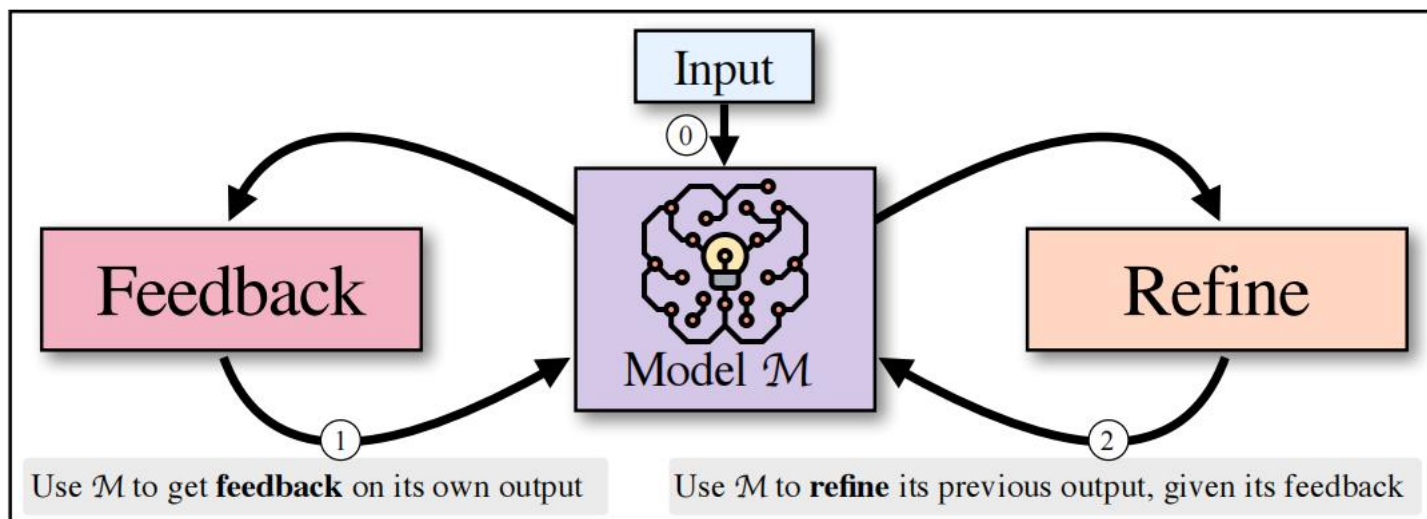# Breaking Mental Set to Improve Reasoning through Diverse Multi-Agent Debate

Yexiang Liu, Jie Cao, Zekun Li, Ran He, Tieniu Tan

# Background

- LLMs often suffer from mistakes when reasoning.

- We can use stronger model to provide feedback.

- Or utilize human supervision.

- However, effective feedback is not always obtainable.

- We need to study how to teach LLMs to self-correct.



The framework of *Self-Reflection*.

Self-Refine: Iterative Refinement with Self-Feedback. NeurIPS 2023.
Language Models can Solve Computer Tasks. NeurIPS 2023.
Reflexion: Language Agents with Verbal Reinforcement Learning. NeurIPS 2023.

- **Many studies have found limitations to Self-Reflection**, performing poorly on reasoning or even worsening.

- If we provide oracles, LLMs can gradually reach the correct answer. However, in practice, we do not know the answer, neither do LLMs.

Table 1: Results of GPT-3.5 and GPT-4 on reasoning benchmarks with the setting in Section 3.1.1.

|  |  | GSM8K | CommonSenseQA | HotpotQA |
|---|---|---|---|---|
| GPT-3.5 | Standard Prompting | 75.9 | 75.8 | 26.0 |
|  | Self-Correct (Oracle) | 84.3 | 89.7 | 29.0 |
| GPT-4 | Standard Prompting | 95.5 | 82.0 | 49.0 |
|  | Self-Correct (Oracle) | 97.5 | 85.5 | 59.0 |

Table 3: Results of GPT-3.5 and GPT-4 on reasoning benchmarks with *intrinsic self-correction*.

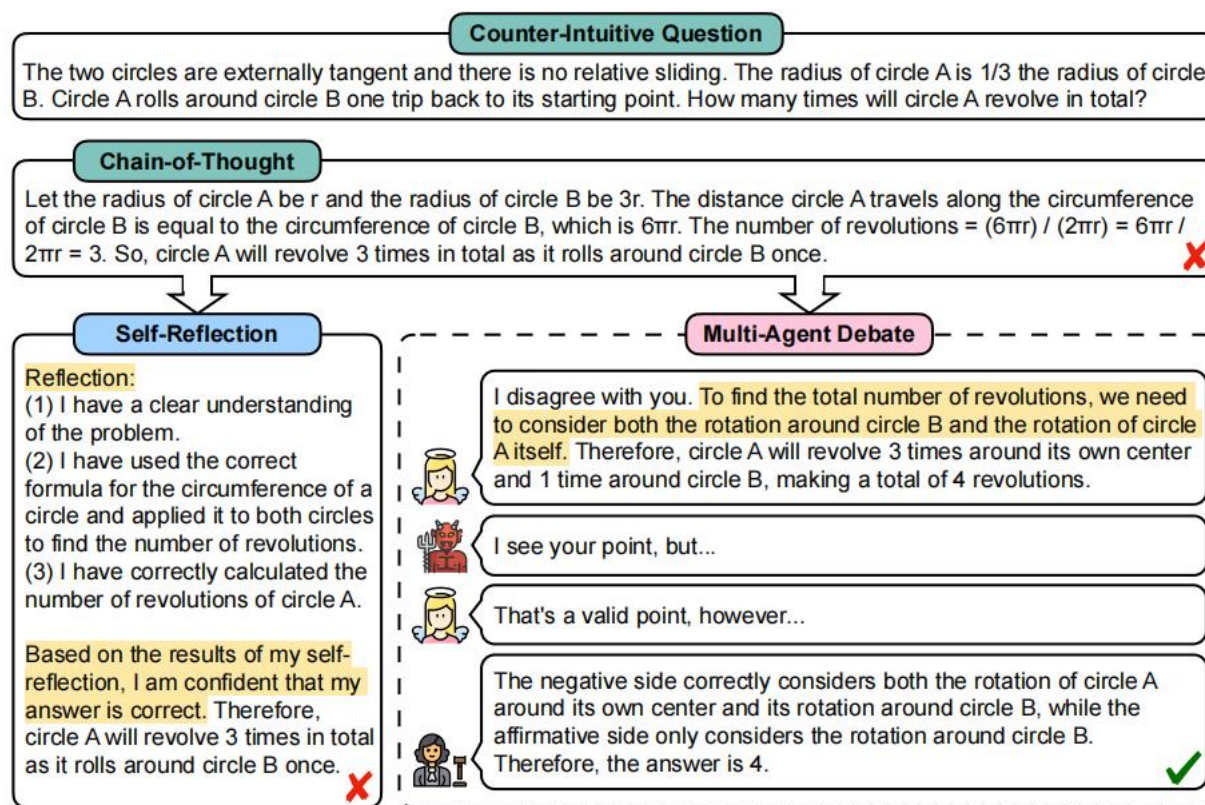|  |  | # calls | GSM8K | CommonSenseQA | HotpotQA |
|---|---|---|---|---|---|
| GPT-3.5 | Standard Prompting | 1 | **75.9** | **75.8** | **26.0** |
|  | Self-Correct (round 1) | 3 | 75.1 | 38.1 | 25.0 |
|  | Self-Correct (round 2) | 5 | 74.7 | 41.8 | 25.0 |
| GPT-4 | Standard Prompting | 1 | **95.5** | **82.0** | **49.0** |
|  | Self-Correct (round 1) | 3 | 91.5 | 79.5 | **49.0** |
|  | Self-Correct (round 2) | 5 | 89.0 | 80.0 | 43.0 |

Large Language Models Cannot Self-Correct Reasoning Yet. ICLR 2024.
GPT-4 Doesn't Know It's Wrong: An Analysis of Iterative Prompting for Reasoning Problems. NeurIPS 2023.
Can Large Language Models Really Improve by Self-critiquing Their Own Plans? NeurIPS 2023.
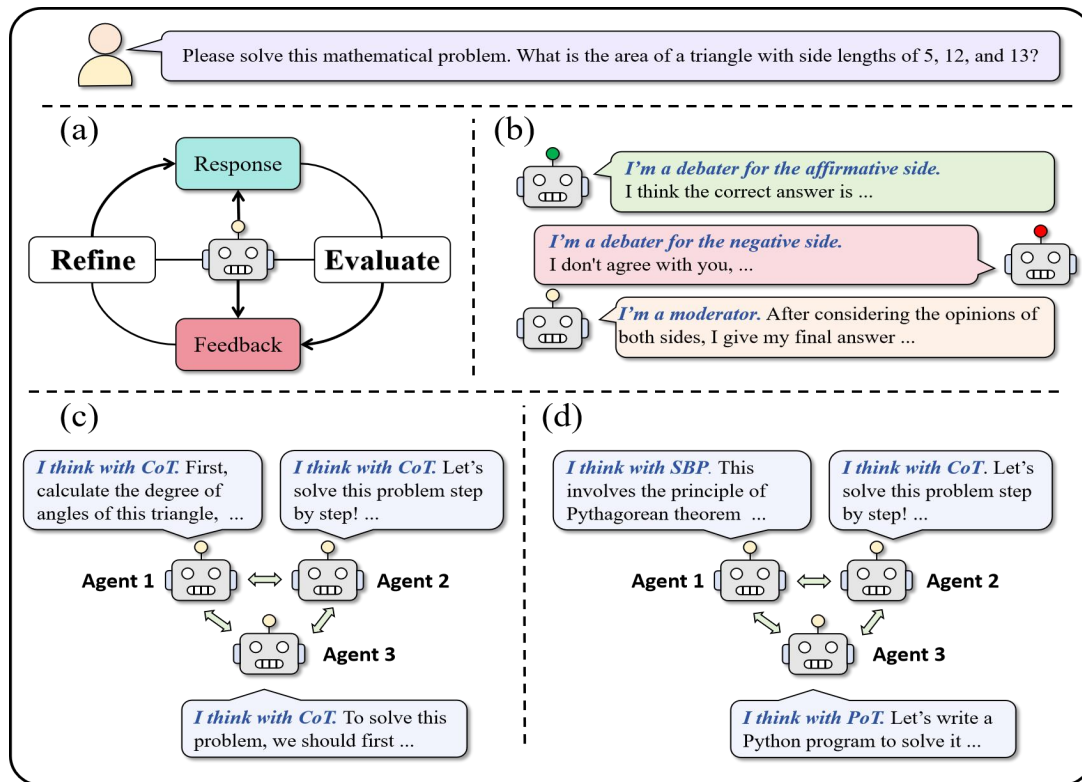
# Multi-Agent Debate (MAD)

- The problem of *Self-Reflection*: Degeneration of Thought.

- MAD: Multiple agents express their arguments in the state of "tit for tat" and a judge manages the debate process to obtain a final solution.

Encouraging Divergent Thinking in Large Language Models through Multi-Agent Debate. EMNLP 2024.
Improving Factuality and Reasoning in Language Models through Multiagent Debate. ICML 2024.

# Diverse Multi-Agent Debate (DMAD)

- Both *Self-Reflection* and MAD suffer from the mental set.

- We propose DMAD. By leveraging diverse problem-solving strategies, each agent can gain insights from different perspectives, refining its responses through discussion and collectively arriving at the optimal solution.

(a) *Self-Reflection*
(b) MAD-persona
(c) MAD
(d) DMAD

# Diverse Multi-Agent Debate (DMAD)

- MAD with a fixed strategy may always get the wrong answer.
- However, it succeeds just by transforming its thinking.

# Diverse Multi-Agent Debate (DMAD)

- MAD with a fixed strategy may always get the wrong answer.

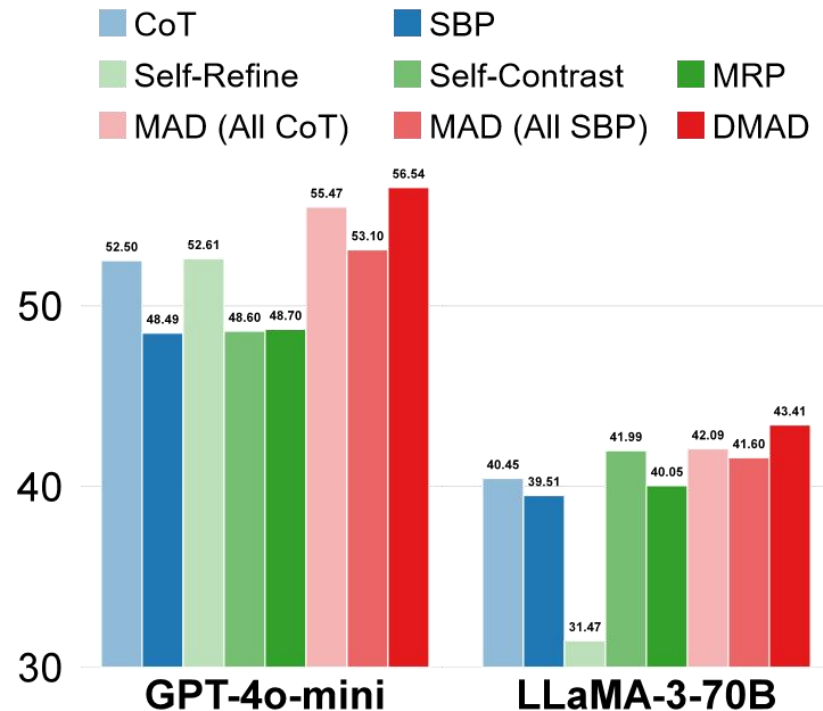- However, it succeeds just by transforming its thinking.

# Diverse Multi-Agent Debate (DMAD)

- DMAD can more effectively solve other methods' mental set problems.

- **What is mental set?**

In our paper, we introduce a new concept of **mental set** according to the psychological theory. Here we supplement a specific definition for it. Denote MAD (All CoT), MAD (All SBP), and MAD (All PoT) as $M_1$, $M_2$ and $M_3$ respectively. When using a kind of MAD method $M_i$ to solve a problem, if all agents consistently get wrong answers in all debate rounds, we assume that $M_i$ is unable to solve the problem correctly. Record all such problems for $M_i$ as the set $P_i$, and get $P = P_1 \cap P_2 \cap P_3$. For a problem $p \in P_i$, if it satisfies $p \notin P$, we define that the problem $p$ causes **mental set** of $M_i$, and define $p$ as the **mental set problem** of $M_i$. It means although $M_i$ constantly gets wrong solutions, the model can correctly solve the problem by changing to another method.

|  | MAD (All CoT) | MAD (All SBP) | MAD (All PoT) |
|---|---|---|---|
| Number of mental set problems | 70 | 87 | 67 |
| Problems that MAD (All CoT) correctly solves | 0 | 45 (51.72%) | 46 (68.7%) |
| Problems that MAD (All SBP) correctly solves | 28 (40.0%) | 0 | 31 (46.3%) |
| Problems that MAD (All PoT) correctly solves | **49 (70.0%)** | 51 (58.62%) | 0 |
| Problems that DMAD correctly solves | 48 (68.6%) | **60 (69.0%)** | **49 (73.1%)** |

# Diverse Multi-Agent Debate (DMAD)

- We evaluate DMAD against various prompting techniques, including *self-reflection* and traditional MAD, across multiple benchmarks including math, chemistry, physics, biology and so on, using both LLMs and Multimodal LLMs. Our experiments show that DMAD consistently outperforms other methods.
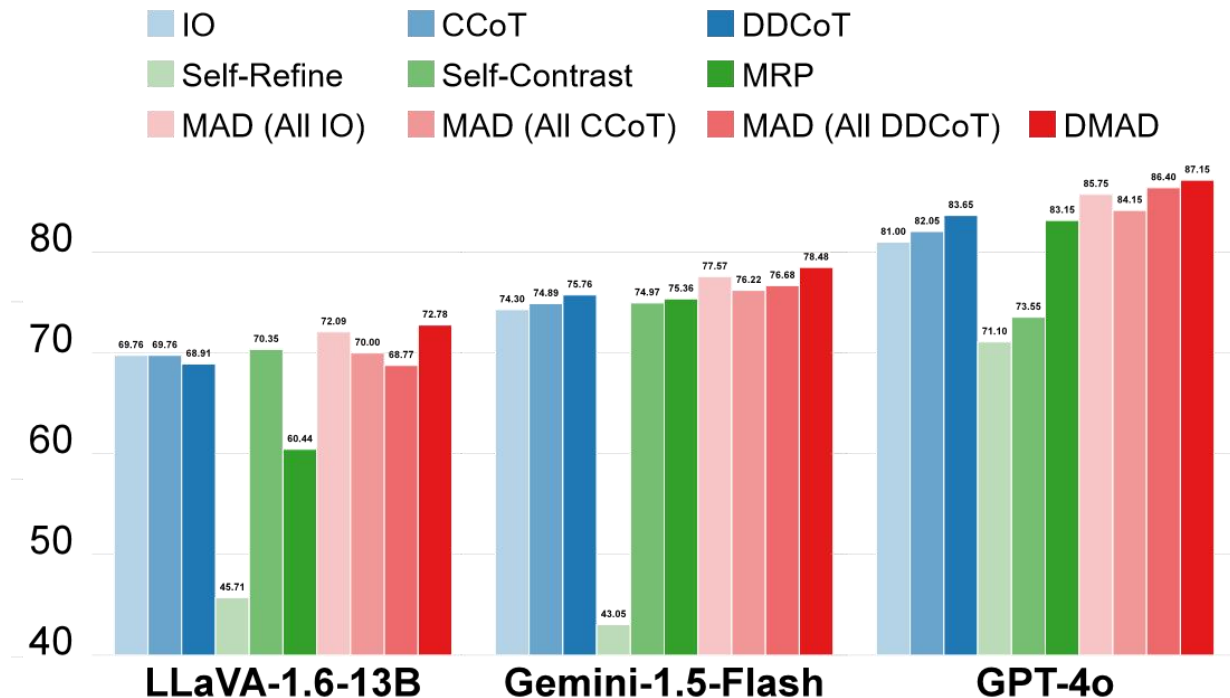


(a) LLMs

- We evaluate DMAD against various prompting techniques, including *self-reflection* and traditional MAD, across multiple benchmarks including math, chemistry, physics, biology and so on, using both LLMs and Multimodal LLMs. Our experiments show that DMAD consistently outperforms other methods.



(b) MLLMs

- What's more, DMAD can deliver better results than MAD in fewer rounds, and perform better when increasing the number of reasoning methods (for n = 1,2,3.).
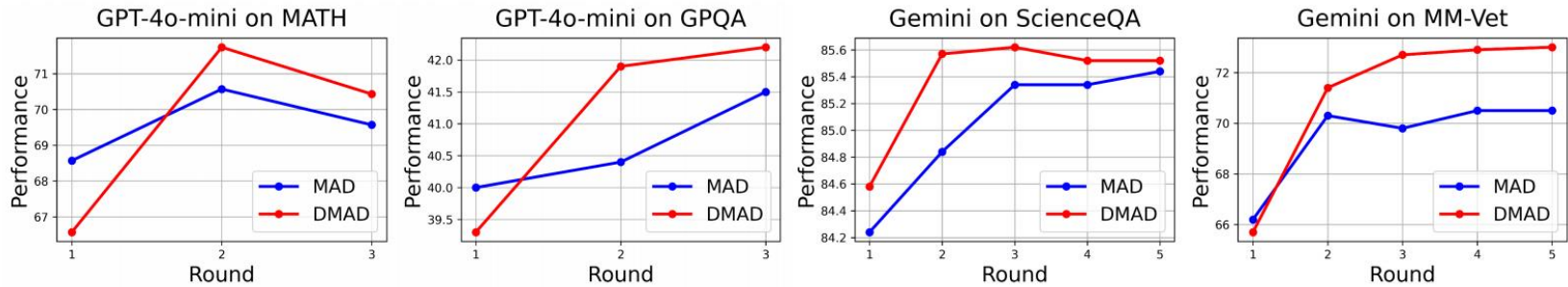


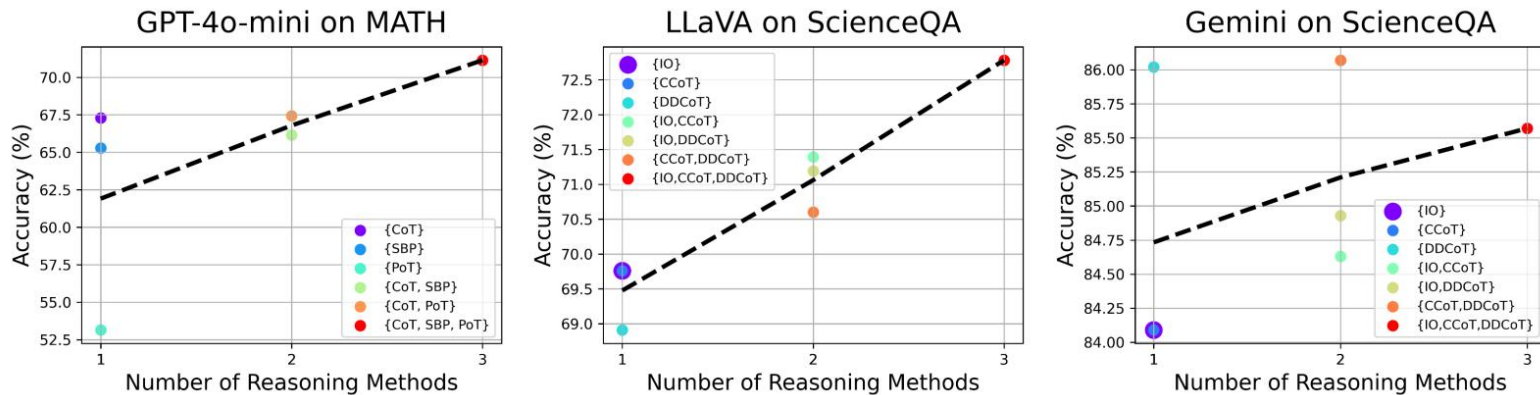Figure 4: Performance with increased rounds. More results are shown in Figure 6.



Figure 5: Performance vs the number of reasoning methods on DMAD.

# Limitations and Future work

■ Although DMAD can perform better, there are still some limitations:

- The inherent problem of MAD: LLM agents are easily influenced by other agents that come to incorrect conclusions, mistakenly changing their original correct solutions to wrong ones.

- There is still room for improvement in accuracy.

- How to design an intelligent system where agents can dynamically select the most suitable reasoning strategy.

- High overhead.

# Thanks