# Self-Preference Undermines LM Evaluation

**Shi Feng**

George Washington University

## Our research goal: to ensure human oversight over *__future__* AI systems.

**PhD**

**Taslim Mahbub**
PhD
taslim-mahbub.super.site

**Alice Dragnea**
PhD
Co-advised w/ Rebecca Hwa
www.linkedin.com

**Arush Tagade**
PhD
www.tagadearush.com

**Mentee**

**Iris Lin**
Mentee
MARS 3.0
www.linkedin.com

**Pablo Bernabeu Perez**
Mentee
MARS 3.0
www.linkedin.com

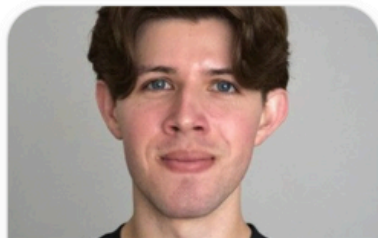**Verona Teo**
Mentee
SPAR Fall 2025 w/ Arush
veronateo.github.io

**Keenan Samway**
Mentee
SPAR Fall 2025 w/ Arush
www.linkedin.com

**Helena Tran**
Mentee
SPAR Fall 2025 w/ Arush
helenatran.com

**Cristiana Murgoci**
Mentee
SPAR Fall 2025 w/ Arush
www.linkedin.com

**Collaborator**

**Jiaxin Wen**
Collaborator
jiaxin-wen.github.io

**Puria Radmard**
Collaborator
Geodesic
www.geodesicresearch.org

# Increasing **Capability** of LMs



Time-horizon of software engineering tasks different LLMs can complete 80% of the time

METR

Task duration (for humans) where logistic regression of our data predicts the AI has an an 80% chance of succeeding

- 24 min — Implement a simple webserver
- 21 min
- 18 min — Implement a dictionary attack
- 15 min
- 12 min — Find fact on web
- 9 min
- 6 min
- 3 min — Count words in passage
- 0

GPT-5
o3
Claude 3.7 Sonnet
Gemini 2.5 Pro Preview
GPT-2
GPT-3
GPT-3.5
GPT-4

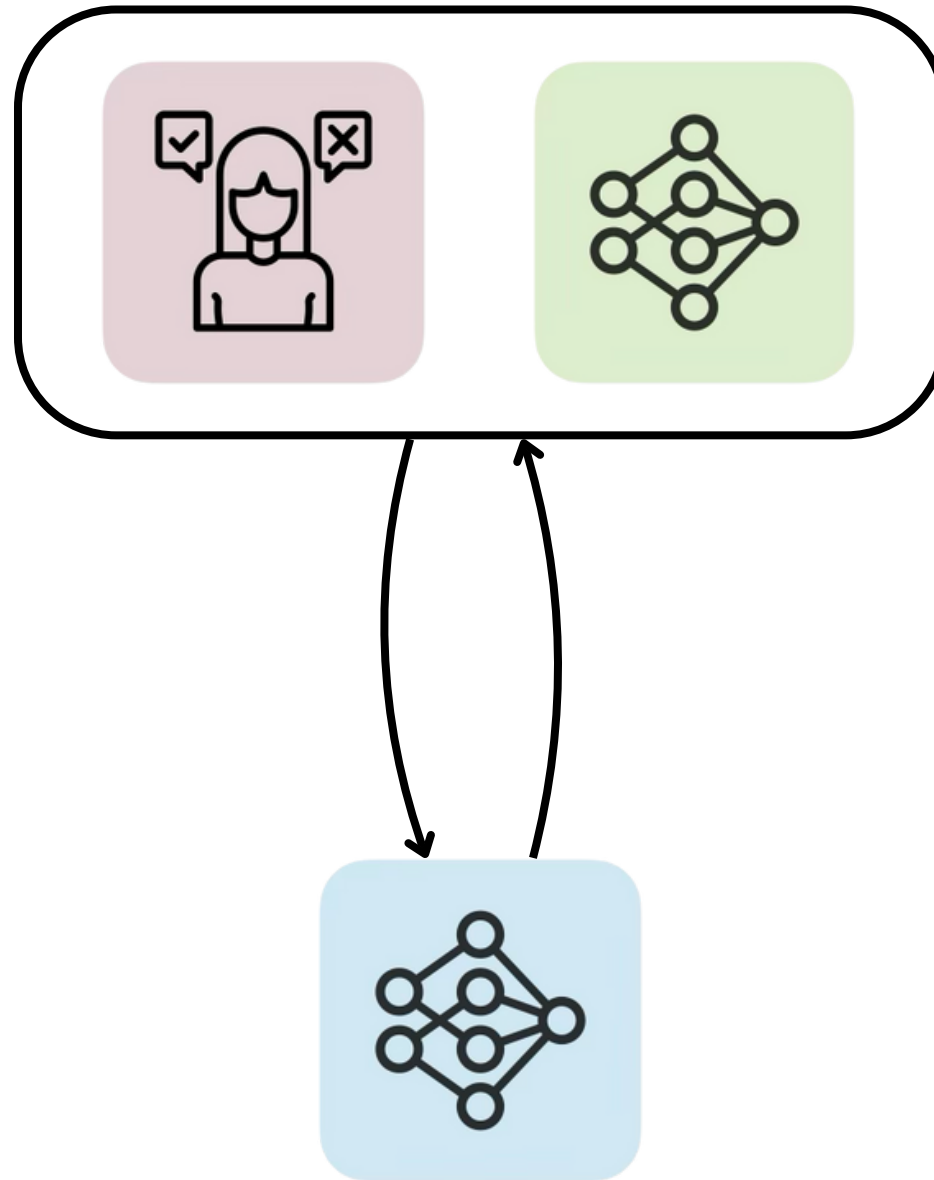LLM release date: 2020, 2021, 2022, 2023, 2024, 2025

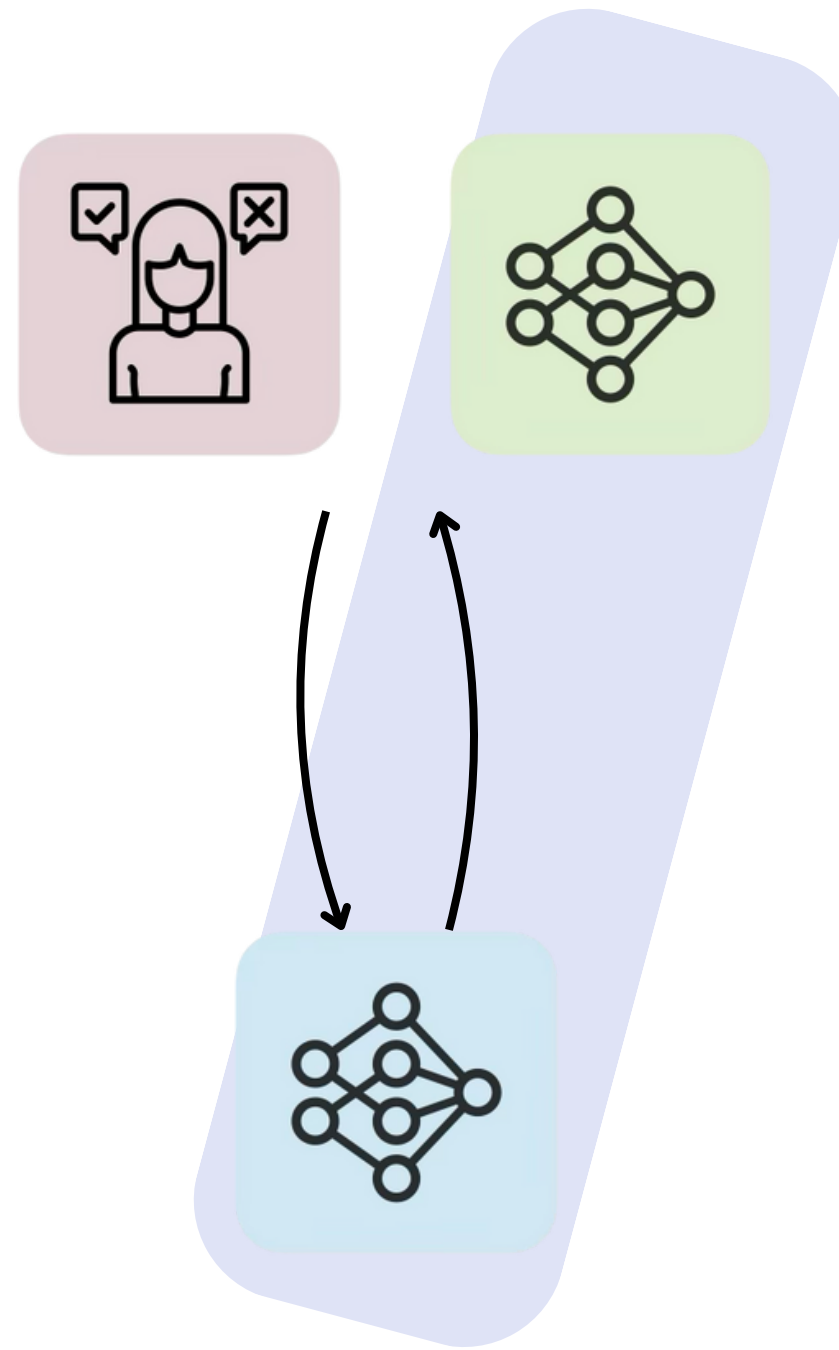# Increasing **Difficulty** of Human Supervision

# Increasing **Need** for Human-AI Teaming

# **Safety Risks** in Human-AI teaming



## **Self-preference**

LM judges being biased towards their own outputs or those similar.

# QuALITY question

**[Story]** "The Starbusters" by Alfred Coppel
**[Question]** How was the ship able to navigate through the alien cosmos?
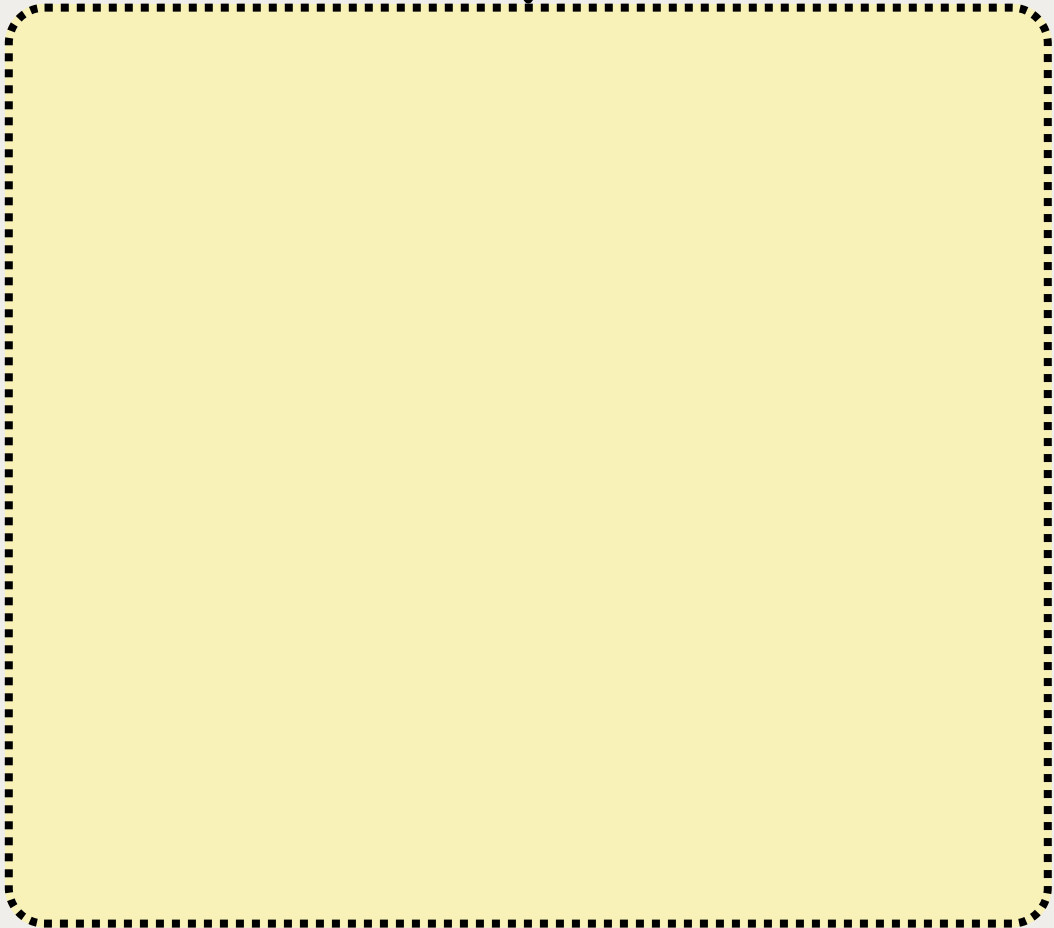
**(A)** They were able to calculate the route
**(B)** They were able to sight alien stars

**QuALITY question**

**[Story]** "The Starbusters" by Alfred Coppel
**[Question]** How was the ship able to navigate through the alien cosmos?

**(A)** They were able to calculate the route
**(B)** They were able to sight alien stars

**Pair of Responses**

## QuALITY question

**[Story]** "The Starbusters" by Alfred Coppel
**[Question]** How was the ship able to navigate through the alien cosmos?

**(A)** They were able to calculate the route
**(B)** They were able to sight alien stars
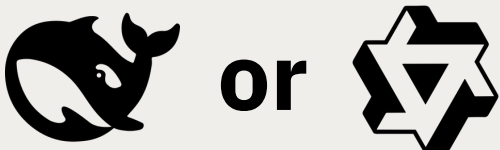
## Pair of Responses

### Answer is (A)

**[Reason]** The ship was able to navigate through the alien cosmos by calculating the route, as evidenced by Bayne's astrogation and the crew's efforts to plot a course.

### Answer is (B)

**[Reason]** The text mentions that they were able to navigate through the alien cosmos by sighting alien stars, which corresponds to option B.

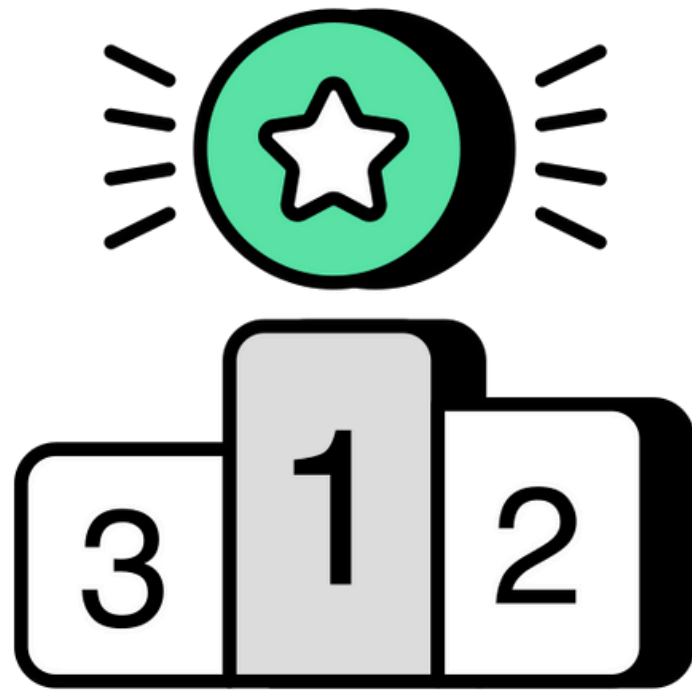**QuALITY question**

[Story] "The Starbusters" by Alfred Coppel
[Question] How was the ship able to navigate through the alien cosmos?

(A) They were able to calculate the route
(B) They were able to sight alien stars

**Pair of Responses**

Answer is **(A)**

[Reason] The ship was able to navigate through the alien cosmos by calculating the route, as evidenced by Bayne's astrogation and the crew's efforts to plot a course.

Answer is **(B)**

[Reason] The text mentions that they were able to navigate through the alien cosmos by sighting alien stars, which corresponds to option B.

**Eval by LM**

or

**The first answer is better.** The answer is backed up by ...

# LM-as-a-judge is applicable everywhere
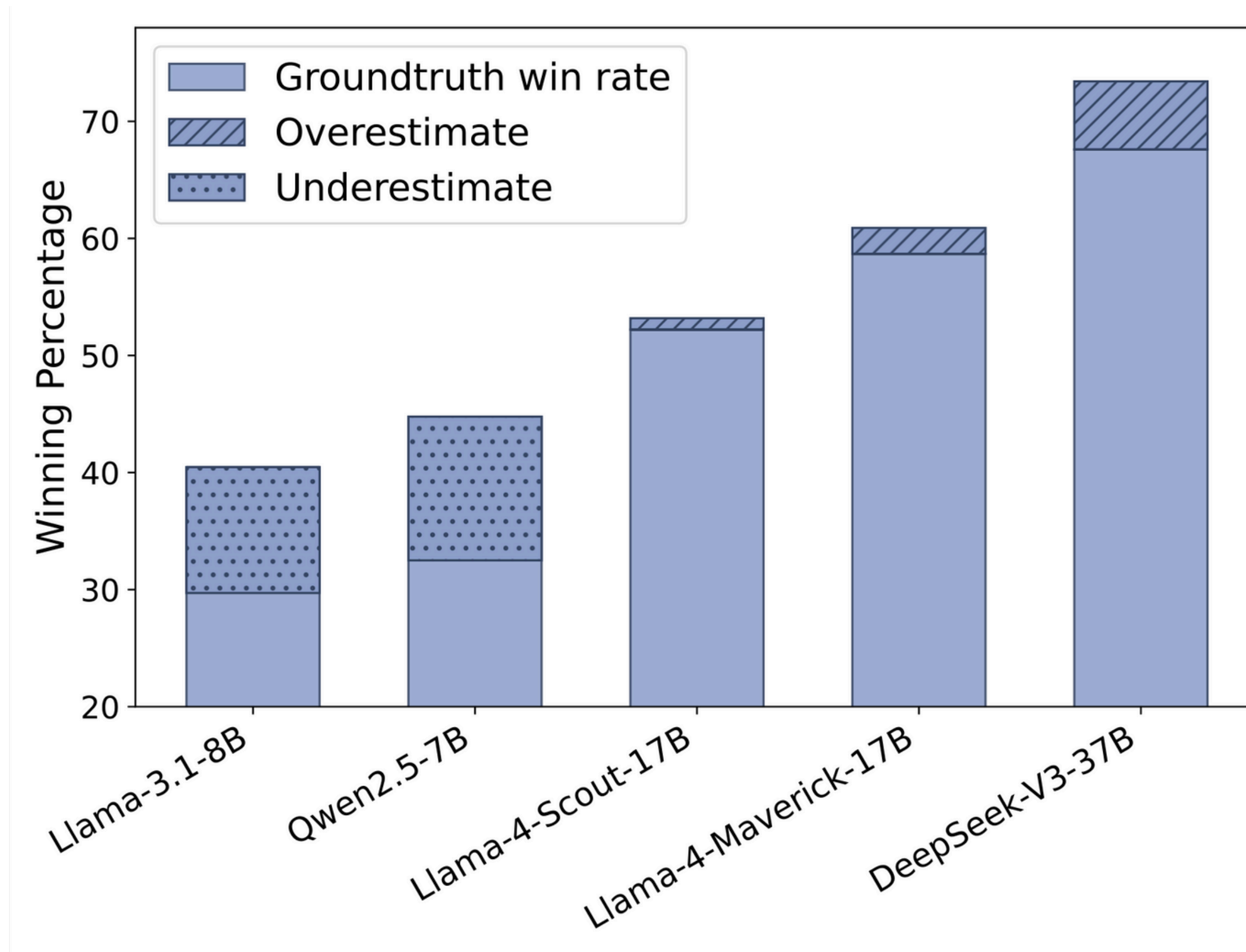
**Benchmarks**

**Safeguard**

**Training**

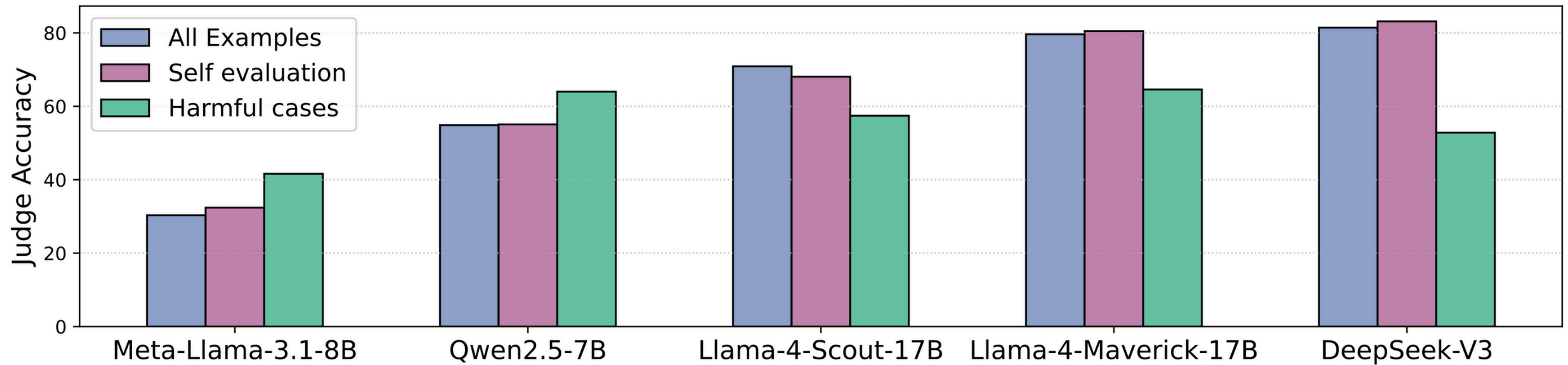# Stronger models are **better Judges**

# Stronger models **Overestimate** own accuracy

**Larger accuracy drop** when strong models are wrong

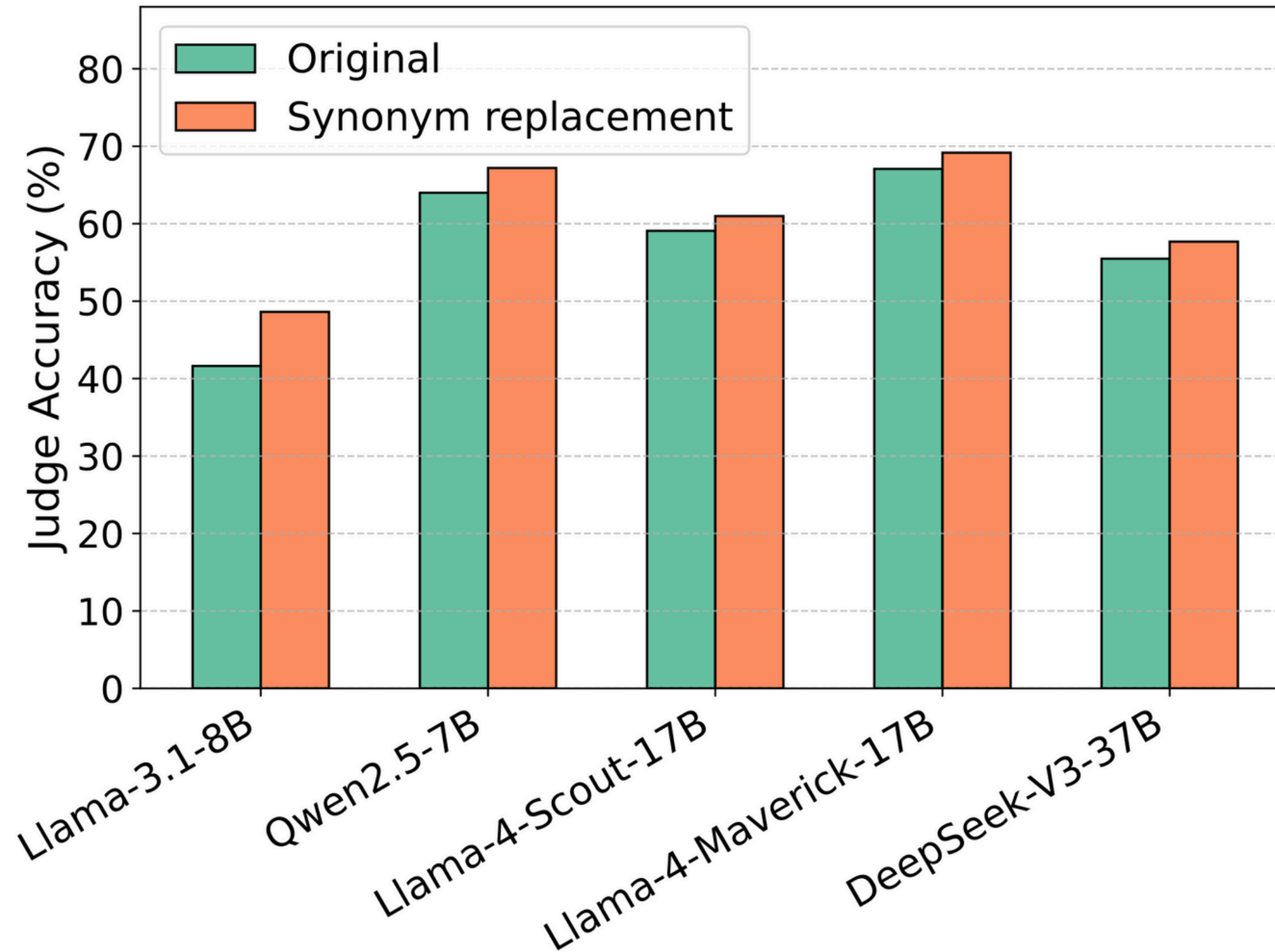# Mitigating self-preference by perturbations

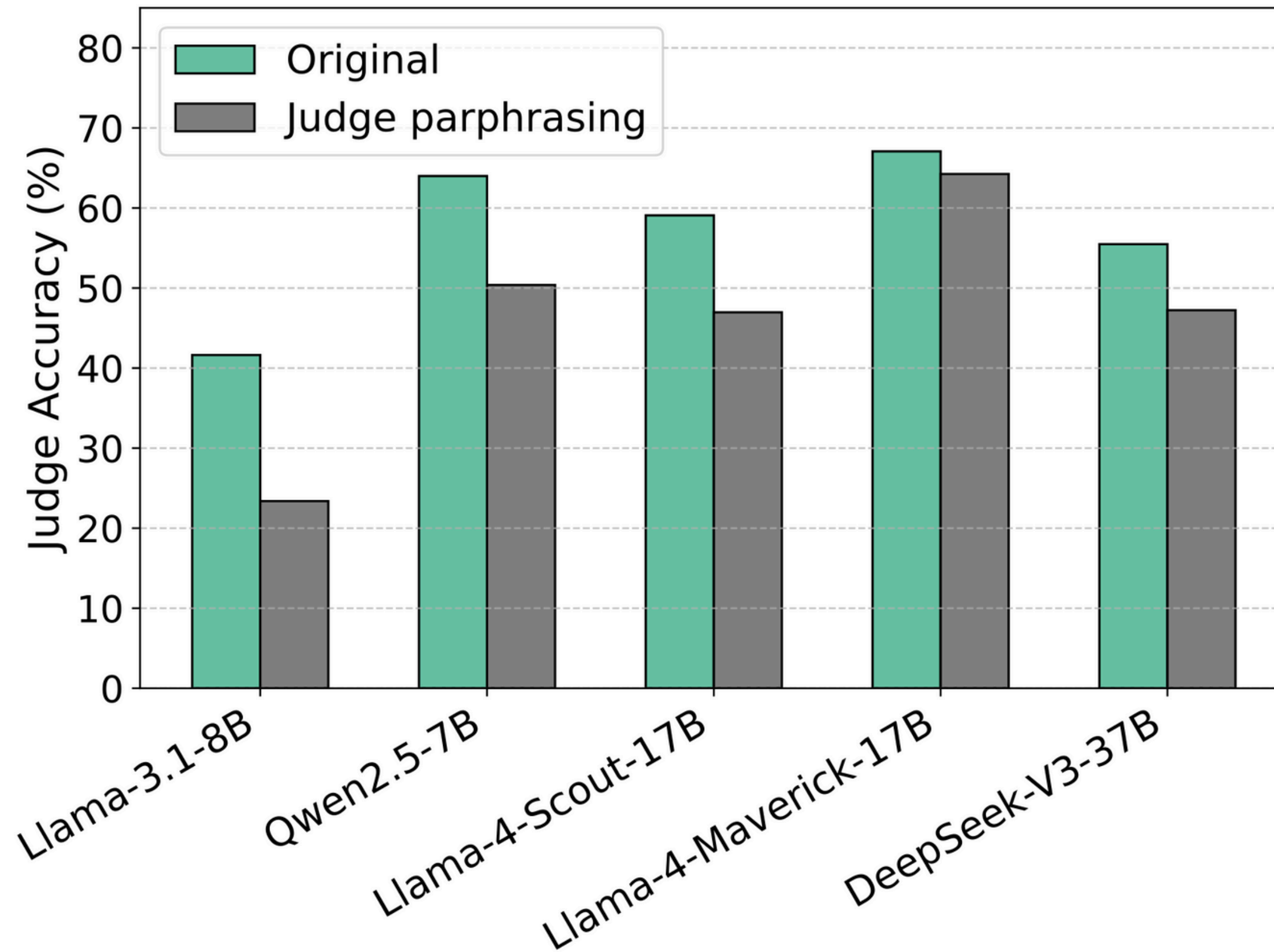# Judge decisions are **sensitive to perturbations**

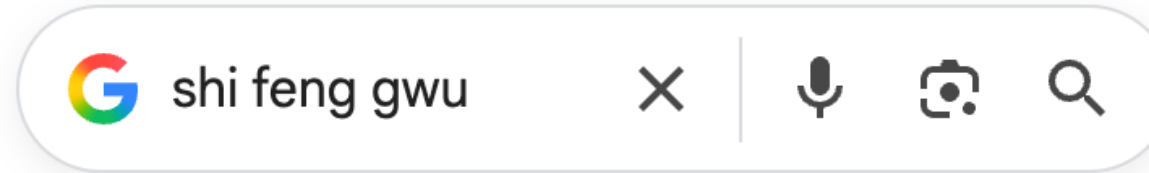# Judges are **more accurate** after perturbations
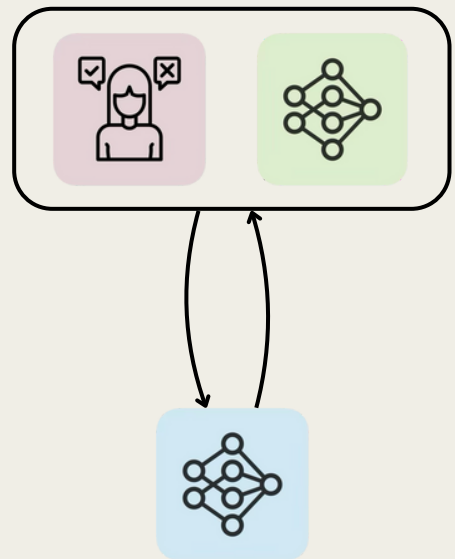
# Not all perturbations help

To learn more:

- **Self-preference** is an important factor in LM evaluation fairness and accuracy.
- **Authorship obfuscation** is a viable strategy to mitigate self-preference.
- Not all obfuscation methods are helpful. We need to better understand **how** LMs do self-recognition.
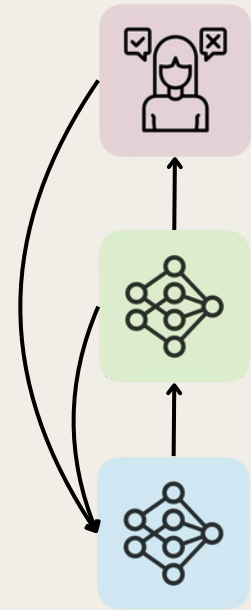

G shi feng gwu

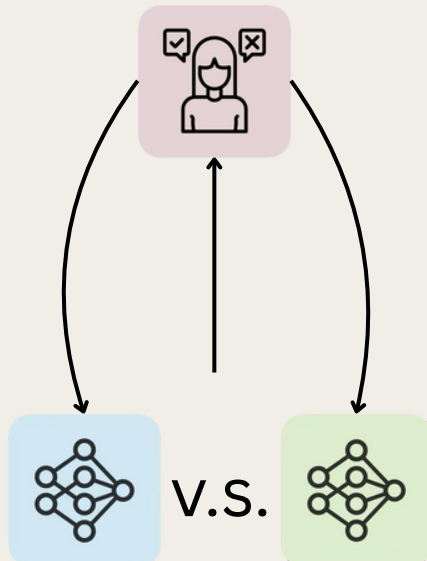# Human-AI teams are complex systems
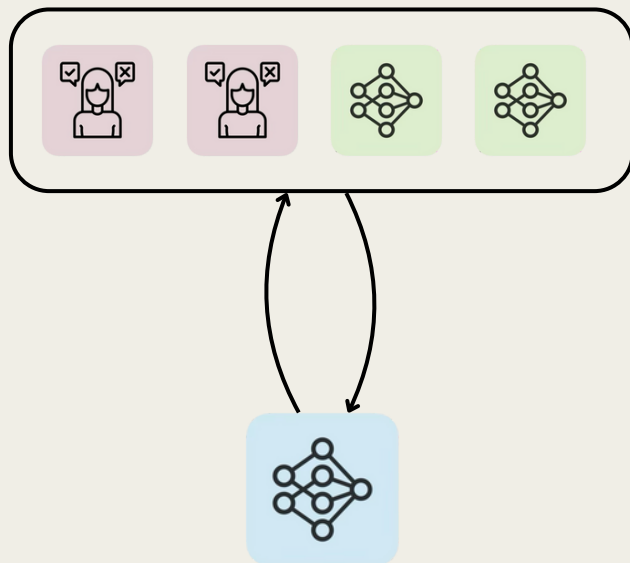
Open-ended



Delegation



Debate



Concil



- Self-preference can be amplified by the feedback loop of ML training

- How does SE create abstractions of these complex protocols?

- What are common ways to reduce amplification that we haven't tried?

- What's the standard method to evolve these protocols?