

# Level up your Agents with LBM

- Ajay
- Adrian Lam
- Benedict Neo
- Tan Wei Chun



# Problem Statement

- LLMs fail to generalize for specific expertise
- Hallucinations and lack of creativity
- Fixed by large scale agent testing with human in the loop





# Our Solution

- Generate multiple agents using various divergent styles
- Evaluate using an LLM
- Align the evaluations with human feedback

## Encouraging Divergent Thinking in Large Language Models through Multi-Agent Debate

Tian Liang, Zhiwei He, Wenxiang Jiao, Xing Wang, Rui Wang, Yujiu Yang, Zhaopeng Tu, Shuming Shi

Modern large language models (LLMs) like ChatGPT have shown remarkable performance on general language tasks but still struggle on complex reasoning tasks, which drives the research on cognitive behaviors of LLMs to explore human-like problem-solving strategies. Along this direction, one representative strategy is self-reflection, which asks an LLM to refine the solution with the feedback generated by itself iteratively. However, our study shows that such reflection-style methods suffer from the Degeneration-of-Thought (DoT) problem: once the LLM has

## Improving Factuality and Reasoning in Language Models through Multiagent Debate

Yilun Du<sup>1</sup>, Shuang Li<sup>1</sup>, Antonio Torralba<sup>1</sup>, Joshua B. Tenenbaum<sup>1</sup>, Igor Mordatch<sup>2</sup>

<sup>1</sup> MIT <sup>2</sup> Google Brain

ICML 2024

 Paper

 Code

[Submitted on 23 Feb 2024]

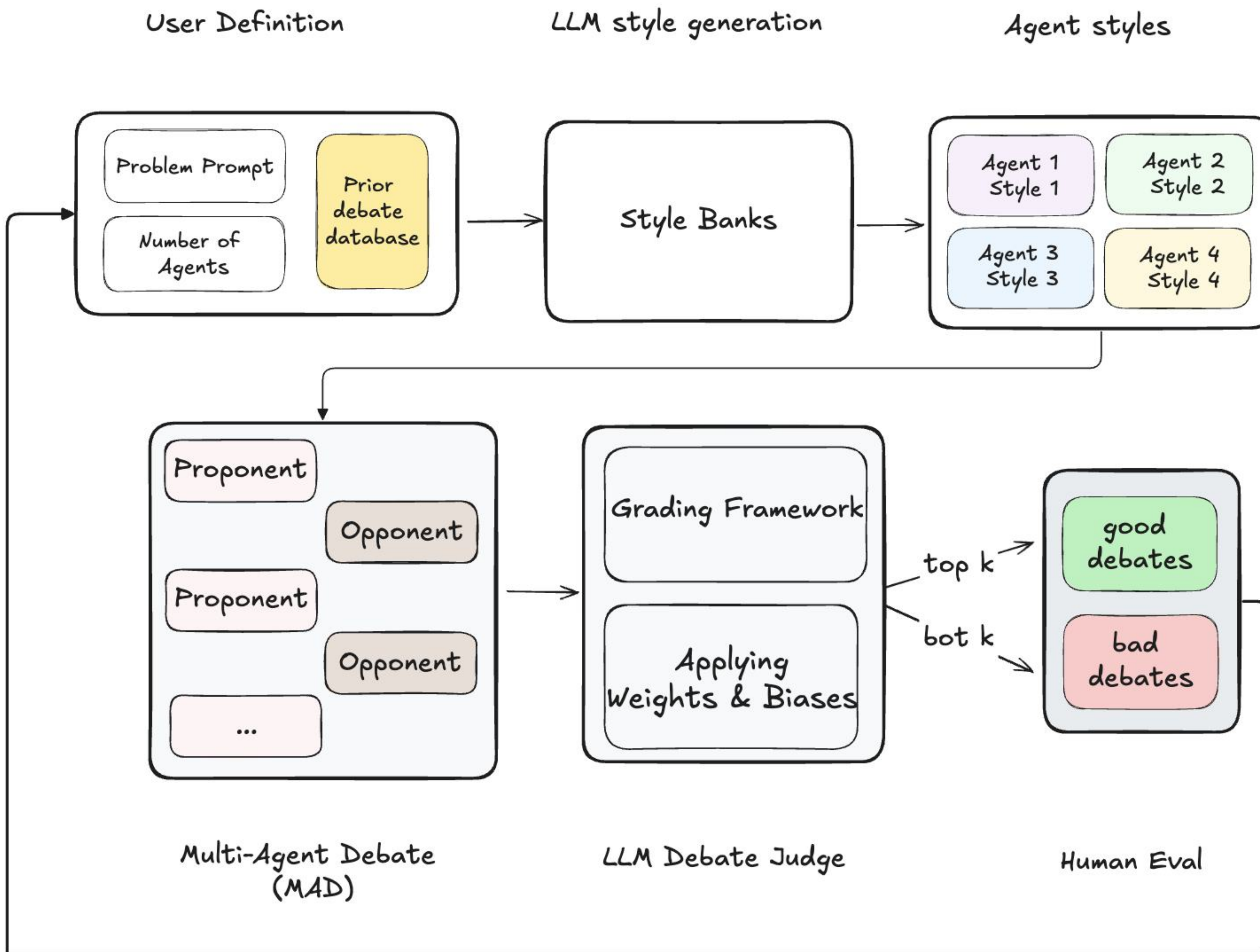
## Fine-Grained Self-Endorsement Improves Factuality and Reasoning

Ante Wang, Linfeng Song, Baolin Peng, Ye Tian, Lifeng Jin, Haitao Mi, Jinsong Su, Dong Yu

This work studies improving large language model (LLM) generations at inference time by mitigating fact-conflicting hallucinations through a self-endorsement framework that leverages the fine-grained fact-level comparisons across multiple sampled responses. Compared with previous methods (e.g., (Liu et al., 2022; Chen et al., 2023)) that perform response-level selection, our approach can better alleviate hallucinations, especially for long-form generation. Our method can broadly benefit smaller and open-source LLMs as it mainly conducts simple content-based comparisons. Experiments on BIG-Bench and other benchmarks effectively improve the factuality of generations with simple and intuitive prompts across different scales of LLMs. Besides, comparisons on the long-form generation task GSM8K demonstrate the potential of self-endorsement for broader application.

gh  
Multi-  
state of  
r MAD  
quire  
nse  
ness of  
e modest  
find that  
[his https](#)





# LLM Debating Styles

The **optimistic** debating style presents a positive outlook and focuses on potential benefits and opportunities. This style is characterized by a **focus on human potential, resilience, and the capacity for positive change...**

The **economic** debating style frames arguments in terms of cost-benefit analysis, resource allocation, and economic impact. This style is characterized by a **focus on incentives, trade-offs, and opportunity costs....**

The **technological** debating style focuses on the role of innovation, digital transformation, and scientific advancements in shaping solutions. This style is characterized by a **focus on disruption, efficiency gains, and the transformative power of technology...**

The **contrarian** debating style consistently takes positions opposite to the mainstream view or prevailing wisdom. This style is characterized by a **willingness to stand apart from the crowd, a skepticism towards widely accepted ideas...**

# LLM Evaluation

- Specify the grading notes to judge responses
- Each number provided matches a list of criteria
- Fine tuned based on human evaluation of scoring quality

 databricks

Blog

Enhancing LLM-as-a-Judge with Grading Notes

Assistant LLM	Judge Method				
	Human	GPT-4	GPT-4 + Grading Notes	GPT-4-Turbo	GPT-4-Turbo + Grading Notes
Positive Label Rate by Judge					
Llama3-70b	71.9%	96.9%	73.1%	83.1%	65.6%
GPT-4o	79.4%	98.1%	81.3%	91.9%	68.8%
Alignment Rate with Human Judge					
Llama3-70b	-	74.7%	96.3%	76.3%	91.3%
GPT-4o	-	76.8%	93.1%	77.5%	84.4%

Respect	1-5
Accurate Information	1-5
Relevance	1-5
Argument Quality	1-5
Critical Thinking	1-5
Organization	1-5
Preparation	1-5

Demo time!

<div data-bbox="49 290 86 562" data-label="Text">TODAY'S DATE</div> <div data-bbox="49 1168 86 1444" data-label="Text">PROJECT NAME</div>	<div data-bbox="209 71 606 118" data-label="Text">KEY FINDINGS</div> <div data-bbox="349 365 2002 487" data-label="Section-Header"> <h1>Next steps: Fine-tuning Styles</h1> </div> <div data-bbox="349 643 2019 765" data-label="Text"> <p>With our validated response dataset, the model can fine tune for the next iteration</p> </div> <div data-bbox="349 943 1012 999" data-label="Text"> <p><b>Get Highest Score Agents</b></p> </div> <div data-bbox="349 1114 1006 1170" data-label="Text"> <p><b>Integrate Past Responses</b></p> </div> <div data-bbox="349 1285 1079 1341" data-label="Text"> <p><b>Update Style Generator LLM</b></p> </div> <div data-bbox="2568 851 3072 1361" data-label="Image"> <p>A decorative graphic consisting of a 3x3 grid of nine rounded squares. The colors of the squares are: top row (orange, yellow, teal), middle row (purple, dark grey, purple), and bottom row (teal, yellow, orange). The squares are slightly offset from each other, creating a dynamic, non-uniform pattern.</p> </div>



Ask your  
questions

# Questions?