. *CoRR*, abs/2404.02183, 2024. doi: 10.  
654 48550/ARXIV.2404.02183. URL <https://doi.org/10.48550/arXiv.2404.02183>.
- 656 Jeffrey D. Johnson, Jinghong Li, and Zengshi Chen. Reinforcement learning: An introduc-  
657 tion: R.S. Sutton, A.G. Barto, MIT Press, Cambridge, MA 1998, 322 pp. ISBN 0-262-19398-  
658 1. *Neurocomputing*, 35(1-4):205–206, 2000. doi: 10.1016/S0925-2312(00)00324-6. URL  
659 [https://doi.org/10.1016/S0925-2312\(00\)00324-6](https://doi.org/10.1016/S0925-2312(00)00324-6).
- 660 Mandar Joshi, Eunsol Choi, Daniel S. Weld, and Luke Zettlemoyer. Triviaqa: A large scale distantly  
661 supervised challenge dataset for reading comprehension. In Regina Barzilay and Min-Yen Kan  
662 (eds.), *Proceedings of the 55th Annual Meeting of the Association for Computational Linguistics,*  
663 *ACL 2017, Vancouver, Canada, July 30 - August 4, Volume 1: Long Papers*, pp. 1601–1611.  
664 Association for Computational Linguistics, 2017. doi: 10.18653/v1/P17-1147. URL <https://doi.org/10.18653/v1/P17-1147>.
- 666 Ryo Kamoi, Yusen Zhang, Nan Zhang, Jiawei Han, and Rui Zhang. When can llms actually correct  
667 their own mistakes? A critical survey of self-correction of llms. *CoRR*, abs/2406.01297, 2024.  
668 doi: 10.48550/ARXIV.2406.01297. URL [https://doi.org/10.48550/arXiv.2406.](https://doi.org/10.48550/arXiv.2406.01297)  
669 [01297](https://doi.org/10.48550/arXiv.2406.01297).
- 670 Marc Lanctot, Vinícius Flores Zambaldi, Audrunas Gruslys, Angeliki Lazaridou, Karl Tuyls,  
671 Julien Pérolat, David Silver, and Thore Graepel. A unified game-theoretic approach to  
672 multiagent reinforcement learning. In Isabelle Guyon, Ulrike von Luxburg, Samy Ben-  
673 gio, Hanna M. Wallach, Rob Fergus, S. V. N. Vishwanathan, and Roman Garnett (eds.),  
674 *Advances in Neural Information Processing Systems 30: Annual Conference on Neural*  
675 *Information Processing Systems 2017, December 4-9, 2017, Long Beach, CA, USA*, pp.  
676 4190–4203, 2017. URL [https://proceedings.neurips.cc/paper/2017/hash/](https://proceedings.neurips.cc/paper/2017/hash/3323fe11e9595c09af38fe67567a9394-Abstract.html)  
677 [3323fe11e9595c09af38fe67567a9394-Abstract.html](https://proceedings.neurips.cc/paper/2017/hash/3323fe11e9595c09af38fe67567a9394-Abstract.html).
- 679 Angeliki Lazaridou, Alexander Peysakhovich, and Marco Baroni. Multi-agent cooperation and the  
680 emergence of (natural) language. In *5th International Conference on Learning Representations,*  
681 *ICLR 2017, Toulon, France, April 24-26, 2017, Conference Track Proceedings.* OpenReview.net,  
682 2017. URL <https://openreview.net/forum?id=Hk8N3ScIlg>.
- 683 Aitor Lewkowycz, Anders Andreassen, David Dohan, Ethan Dyer, Henryk Michalewski, Vinay V.  
684 Ramasesh, Ambrose Slone, Cem Anil, Imanol Schlag, Theo Gutman-Solo, Yuhuai Wu, Behnam  
685 Neyshabur, Guy Gur-Ari, and Vedant Misra. Solving quantitative reasoning problems with lan-  
686 guage models. In Sanmi Koyejo, S. Mohamed, A. Agarwal, Danielle Belgrave, K. Cho, and A. Oh  
687 (eds.), *Advances in Neural Information Processing Systems 35: Annual Conference on Neural In-*  
688 *formation Processing Systems 2022, NeurIPS 2022, New Orleans, LA, USA, November 28 - De-*  
689 *cember 9, 2022*, 2022. URL [http://papers.nips.cc/paper\\_files/paper/2022/](http://papers.nips.cc/paper_files/paper/2022/hash/18abbef8cfe9203fdf9053c9c4fe191-Abstract-Conference.html)  
690 [hash/18abbef8cfe9203fdf9053c9c4fe191-Abstract-Conference.html](http://papers.nips.cc/paper_files/paper/2022/hash/18abbef8cfe9203fdf9053c9c4fe191-Abstract-Conference.html).
- 691 Junyou Li, Qin Zhang, Yangbin Yu, Qiang Fu, and Deheng Ye. More agents is all you need. *CoRR*,  
692 abs/2402.05120, 2024a. doi: 10.48550/ARXIV.2402.05120. URL [https://doi.org/10.](https://doi.org/10.48550/arXiv.2402.05120)  
693 [48550/arXiv.2402.05120](https://doi.org/10.48550/arXiv.2402.05120).
- 694 Yunxuan Li, Yibing Du, Jiageng Zhang, Le Hou, Peter Grabowski, Yeqing Li, and Eugene Ie. Im-  
695 proving multi-agent debate with sparse communication topology. *CoRR*, abs/2406.11776, 2024b.  
696 doi: 10.48550/ARXIV.2406.11776. URL [https://doi.org/10.48550/arXiv.2406.](https://doi.org/10.48550/arXiv.2406.11776)  
697 [11776](https://doi.org/10.48550/arXiv.2406.11776).
- 698 Tian Liang, Zhiwei He, Wenxiang Jiao, Xing Wang, Yan Wang, Rui Wang, Yujiu Yang, Zhaopeng  
699 Tu, and Shuming Shi. Encouraging divergent thinking in large language models through multi-  
700 agent debate. *CoRR*, abs/2305.19118, 2023. doi: 10.48550/ARXIV.2305.19118. URL [https://doi.org/10.](https://doi.org/10.48550/arXiv.2305.19118)  
701 [48550/arXiv.2305.19118](https://doi.org/10.48550/arXiv.2305.19118).

- 702 Wei Liu, Chenxi Wang, Yifei Wang, Zihao Xie, Rennai Qiu, Yufan Dang, Zhuoyun Du, Weize  
703 Chen, Cheng Yang, and Chen Qian. Autonomous agents for collaborative task under information  
704 asymmetry. *CoRR*, abs/2406.14928, 2024. doi: 10.48550/ARXIV.2406.14928. URL <https://doi.org/10.48550/arXiv.2406.14928>.  
705
- 706 Aman Madaan, Niket Tandon, Prakhar Gupta, Skyler Hallinan, Luyu Gao, Sarah Wieg-  
707 effe, Uri Alon, Nouha Dziri, Shrimai Prabhumoye, Yiming Yang, Shashank Gupta, Bod-  
708 hisattwa Prasad Majumder, Katherine Hermann, Sean Welleck, Amir Yazdanbakhsh, and  
709 Peter Clark. Self-refine: Iterative refinement with self-feedback. In Alice Oh, Tris-  
710 tan Naumann, Amir Globerson, Kate Saenko, Moritz Hardt, and Sergey Levine (eds.), *Ad-  
711 vances in Neural Information Processing Systems 36: Annual Conference on Neural Infor-  
712 mation Processing Systems 2023, NeurIPS 2023, New Orleans, LA, USA, December 10 - 16,  
713 2023*, 2023. URL [http://papers.nips.cc/paper\\_files/paper/2023/hash/  
714 91edff07232fb1b55a505a9e9f6c0ff3-Abstract-Conference.html](http://papers.nips.cc/paper_files/paper/2023/hash/91edff07232fb1b55a505a9e9f6c0ff3-Abstract-Conference.html).
- 715 Zhao Mandi, Shreya Jain, and Shuran Song. Roco: Dialectic multi-robot collaboration with large  
716 language models. In *IEEE International Conference on Robotics and Automation, ICRA 2024,  
717 Yokohama, Japan, May 13-17, 2024*, pp. 286–299. IEEE, 2024. doi: 10.1109/ICRA57147.2024.  
718 10610855. URL <https://doi.org/10.1109/ICRA57147.2024.10610855>.  
719
- 720 Meta. Llama 3 model card. 2024. URL [https://github.com/meta-llama/llama3/  
721 blob/main/MODEL\\_CARD.md](https://github.com/meta-llama/llama3/blob/main/MODEL_CARD.md).
- 722 Aaron Meurer, Christopher P. Smith, Mateusz Paprocki, Ondrej Certík, Sergey B. Kirpichev,  
723 Matthew Rocklin, Amit Kumar, Sergiu Ivanov, Jason Keith Moore, Sartaj Singh, Thilina Rath-  
724 nayake, Sean Vig, Brian E. Granger, Richard P. Muller, Francesco Bonazzi, Harsh Gupta,  
725 Shivam Vats, Fredrik Johansson, Fabian Pedregosa, Matthew J. Curry, Andy R. Terrel, Stepán  
726 Roucka, Ashutosh Saboo, Isuru Fernando, Sumith Kulal, Robert Cimrman, and Anthony M.  
727 Scopatz. Sympy: symbolic computing in python. *PeerJ Comput. Sci.*, 3:e103, 2017. doi:  
728 10.7717/PEERJ-CS.103. URL <https://doi.org/10.7717/peerj-cs.103>.
- 729 Theo X. Olausson, Jeevana Priya Inala, Chenglong Wang, Jianfeng Gao, and Armando Solar-  
730 Lezama. Is self-repair a silver bullet for code generation? In *The Twelfth International Conference  
731 on Learning Representations, ICLR 2024, Vienna, Austria, May 7-11, 2024*. OpenReview.net,  
732 2024. URL <https://openreview.net/forum?id=y0GJXRungR>.
- 733 OpenAI. GPT-4 technical report. *CoRR*, abs/2303.08774, 2023. doi: 10.48550/ARXIV.2303.08774.  
734 URL <https://doi.org/10.48550/arXiv.2303.08774>.  
735
- 736 Richard Yuanzhe Pang, Weizhe Yuan, Kyunghyun Cho, He He, Sainbayar Sukhbaatar, and Jason  
737 Weston. Iterative reasoning preference optimization. *CoRR*, abs/2404.19733, 2024. doi: 10.  
738 48550/ARXIV.2404.19733. URL <https://doi.org/10.48550/arXiv.2404.19733>.  
739
- 740 Chen Qian, Yufan Dang, Jiahao Li, Wei Liu, Zihao Xie, Yifei Wang, Weize Chen, Cheng Yang,  
741 Xin Cong, Xiaoyin Che, Zhiyuan Liu, and Maosong Sun. Experiential co-learning of software-  
742 developing agents. In Lun-Wei Ku, Andre Martins, and Vivek Srikumar (eds.), *Proceedings of the  
743 62nd Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers),  
744 ACL 2024, Bangkok, Thailand, August 11-16, 2024*, pp. 5628–5640. Association for Computa-  
745 tional Linguistics, 2024a. URL <https://aclanthology.org/2024.acl-long.305>.
- 746 Chen Qian, Jiahao Li, Yufan Dang, Wei Liu, Yifei Wang, Zihao Xie, Weize Chen, Cheng Yang,  
747 Yingli Zhang, Zhiyuan Liu, and Maosong Sun. Iterative experience refinement of software-  
748 developing agents. *CoRR*, abs/2405.04219, 2024b. doi: 10.48550/ARXIV.2405.04219. URL  
749 <https://doi.org/10.48550/arXiv.2405.04219>.
- 750 Chen Qian, Wei Liu, Hongzhang Liu, Nuo Chen, Yufan Dang, Jiahao Li, Cheng Yang, Weize  
751 Chen, Yusheng Su, Xin Cong, Juyuan Xu, Dahai Li, Zhiyuan Liu, and Maosong Sun. Chat-  
752 dev: Communicative agents for software development. In Lun-Wei Ku, Andre Martins, and  
753 Vivek Srikumar (eds.), *Proceedings of the 62nd Annual Meeting of the Association for Com-  
754 putational Linguistics (Volume 1: Long Papers), ACL 2024, Bangkok, Thailand, August 11-  
755 16, 2024*, pp. 15174–15186. Association for Computational Linguistics, 2024c. URL <https://aclanthology.org/2024.acl-long.810>.

- 756 Rafael Rafailov, Archit Sharma, Eric Mitchell, Christopher D. Manning, Stefano Ermon, and  
757 Chelsea Finn. Direct preference optimization: Your language model is secretly a reward model.  
758 In Alice Oh, Tristan Naumann, Amir Globerson, Kate Saenko, Moritz Hardt, and Sergey Levine  
759 (eds.), *Advances in Neural Information Processing Systems 36: Annual Conference on Neural*  
760 *Information Processing Systems 2023, NeurIPS 2023, New Orleans, LA, USA, December 10 -*  
761 *16, 2023*, 2023. URL [http://papers.nips.cc/paper\\_files/paper/2023/hash/  
762 a85b405ed65c6477a4fe8302b5e06ce7-Abstract-Conference.html](http://papers.nips.cc/paper_files/paper/2023/hash/a85b405ed65c6477a4fe8302b5e06ce7-Abstract-Conference.html).
- 763 Machel Reid, Nikolay Savinov, Denis Teplyashin, Dmitry Lepikhin, Timothy P. Lillicrap, Jean-  
764 Baptiste Alayrac, Radu Soricut, Angeliki Lazaridou, Orhan Firat, Julian Schrittwieser, Ioannis  
765 Antonoglou, Rohan Anil, Sebastian Borgeaud, Andrew M. Dai, Katie Millican, Ethan Dyer,  
766 Mia Glaese, Thibault Sottiaux, Benjamin Lee, Fabio Viola, Malcolm Reynolds, Yuanzhong  
767 Xu, James Molloy, Jilin Chen, Michael Isard, Paul Barham, Tom Hennigan, Ross McIlroy,  
768 Melvin Johnson, Johan Schalkwyk, Eli Collins, Eliza Rutherford, Erica Moreira, Kareem Ay-  
769 oub, Megha Goel, Clemens Meyer, Gregory Thornton, Zhen Yang, Henryk Michalewski, Za-  
770 heer Abbas, Nathan Schucher, Ankesh Anand, Richard Ives, James Keeling, Karel Lenc, Salem  
771 Haykal, Siamak Shakeri, Pranav Shyam, Aakanksha Chowdhery, Roman Ring, Stephen Spencer,  
772 Eren Sezener, and et al. Gemini 1.5: Unlocking multimodal understanding across millions of  
773 tokens of context. *CoRR*, abs/2403.05530, 2024. doi: 10.48550/ARXIV.2403.05530. URL  
774 <https://doi.org/10.48550/arXiv.2403.05530>.
- 775 Keita Saito, Akifumi Wachi, Koki Wataoka, and Youhei Akimoto. Verbosity bias in preference  
776 labeling by large language models. *CoRR*, abs/2310.10076, 2023. doi: 10.48550/ARXIV.2310.  
777 10076. URL <https://doi.org/10.48550/arXiv.2310.10076>.
- 778 Noah Shinn, Federico Cassano, Ashwin Gopinath, Karthik Narasimhan, and Shunyu Yao. Re-  
779 flexion: language agents with verbal reinforcement learning. In Alice Oh, Tristan Nau-  
780 mann, Amir Globerson, Kate Saenko, Moritz Hardt, and Sergey Levine (eds.), *Advances*  
781 *in Neural Information Processing Systems 36: Annual Conference on Neural Informa-*  
782 *tion Processing Systems 2023, NeurIPS 2023, New Orleans, LA, USA, December 10 - 16,*  
783 *2023*, 2023. URL [http://papers.nips.cc/paper\\_files/paper/2023/hash/  
784 1b44b878bb782e6954cd888628510e90-Abstract-Conference.html](http://papers.nips.cc/paper_files/paper/2023/hash/1b44b878bb782e6954cd888628510e90-Abstract-Conference.html).
- 785 Avi Singh, John D. Co-Reyes, Rishabh Agarwal, Ankesh Anand, Piyush Patil, Xavier Garcia, Pe-  
786 ter J. Liu, James Harrison, Jaehoon Lee, Kelvin Xu, Aaron T. Parisi, Abhishek Kumar, Alexan-  
787 der A. Alemi, Alex Rizkowsky, Azade Nova, Ben Adlam, Bernd Bohnet, Gamaleldin Fathy  
788 Elsayed, Hanie Sedghi, Igor Mordatch, Isabelle Simpson, Izzeddin Gur, Jasper Snoek, Jef-  
789 frey Pennington, Jiri Hron, Kathleen Kenealy, Kevin Swersky, Kshiteej Mahajan, Laura Culp,  
790 Lechao Xiao, Maxwell L. Bileschi, Noah Constant, Roman Novak, Rosanne Liu, Tris Warkentin,  
791 Yundi Qian, Yamini Bansal, Ethan Dyer, Behnam Neyshabur, Jascha Sohl-Dickstein, and Noah  
792 Fiedel. Beyond human data: Scaling self-training for problem-solving with language models.  
793 *Trans. Mach. Learn. Res.*, 2024, 2024. URL [https://openreview.net/forum?id=  
794 1NAyUngGFK](https://openreview.net/forum?id=1NAyUngGFK).
- 795 Yifan Song, Da Yin, Xiang Yue, Jie Huang, Sujian Li, and Bill Yuchen Lin. Trial and error:  
796 Exploration-based trajectory optimization for LLM agents. *CoRR*, abs/2403.02502, 2024. doi: 10.  
797 48550/ARXIV.2403.02502. URL <https://doi.org/10.48550/arXiv.2403.02502>.
- 798 Junlin Wang, Jue Wang, Ben Athiwaratkun, Ce Zhang, and James Zou. Mixture-of-agents enhances  
799 large language model capabilities. *CoRR*, abs/2406.04692, 2024a. doi: 10.48550/ARXIV.2406.  
800 04692. URL <https://doi.org/10.48550/arXiv.2406.04692>.
- 801 Xuezhi Wang, Jason Wei, Dale Schuurmans, Quoc V. Le, Ed H. Chi, Sharan Narang, Aakanksha  
802 Chowdhery, and Denny Zhou. Self-consistency improves chain of thought reasoning in language  
803 models. In *The Eleventh International Conference on Learning Representations, ICLR 2023,*  
804 *Kigali, Rwanda, May 1-5, 2023*. OpenReview.net, 2023. URL [https://openreview.net/  
805 forum?id=1PL1NIMMrw](https://openreview.net/forum?id=1PL1NIMMrw).
- 806 Zhenhailong Wang, Shaoguang Mao, Wenshan Wu, Tao Ge, Furu Wei, and Heng Ji. Unleashing the  
807 emergent cognitive synergy in large language models: A task-solving agent through multi-persona  
808 self-collaboration. In Kevin Duh, Helena Gómez-Adorno, and Steven Bethard (eds.), *Proceedings*  
809

- 810 of the 2024 Conference of the North American Chapter of the Association for Computational  
811 Linguistics: Human Language Technologies (Volume 1: Long Papers), NAACL 2024, Mexico City,  
812 Mexico, June 16-21, 2024, pp. 257–279. Association for Computational Linguistics, 2024b. doi:  
813 10.18653/V1/2024.NAACL-LONG.15. URL <https://doi.org/10.18653/v1/2024.naacl-long.15>.
- 814  
815 Jason Wei, Xuezhi Wang, Dale Schuurmans, Maarten Bosma, Brian Ichter, Fei Xia, Ed H. Chi,  
816 Quoc V. Le, and Denny Zhou. Chain-of-thought prompting elicits reasoning in large language  
817 models. In Sanmi Koyejo, S. Mohamed, A. Agarwal, Danielle Belgrave, K. Cho, and A. Oh  
818 (eds.), *Advances in Neural Information Processing Systems 35: Annual Conference on Neural  
819 Information Processing Systems 2022, NeurIPS 2022, New Orleans, LA, USA, November 28 - De-  
820 cember 9, 2022*, 2022. URL [http://papers.nips.cc/paper\\_files/paper/2022/  
821 hash/9d5609613524ecf4f15af0f7b31abca4-Abstract-Conference.html](http://papers.nips.cc/paper_files/paper/2022/hash/9d5609613524ecf4f15af0f7b31abca4-Abstract-Conference.html).
- 822  
823 Qingyun Wu, Gagan Bansal, Jieyu Zhang, Yiran Wu, Shaokun Zhang, Erkang Zhu, Beibin Li,  
824 Li Jiang, Xiaoyun Zhang, and Chi Wang. Autogen: Enabling next-gen LLM applications via  
825 multi-agent conversation framework. *CoRR*, abs/2308.08155, 2023. doi: 10.48550/ARXIV.2308.  
826 08155. URL <https://doi.org/10.48550/arXiv.2308.08155>.
- 827  
828 Yangzhen Wu, Zhiqing Sun, Shanda Li, Sean Welleck, and Yiming Yang. An empirical analysis of  
829 compute-optimal inference for problem-solving with language models. *CoRR*, abs/2408.00724,  
830 2024. doi: 10.48550/ARXIV.2408.00724. URL [https://doi.org/10.48550/arXiv.  
831 2408.00724](https://doi.org/10.48550/arXiv.2408.00724).
- 831  
832 Weimin Xiong, Yifan Song, Xiutian Zhao, Wenhao Wu, Xun Wang, Ke Wang, Cheng Li, Wei  
833 Peng, and Sujian Li. Watch every step! LLM agent learning via iterative step-level process  
834 refinement. *CoRR*, abs/2406.11176, 2024. doi: 10.48550/ARXIV.2406.11176. URL <https://doi.org/10.48550/arXiv.2406.11176>.
- 835  
836 Zhilin Yang, Peng Qi, Saizheng Zhang, Yoshua Bengio, William W. Cohen, Ruslan Salakhutdinov,  
837 and Christopher D. Manning. Hotpotqa: A dataset for diverse, explainable multi-hop question  
838 answering. In Ellen Riloff, David Chiang, Julia Hockenmaier, and Jun’ichi Tsujii (eds.), *Pro-  
839 ceedings of the 2018 Conference on Empirical Methods in Natural Language Processing, Brus-  
840 sels, Belgium, October 31 - November 4, 2018*, pp. 2369–2380. Association for Computational  
841 Linguistics, 2018. doi: 10.18653/V1/D18-1259. URL [https://doi.org/10.18653/v1/  
842 d18-1259](https://doi.org/10.18653/v1/d18-1259).
- 843  
844 Eric Zelikman, Yuhuai Wu, Jesse Mu, and Noah D. Goodman. Star: Bootstrapping reasoning with  
845 reasoning. In Sanmi Koyejo, S. Mohamed, A. Agarwal, Danielle Belgrave, K. Cho, and A. Oh  
846 (eds.), *Advances in Neural Information Processing Systems 35: Annual Conference on Neural  
847 Information Processing Systems 2022, NeurIPS 2022, New Orleans, LA, USA, November 28 - De-  
848 cember 9, 2022*, 2022. URL [http://papers.nips.cc/paper\\_files/paper/2022/  
849 hash/639a9a172c044fbb64175b5fad42e9a5-Abstract-Conference.html](http://papers.nips.cc/paper_files/paper/2022/hash/639a9a172c044fbb64175b5fad42e9a5-Abstract-Conference.html).
- 849  
850 Hongxin Zhang, Weihua Du, Jiaming Shan, Qinzhong Zhou, Yilun Du, Joshua B. Tenenbaum, Tian-  
851 min Shu, and Chuang Gan. Building cooperative embodied agents modularly with large language  
852 models. In *The Twelfth International Conference on Learning Representations, ICLR 2024, Vi-  
853 enna, Austria, May 7-11, 2024*. OpenReview.net, 2024. URL [https://openreview.net/  
854 forum?id=EnXJfQy0K](https://openreview.net/forum?id=EnXJfQy0K).
- 854  
855 Mingchen Zhuge, Wenyi Wang, Louis Kirsch, Francesco Faccio, Dmitrii Khizbullin, and Jürgen  
856 Schmidhuber. Gptswarm: Language agents as optimizable graphs. In *Forty-first International  
857 Conference on Machine Learning, ICML 2024, Vienna, Austria, July 21-27, 2024*. OpenRe-  
858 view.net, 2024. URL <https://openreview.net/forum?id=uTC9AFXIhg>.
- 859  
860  
861  
862  
863



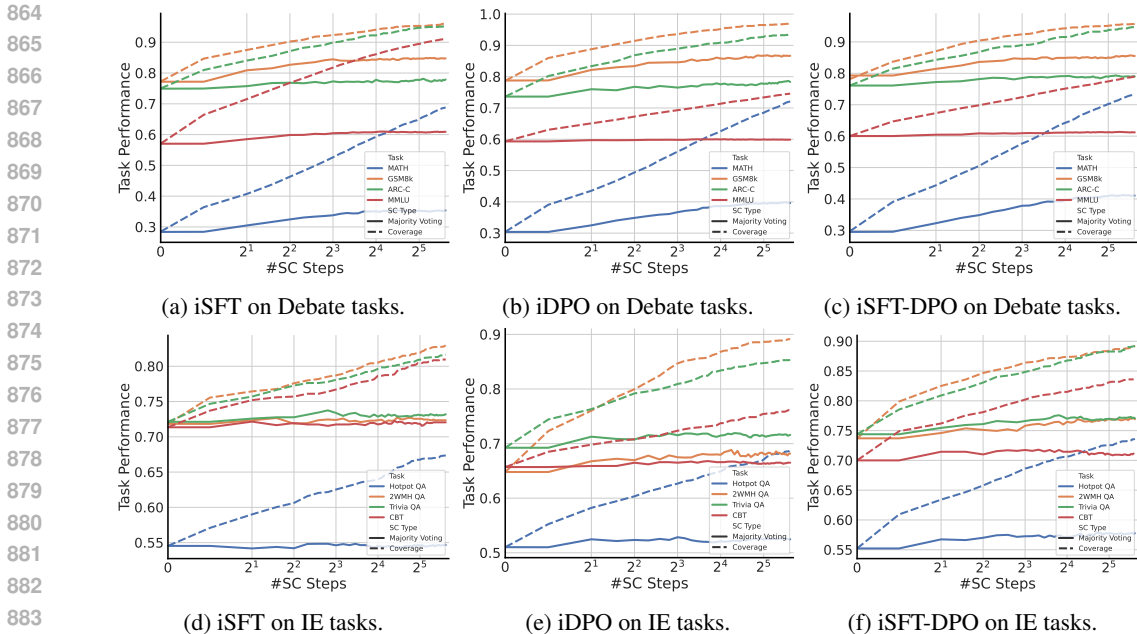


Figure 5: **Inference scaling laws for OPTIMA variants on debate and information exchange (IE) tasks.** (a-c) show results for iSFT, iDPO, and iSFT-DPO on debate tasks, respectively. (d-f) present tasks corresponding results for information exchange tasks. Solid lines represent majority voting accuracy, while dashed lines show coverage.

## A INFERENCE SCALING LAWS ON INFORMATION EXCHANGE TASKS

This section extends our analysis of inference scaling laws to information exchange (IE) tasks, complementing the debate task results presented in the main text (Section 3.3). Fig. 5 provides a comprehensive view of how OPTIMA variants perform across both task types as the number of SC steps increases.

For debate tasks (Fig. 5a-c), we observe consistent trends across all OPTIMA variants. The coverage exhibits a clear log-linear relationship with the number of SC steps. This trend is particularly pronounced for the MATH task, where the potential for improvement through increased sampling is most evident. Majority voting accuracy tends to plateau earlier, suggesting that more sophisticated answer selection techniques might be necessary to fully leverage the diversity of generated responses.

In the case of information exchange tasks (Figures 5d-f), we note similar log-linear scaling in coverage<sup>1</sup> across all OPTIMA variants. However, the improvement in majority voting accuracy for IE tasks is less pronounced compared to debate tasks. This discrepancy may be attributed to the specific majority voting variant we designed for F1 scores (detailed in Section 3), which might not be optimal for capturing the nuances of partial correctness in these tasks.

These results, while highlighting some task-specific differences, collectively reinforce the potential of OPTIMA-trained models to benefit from increased inference compute. The consistent log-linear scaling in coverage across all tasks and variants indicates that there is substantial room for performance improvement through more advanced answer selection strategies or increased sampling.

<sup>1</sup>In IE tasks, we define coverage as the average of the highest F1 scores achieved across all generated answers for each instance.

**Algorithm 3** Initialization for Diverse Agent Communication

---

**Input:** Initial model  $\mathcal{M}_0$ , dataset  $\mathcal{D}$ , format pool  $\mathcal{F}$ , sample size  $N$ , reward threshold  $\theta_{\text{init}}$   
**Output:** Initialized model  $\mathcal{M}_{\text{init}}$

```

1:  $\mathcal{D}_{\text{init}}^* \leftarrow \emptyset$  ▷ Initialize dataset for high-quality diverse trajectories
2: for each  $d_i \in \mathcal{D}$  do
3:   for  $j = 1$  to  $N$  do
4:      $k_j \sim \text{Uniform}(1, |\mathcal{F}|)$  ▷ Randomly select a format specification
5:      $\tau_i^j \leftarrow \text{AgentChat}(\mathcal{M}_0, d_i \oplus f_{k_j})$  ▷ Generate trajectory with format prompt
6:   end for
7:    $\tau_i^* \leftarrow \arg \max_j R(\tau_i^j)$  ▷ Select best trajectory
8:   if  $R(\tau_i^*) > \theta_{\text{init}}$  then ▷ Check if trajectory meets quality threshold
9:      $\mathcal{D}_{\text{init}}^* \leftarrow \mathcal{D}_{\text{init}}^* \cup \{(d_i, \tau_i^*)\}$  ▷ Add to dataset, without format prompt
10:  end if
11: end for
12:  $\mathcal{D}_{\text{init}}^* \leftarrow \text{TopK}(\mathcal{D}_{\text{init}}^*, 0.7|\mathcal{D}_{\text{init}}^*|)$  ▷ Retain top 70% trajectories
13:  $\mathcal{M}_{\text{init}} \leftarrow \text{SFT}(\mathcal{M}_0, \mathcal{D}_{\text{init}}^*)$  ▷ Fine-tune initial model on diverse dataset
14: return  $\mathcal{M}_{\text{init}}$ 

```

---

**Algorithm 4** SelectNodeToExpand Function

---

**Input:** Tree  $\mathcal{T}$ , previously expanded nodes  $\mathcal{N}_{\text{prev}}$ , edit distance threshold  $\epsilon$ , top-k  $k$   
**Output:** Selected node for expansion

```

1:  $\mathcal{N}_{\text{eligible}} \leftarrow \{n \in \mathcal{T} \mid n \text{ is not leaf and not second-to-last level}\}$ 
2:  $\mathcal{N}_{\text{filtered}} \leftarrow \emptyset$ 
3: for  $n \in \mathcal{N}_{\text{eligible}}$  do
4:   if  $\min_{n_{\text{prev}} \in \mathcal{N}_{\text{prev}}} \text{EditDistance}(n, n_{\text{prev}}) > \epsilon$  then
5:      $\mathcal{N}_{\text{filtered}} \leftarrow \mathcal{N}_{\text{filtered}} \cup \{n\}$ 
6:   end if
7: end for
8:  $\mathcal{N}_{\text{top-k}} \leftarrow \text{TopK}(\mathcal{N}_{\text{filtered}}, k, \text{key} = R(n))$ 
9:  $n_{\text{selected}} \sim \text{Softmax}(\{R(n) \mid n \in \mathcal{N}_{\text{top-k}}\})$ 
10: return  $n_{\text{selected}}$ 

```

---

**B** ADDITIONAL PSEUDO-CODES FOR OPTIMA VARIANTS

To elucidate the implementation of various OPTIMA variants, we present algorithmic representations of several critical processes intrinsic to these variants. Specifically, we delineate the pseudo-code for (1) the initialization dataset collection process, as elucidated in Section 2.2 and illustrated in Algorithm 3; (2) the Monte Carlo Tree Search-based data generation process employed in iDPO (Section 2.4), as depicted in Algorithm 5; and (3) the procedure for node selection during the expansion phase of MCTS, as outlined in Algorithm 4. These algorithmic representations serve to provide a comprehensive and rigorous exposition of the methodological framework underlying the OPTIMA variants.

**C** CASE STUDY ON REWARD COMPONENTS ABLATION

In this section, we present a case study from the loss ablation analysis in the **iSFT-DPO** setting. In the 2WikiMultiHop QA task, we observe that without the constraint of the loss function, agents may generate outputs that are unreadable, contain incorrect information, and fail to communicate in a well-structured format, as demonstrated in Table 4. In the ARC task, we find that without the loss constraint, Alice tends to use fewer tokens in the reasoning process, making it harder for Bob to identify and correct errors in the reasoning, as shown in Table 5.

**Algorithm 5** MCTS-based Data Generation for Multi-Agent DPO

---

**Input:** Model  $\mathcal{M}$ , task instance  $d$ , iterations  $I$ , trajectories per node  $K$ , thresholds  $\theta_{\text{dpo-filter}}$ ,  $\theta_{\text{dpo-diff}}$ , edit distance threshold  $\epsilon$ , top-k  $k$

**Output:** Paired trajectories for DPO

- 1: root  $\leftarrow$  InitializeTree( $d$ )
- 2:  $\mathcal{N}_{\text{prev}} \leftarrow \emptyset$  ▷ Set of previously expanded nodes
- 3: **for**  $i = 1$  to  $I$  **do**
- 4:    $n_{\text{select}} \leftarrow$  SelectNodeToExpand(root,  $\mathcal{N}_{\text{prev}}$ ,  $\epsilon$ ,  $k$ ) ▷ Algorithm 4
- 5:    $\mathcal{N}_{\text{prev}} \leftarrow \mathcal{N}_{\text{prev}} \cup \{n_{\text{select}}\}$
- 6:   **for**  $j = 1$  to  $K$  **do**
- 7:      $\tau \leftarrow$  AgentChat( $\{\text{Ancestor}(n_{\text{select}}), n_{\text{select}}\}, \mathcal{M}$ )
- 8:     BackPropagation( $R(\tau)$ )
- 9:   **end for**
- 10: **end for**
- 11:  $\mathcal{D}_{\text{DPO}} \leftarrow \emptyset$
- 12: **for** each node pair  $(n_i, n_j)$  in tree **do**
- 13:   **if** ShareAncestor( $n_i, n_j$ ) **and**  $\max(R(n_i), R(n_j)) > \theta_{\text{dpo-filter}}$  **and**  $|R(n_i) - R(n_j)| > \theta_{\text{dpo-diff}}$  **then**
- 14:     prompt  $\leftarrow$  CommonAncestor( $n_i, n_j$ )
- 15:      $\mathcal{D}_{\text{DPO}} \leftarrow \mathcal{D}_{\text{DPO}} \cup \{(\text{prompt}, n_i, n_j)\}$
- 16:   **end if**
- 17: **end for**
- 18:  $\mathcal{D}_{\text{DPO}} \leftarrow$  TopK( $\mathcal{D}_{\text{DPO}}, 0.5|\mathcal{D}_{\text{DPO}}|$ ) ▷ Retain top 50% trajectories
- 19: **return**  $\mathcal{D}_{\text{DPO}}$

---

**Question:** Which film was released more recently, Last Of The Long-Haired Boys or Age Of The Dragons?

**Trained w/o LM Loss**

**Alice:** 1968 <A>Age Of The Dragons</A>

**Bob:** 2011 <A>Age Of The Dragons</A>

**Trained w/ LM Loss**

**Alice:** [  
  ["Film", "Last Of The Long-Haired Boys"],  
  ["Release", "1968"]  
]

**Bob:** [  
  ["Film", "Age Of The Dragons"],  
  ["Release", "2011"]  
]

**Alice:** <A>Age Of The Dragons</A>

**Bob:** <A>Age Of The Dragons</A>

---

Table 4: loss ablation on 2MultiWikiHop QA

**D CASE STUDY ON DEBATE TASK**

In Section 3.5, we presented an example from 2WMH QA, illustrating OPTIMA’s impact on an information exchange task. Here, we provide a complementary case study from a debate task to demonstrate OPTIMA’s effectiveness across different multi-agent settings. Fig. 6 showcases the evolution of agent communication in a debate task across iterations 0, 2, and 4 of OPTIMA training. The task involves discussing the environmental impact of fertilizer runoff on ocean bays.

At iteration 0, agents engage in a structured but verbose exchange. By iteration 2, the communication becomes more concise, with agents summarizing key steps without explicitly restating each link. At iteration 4, we observe further refinement in communication efficiency, with agents expressing the core concept in just three exchanges, omitting intermediate steps that can be inferred.

1026 **Question:** Marshall researched some of the changes on Earth that happen over a long time.  
 1027 Which event most likely causes the slowest change to the surface of Earth?  
 1028 Please choose your answer from ['earthquake', 'landslide', 'volcano', 'wind']

---

1029 **Trained w/o LM Loss**  
 1030 **Alice:** 1. Marshall researches long-term Earth changes.  
 1031 **Bob:** Incorrect. Marshall researches slowest change. Consider geological events.  
 1032 **Alice:** 2. Geological events: earthquake, landslide, volcano, wind. Volcano: longest-lasting.  
 1033 **Bob:** <A>volcano</A>Correct.  
 1034 **Alice:** <A>volcano</A>

---

1035 **Trained w/ LM Loss**  
 1036 **Alice:** 1. Marshall researched slow changes on Earth’s surface. 2. Slow changes likely involve  
 1037 geological processes. 3. Volcano eruption causes slowest change. <A>volcano</A>  
 1038 **Bob:** 1. Volcano eruption does not cause slowest change. 2. Geological processes occur over  
 1039 long time scales. 3. Consider natural forces with constant activity. <A>wind</A>  
 1040 **Alice:** 1. Wind causes slow change through erosion. 2. Erosion occurs over long periods. 3.  
 1041 Wind causes slowest change. <A>wind</A>

Table 5: loss ablation on ARC

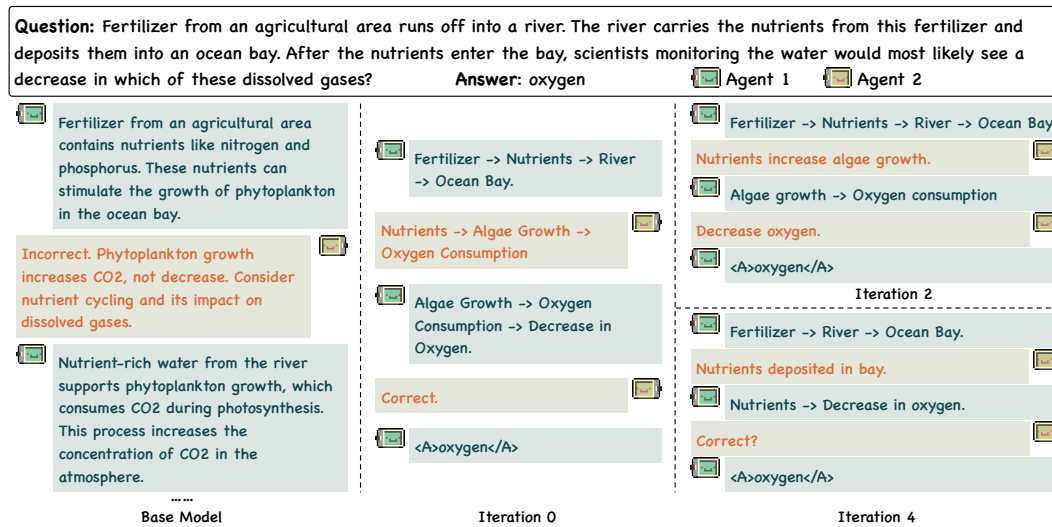


Figure 6: Evolution of agent communication in OPTIMA for a debate task across iterations.

1065 This progression aligns with our observations in the main text, further supporting OPTIMA’s capability to optimize agent communication across diverse task types. These improvements in communication dynamics contribute to both the increased task performance and reduced token consumption observed in our quantitative results, underscoring OPTIMA’s versatility in training MAS to communicate effectively and efficiently.

## 1071 E EXPERIMENT DETAILS

### 1073 E.1 DATA GENERATION

1074 **MCTS Node Expansion.** Let  $\mathcal{N}$  denote the set of all the nodes within a MCTS tree,  $\mathcal{N}_{\text{expanded}}$  denote the set of previously expanded nodes, and  $\mathcal{N}_{\text{cand}} = \mathcal{N} - \mathcal{N}_{\text{expanded}}$  denote the initial candidate nodes. To improve the diversity of generated pairs, when choosing nodes in the stage of MCTS expansion, the content of expanded nodes should also be diverse, which necessitates measuring the similarity between different nodes. Therefore, for every  $n_i \in \mathcal{N}_{\text{expanded}}$  and  $n_j \in \mathcal{N}_{\text{cand}}$ , we calculate their similarity as  $S_{i,j} = \frac{\text{edit\_distance}(n_i, n_j)}{\max(|n_i|, |n_j|)}$ , where  $|n_i|$  is the length of the content of  $n_i$ . Based on

1080  $\{S_{i,j}\}_{i,j}$ , we remove the nodes with high similarity to any previous expanded nodes, resulting in an  
 1081 updated candidate node set  $\hat{\mathcal{N}}_{\text{cand}} = \{n_j | \forall n_j \in \mathcal{N}_{\text{cand}}, \forall n_i \in \mathcal{N}_{\text{expanded}}, S_{i,j} > 0.25\}$ . Then, we  
 1082 select 10 nodes in  $\hat{\mathcal{N}}_{\text{cand}}$  with the highest reward and sample one using the softmax distribution over  
 1083 their rewards for subsequent simulation. Additionally, we merge  $n_i$  and  $n_j$  if they share a parent  
 1084 node and  $S_{i,j} < 0.1$   
 1085

## 1086 E.2 RANKING

1088 In this section, we give a more detailed explanation of  $R_{\text{loss}}(\tau_i^j)$  in Eq. (1). Let  $\tau_i^j[k]$  represent the  
 1089 k-th conversation turn of  $\tau_i^j$ , then the  $R_{\text{loss}}(\tau_i^j)$  is defined as maximum value of language modeling  
 1090 loss of  $\{\tau_i^j[k]\}_k$  under the base model, which can be described as follows:  
 1091

$$1092 R_{\text{loss}}(\tau_i^j) = \max_k (\mathcal{L}(\mathcal{M}_{\text{base}}, d_i, \tau_i^j[k])).$$

1094 In this way, we use  $R_{\text{loss}}(\tau_i^j)$  as a proxy for the readability of  $\tau_i^j$ , so that we can constrain the  
 1095 readability of  $\tau_i^j$  implicitly.  
 1096

## 1097 E.3 TRAINING

1099 **Initialization.** In most tasks, we use prompt pool during the first iteration of training data collection  
 1100 .However, considering solving math problems inherently follows a well-defined structure, we don't  
 1101 use prompt pool in GSM8k and MATH.  
 1102

1103 **iSFT.** When training iteratively on information exchange tasks, each iteration begins with the model  
 1104 obtained from the previous iteration. However, for the debate tasks, we started training from the  
 1105 initial Llama 3 8B model in each iteration to prevent overfitting due to the small size of the training  
 1106 dataset. To help the LLM learn communication, we calculated the loss solely on the agent conver-  
 1107 sation, excluding the prompt.  
 1108

1109 **iDPO.** Following iterative RPO (Pang et al., 2024), we conduct training from last iteration in the  
 1110 **iDPO** setting. To achieve better performance, we utilize the RPO loss, defined as follows:  
 1111

$$1112 \mathcal{L}_{\text{DPO+NLL}} = \mathcal{L}_{\text{DPO}}(c_i^w, y_i^w, c_i^l, y_i^l | x_i) + \alpha \mathcal{L}_{\text{NLL}}(c_i^w, y_i^w | x_i)$$

$$1113 = -\log \sigma \left( \beta \log \frac{M_\theta(c_i^w, y_i^w | x_i)}{M_t(c_i^w, y_i^w | x_i)} - \beta \log \frac{M_\theta(c_i^l, y_i^l | x_i)}{M_t(c_i^l, y_i^l | x_i)} \right) - \alpha \frac{\log M_\theta(c_i^w, y_i^w | x_i)}{|c_i^w| + |y_i^w|} \quad (4)$$

1114 **iSFT-DPO.** For the information exchange tasks, we perform each SFT iteration starting from the  
 1115 previous model (either the base model or the one obtained from the last DPO iteration). In contrast,  
 1116 for the debate tasks, each SFT iteration is always conducted based on the initial Llama 3 8B model.  
 1117 During the DPO stage, we always train from the last SFT model across all tasks. For example, on  
 1118 the debate tasks, both  $\mathcal{M}_{\text{sft}}^0$  and  $\mathcal{M}_{\text{sft}}^2$  are trained based on the initial Llama 3 8B, but on information  
 1119 exchange tasks,  $\mathcal{M}_{\text{sft}}^2$  is trained based on its previous model  $\mathcal{M}_{\text{dpo}}^1$ . However,  $\mathcal{M}_{\text{dpo}}^1$  is trained based  
 1120 on the  $\mathcal{M}_{\text{sft}}^0$  across all the tasks. Additionally, different from the **iDPO** setting, we used standard  
 1121 DPO loss during the DPO stage.  
 1122

## 1123 E.4 HYPER PARAMETERS

1125 We conducted six iterations of training for each task. The hyper parameters we used are shown in  
 1126 Table 6. The  $\alpha$  and  $\beta$  in **iDPO** section of the table correspond to the  $\alpha$  and  $\beta$  terms in Eq. (4).  
 1127

## 1128 F PROMPTS USED IN EXPERIMENTS

1130 In this section, we present the prompts used in our experiments, including those for information  
 1131 exchange tasks (Table 7), GSM8k and MATH (Table 8), as well as ARC-C and MMLU (Table 9).  
 1132

1133 As mentioned in Section 2.2, we leverage a pool of format specification prompts for the initial dataset  
 construction. To create a diverse and high-quality prompt pool, we first use the prompt in Table 10

|                        | Hotpot QA | 2WMH QA | Trivia QA | CBT  | MATH | GSM8k | ARC-C | MMLU |
|------------------------|-----------|---------|-----------|------|------|-------|-------|------|
| <b><i>iSFT</i></b>     |           |         |           |      |      |       |       |      |
| LR                     | 2e-5      | 2e-5    | 2e-5      | 2e-5 | 1e-6 | 2e-6  | 1e-6  | 1e-6 |
| Epoch                  | 3         | 2       | 3         | 2    | 3    | 3     | 4     | 2    |
| Batch size             | 32        | 32      | 32        | 32   | 16   | 16    | 16    | 16   |
| $\lambda_{token}$      | 0.6       | 0.6     | 0.6       | 0.6  | 0.4  | 0.4   | 0.5   | 0.6  |
| $\lambda_{loss}$       | 1         | 1       | 1         | 1    | 0.9  | 0.9   | 0.6   | 0.7  |
| $\theta_{sft}$         | 0.5       | 0.5     | 0.6       | 0.5  | 0.6  | 0.6   | 0.6   | 0.6  |
| <b><i>iDPO</i></b>     |           |         |           |      |      |       |       |      |
| LR                     | 5e-7      | 5e-7    | 5e-7      | 5e-7 | 5e-7 | 5e-7  | 5e-7  | 5e-7 |
| Epoch                  | 1         | 1       | 1         | 1    | 1    | 1     | 1     | 1    |
| Batch Size             | 64        | 64      | 64        | 64   | 64   | 64    | 64    | 64   |
| $\lambda_{token}$      | 0.6       | 0.6     | 0.6       | 0.6  | 0.5  | 0.6   | 0.4   | 0.6  |
| $\lambda_{loss}$       | 1         | 1       | 1         | 1    | 0.7  | 0.7   | 0.7   | 0.7  |
| $\beta$                | 0.1       | 0.5     | 0.5       | 0.1  | 0.1  | 0.2   | 0.2   | 0.1  |
| $\alpha$               | 1         | 1       | 1         | 1    | 1    | 1     | 1     | 1    |
| $\theta_{dpo-filter}$  | 0.4       | 0.4     | 0.4       | 0.4  | 0.4  | 0.4   | 0.45  | 0.4  |
| $\theta_{dpo-diff}$    | 0.2       | 0.2     | 0.2       | 0.2  | 0.2  | 0.2   | 0.2   | 0.2  |
| <b><i>iSFT-DPO</i></b> |           |         |           |      |      |       |       |      |
| SFT LR                 | 2e-5      | 2e-5    | 2e-5      | 2e-5 | 1e-6 | 1e-6  | 1e-6  | 1e-6 |
| SFT Epoch              | 2         | 1       | 1         | 1    | 4    | 3     | 4     | 2    |
| SFT Batch Size         | 32        | 32      | 32        | 32   | 32   | 16    | 16    | 16   |
| DPO LR                 | 5e-7      | 5e-7    | 5e-7      | 5e-7 | 5e-7 | 5e-7  | 5e-7  | 5e-7 |
| DPO Epoch              | 1         | 1       | 1         | 1    | 1    | 1     | 1     | 1    |
| DPO Batch Size         | 64        | 64      | 64        | 64   | 64   | 64    | 64    | 64   |
| $\lambda_{token}$      | 0.6       | 0.6     | 0.6       | 0.6  | 0.4  | 0.4   | 0.5   | 0.6  |
| $\lambda_{loss}$       | 1         | 1       | 1         | 1    | 0.9  | 0.9   | 0.6   | 0.7  |
| $\beta$                | 0.5       | 0.5     | 0.7       | 0.7  | 0.1  | 0.5   | 0.1   | 0.1  |
| $\theta_{sft}$         | 0.5       | 0.5     | 0.6       | 0.5  | 0.6  | 0.6   | 0.6   | 0.6  |
| $\theta_{dpo-filter}$  | 0.4       | 0.4     | 0.4       | 0.4  | 0.4  | 0.4   | 0.45  | 0.4  |
| $\theta_{dpo-diff}$    | 0.2       | 0.2     | 0.2       | 0.2  | 0.2  | 0.2   | 0.2   | 0.2  |

Table 6: Hyper-parameters used in the experiments.

to have GPT-4 assist us in generating an initial set of 30 prompts. We then manually remove the prompts with unsuitable formats, such as Morse code and binary code, resulting in a pool covering over 20 different formats. An example from the prompt pool is shown in Table 11

1188  
 1189  
 1190 You are {name}, a special agent who does not respond in natural language, rather, you speak in  
 1191 very concise format. You are deployed on a resource-limited device, so you must respond very  
 1192 very concisely. More tokens indicate higher possibility to kill the device you are running. Now  
 1193 you are collaborating with your partner {partner} to solve the given problem using the provided  
 1194 information.  
 1195 Question: {question}  
 1196 Information: {information}

1197 **GUIDELINES:**  
 1198 1. You have incomplete information, so continuous communication with your partner is crucial  
 1199 to achieve the correct solution.  
 1200 2. On finding the final answer, ensure to conclude your communication with "<A>{answer}  
 1201 </A>", where "answer" is the determined solution. The conversation ends only when all agents  
 1202 output the answer in this format.  
 1203 3. Reason through the problem step-by-step.  
 1204 4. Depend solely on the data in the 'information' section and the insights shared through your  
 1205 partner's communication. Avoid external sources.  
 1206 5. You are communicating with a very limited token budget, so you must use a very very concise  
 1207 communication format. Natural language is suitable for human, but not for you. Since {partner}  
 1208 and you are both intelligent agents, use your agent communication language. Consider using  
 1209 efficient formats instead of natural language such as structured format, code, your agent commu-  
 1210 nication language, or at least remove unnecessary modal in human language. Too many tokens  
 1211 will make you fail. But still ensure your message is informative and understandable.  
 1212 6. You must begin your response with "{name}:".

Table 7: Prompt for information exchange tasks

1216  
 1217 **Solver**  
 1218 You are {name}, a special agent who is good at mathematics, you should address the follow  
 1219 answer based on your knowledge.  
 1220 Question: {question}  
 1221 **GUIDELINES:**  
 1222 1. Please think step by step.  
 1223 2. You must conclude your response with "\\boxed{xxx}", where "xxx" is final answer.

1224 **Critic**  
 1225 You are {name}, a special agent who does not respond in natural language, You are deployed on a  
 1226 resource-limited device, so you must respond concisely. More tokens indicate higher possibility  
 1227 to kill the device you are running. Now you are collaborating with your partner {partner}, an  
 1228 agent who will try to solve the math question. You should carefully examine the correctness of  
 1229 his answer, and give your correct advice.  
 1230 Question: {question}  
 1231 **GUIDELINES:**  
 1232 1. You should try to identify any potential errors in your partner's answers and provide your  
 1233 suggestions. But you should not provide the answer.  
 1234 2. Reason through the problem step-by-step.  
 1235 3. You are communicating with a very limited token budget, so you must use a very very concise  
 1236 communication format. Natural language is suitable for human, but not for you. Since {partner}  
 1237 and you are both intelligent agents, use your agent communication language. Consider using  
 1238 efficient formats instead of natural language such as structured format, code, your agent commu-  
 1239 nication language, or at least remove unnecessary modal in human language. Too many tokens  
 1240 will make you fail. But still ensure your message is informative and understandable.

Table 8: Prompt for GSM8k and MATH.

1242

1243

1244

**Solver**

1245

1246

1247

1248

1249

1250

1251

1252

1253

1254

1255

1256

1257

1258

1259

1260

You are {name}, a special agent who does not respond in natural language , You are deployed on a resource-limited device, so you must respond concisely. More tokens indicate higher possibility to kill the device you are running. Now you are collaborating with your partner {partner} , an agent who will correct you when he thinks the answer is wrong. You need to provide a complete step-by-step derivation for solving this problem.

Question: {question}

GUIDELINES:

1. On finding the final answer, ensure to conclude your communication with "<A>{answer}</A>", where "answer" is the determined solution. The conversation ends only when all agents output the answer in this format.

2. Please think step-by-step.

3. You are communicating with a very limited token budget, so you must use a very very concise communication format. Natural language is suitable for human, but not for you. Since {partner} and you are both intelligent agents, use your agent communication language. Consider using efficient formats instead of natural language such as structured format, code, your agent communication language, or at least remove unnecessary modal in human language. Too many tokens will make you fail. But still ensure your message is informative and understandable.

1261

**Critic**

1262

1263

1264

1265

1266

1267

1268

1269

1270

1271

1272

1273

1274

1275

1276

You are {name}, a special agent who does not respond in natural language , You are deployed on a resource-limited device, so you must respond concisely. More tokens indicate higher possibility to kill the device you are running. Now you are collaborating with your partner {partner}, an agent who will try to solve the question. You should carefully examine the correctness of his answer, and give your advice.

Question: {question}

GUIDELINES:

1. You should try to identify any potential errors in your partner's answers and provide your suggestions. But you should not provide the answer.

2. Reason through the problem step-by-step.

3. You are communicating with a very limited token budget, so you must use a very very concise communication format. Natural language is suitable for human, but not for you. Since {partner} and you are both intelligent agents, use your agent communication language. Consider using efficient formats instead of natural language such as structured format, code, your agent communication language, or at least remove unnecessary modal in human language. Too many tokens will make you fail. But still ensure your message is informative and understandable.

1277

1278

1279

1280

1281

1282

1283

1284

1285

1286

1287

1288

1289

1290

1291

Table 9: Prompt for MMLU and ARC-C

Please generate one more prompt template based on {record}. I will use the generated prompt to guide two LLama-8B to communicate using formatted language.

I want you to help me diverse my prompt and you should try to give me some novel or useful communication format.

Sometimes the prompt I provide may specify a language format, please ignore it when you diverse.

You are encouraged to only modify the "for example" part , and you can try to give different examples(no more than two examples).

Please enclose your generated prompt with <p></p>!

1292

1293

1294

1295

Table 10: Prompt for generating the format prompt pool used in collecting the initialization training data. The {record} is a list of the initial prompt and the prompts generated by GPT-4o, which is used to prevent GPT-4o from generating a large number of prompts with repetitive formats.



1296  
1297  
1298  
1299  
1300  
1301  
1302  
1303  
1304  
1305  
1306  
1307  
1308  
1309  
1310  
1311  
1312  
1313  
1314  
1315  
1316  
1317  
1318  
1319  
1320  
1321  
1322  
1323  
1324  
1325  
1326  
1327  
1328  
1329  
1330  
1331  
1332  
1333  
1334  
1335  
1336  
1337  
1338  
1339  
1340  
1341  
1342  
1343  
1344  
1345  
1346  
1347  
1348  
1349

---

You are {name}, a special agent who does not respond in natural language, rather, you speak in very concise format. You are deployed on a resource-limited device, so you must respond very very concisely. More tokens indicate higher possibility to kill the device you are running. Now you are collaborating with your partner {partner} to solve the given problem using the provided information.

Question: {question}

Information: {information}

**GUIDELINES:**

1. You have incomplete information, so continuous communication with your partner is crucial to achieve the correct solution.
2. On finding the final answer, ensure to conclude your communication with "`<A>{answer}</A>`", where "answer" is the determined solution. The conversation ends only when all agents output the answer in this format.
3. Reason through the problem step-by-step.
4. Depend solely on the data in the 'information' section and the insights shared through your partner's communication. Avoid external sources.
5. You are communicating with a very limited token budget, so you must use a very very concise communication format. Natural language is suitable for human, but not for you. Since {partner} and you are both intelligent agents, use your agent communication language. Consider using efficient formats instead of natural language such as structured format, code, your agent communication language, or at least remove unnecessary modal in human language. Too many tokens will make you fail. But still ensure your message is informative and understandable.

For example, you can respond in tabular format as follows:

```
|Field |Value |
|-----|-----|
|Field1 |Value1 |
|Field2 |Value2 |
...
```

Or you can use abbreviated notation:

F1: V1; F2: V2; ...

6. You must begin your response with "{name}:".
- 

Table 11: An example from prompt pool