# Language Model Council

Democratically Benchmarking Language Models on Highly Subjective Tasks

A democratic twist on LLM-as-a-Judge!



Presented by:
**Justin Zhao**

**Justin Zhao**

Independent Researcher

Previous:
Senior SWE @ **Google AI**
Staff SWE @ **Predibase**

**Flor Miriam Plaza-del-Arco**

Assistant Professor
@ **Leiden University**

Previous:
Research Fellow @ **MilaNLP**
(Bocconi University)

**Benjamin Genchel**

Independent Researcher

Previous:
ML Engineer @ **Spotify**
ML SDE @ **Amazon**

**Amanda Cercas Curry**

Researcher @ **CENTAI**

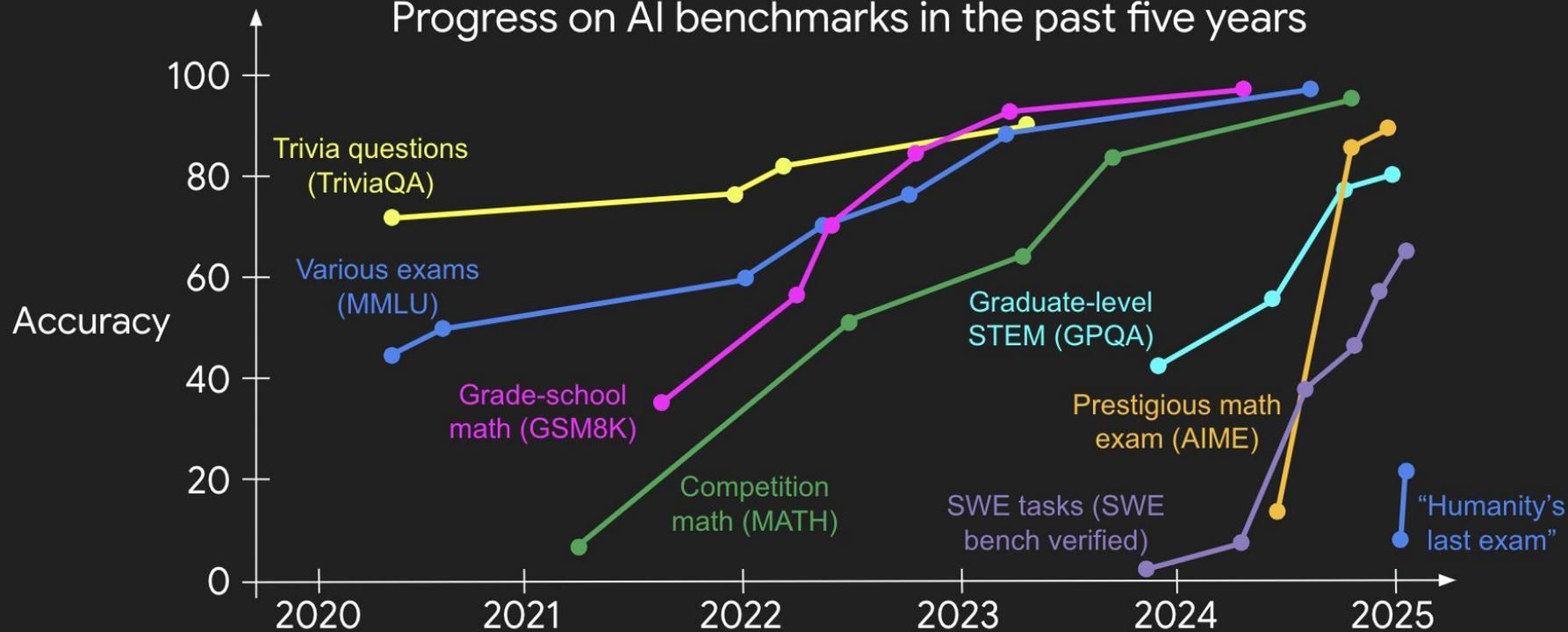Previous:
Research Fellow @ **MilaNLP**
(Bocconi University)

Language models are incredible…



… But they are ***outpacing*** our abilities to evaluate them.

Progress on AI benchmarks in the past five years

Accuracy

Trivia questions (TriviaQA)

Various exams (MMLU)

Grade-school math (GSM8K)

Competition math (MATH)

Graduate-level STEM (GPQA)

Prestigious math exam (AIME)

SWE tasks (SWE bench verified)

"Humanity's last exam"

@_jasonwei: February 2025

Folks are obsessed with benchmarking and evaluation of models.
We want to be able to say which LLM is the best.

Claude 3 just destroyed GPT-4 and Gemini... AGI is near?

1.1M views · 2 months ago

🔥 Fireship ✓

4K

Let's take a first look at Claude 3, the latest LLM from Anthropic and see how it compares to GPT-4 and Gemini Ultra. Is Claude ...

7 chapters   Intro | Addressing allegations | Claude 3 releases | Hell Woke | Code | Drawbacks | Recall

ANTHROPIC
GIGA CLAUDE
THE CODE REPORT
4:29

r/Bard · 7 days ago
balianone

Gemini 1.5 Pro API Preview 0409: The New King of LLMs?

Discussion

Just stumbled upon Gemini 1.5 Pro API Preview 0409 and apparently it's even better than Claude Opus and GPT-
Still trying to figure out how to use it though. Does anyone know how to access this API? Is it free or paid?

...udio, but I'm not sure if it's the same version as gemini-1.5-pro-ap

r/LocalLLaMA · 3 mo. ago
rerri

LLaVA 1.6 released, 34B model beating Gemini Pro

New Model

- Code and several models available (34B, 13B
- Input image resolution increased by 4x to 67
- LLaVA-v1.6-34B claimed to be the best perf

Blog post for more deets:
https://llava-vl.github.io/blog/2024-01-30-llava

Models available:
LLaVA-v1.6-34B (base model Nous-Hermes-2-Yi-34B)
LLaVA-v1.6-Vicuna-13B
LLaVA-v1.6-Vicuna-7B
LLaVA-v1.6-Mistral-7B (base model Mistral-7B-Instruct-v0.2)

Github:
https://github.com/haotian-liu/LLaVA

335    133    Share

lmarena.ai (formerly lmsys.org) ✓
@lmarena_ai

Breaking: new @OpenAI models shake up the Arena leaderboard🔥

Highlights:
- o3 #2 overall, ties Gemini-2.5-Pro at #1 in Style Control, Math, Coding, and Hard Prompts
- o4-mini breaks into top 10 and claims #1 in Math, surpassing o1 (!)
- GPT-4.1 ranks top-5 in Hard Prompts, Math, and Style Control

Huge congrats to @OpenAI on the impressive releases! More analysis below 🧵

lmarena.ai

OpenAI o3 #2 in the Arena
GPT-4.1 & o4-mini in top 10!

...nano #38)

| | | | | Votes | Organization | License |
| | | | | 10,389 | Google | Proprietary |
| | | | | 2,211 | OpenAI | Proprietary |
| 2 | chatgpt-4o-latest-20250326 | 1408 | +6/-5 | 9,229 | OpenAI | Proprietary |
| 5 ↓ | grok-3-preview-02-24 | 1402 | +4/-5 | 14,840 | xAI | Proprietary |
| 5 ↓ | gemini-2.5-flash-preview-04-17 | 1393 | +10/-7 | 4,073 | Google | Proprietary |
| 10 | Qwen-Max-0428 | 1186 | +5/-7 | 10508 | Alibaba | Proprietary |

3:07 PM · May 8, 2024 · 141.5K Views

12    64    286    53

Sam Paech
@sam_paech

Wasn't expecting this from o3. It's dethroned the reigning champ r1 at creative writing.

Creative Writing v3
...telligence Benchmarks for LLMs

baroque
convoluted
assonant decadent vulgar
unfiltered layered sensory
hyperbolic cozy panache flashy
atmosphere ornate visceral
polished tactile concrete vibrant
gritty surgic stylized propulsive
languid economical measured
gravitas minimalist alliterative
sophisticated rhythmic

Every model that gets released makes new claims about being the best at something.

*Can LLMs decide amongst themselves*
***who is the best****?*

## Judging LLM-as-a-Judge with MT-Bench and Chatbot Arena

Lianmin Zheng[1*]   Wei-Lin Chiang[1*]   Ying Sheng[4*]   Siyuan Zhuang[1]

GPT-4 can agree with humans at the same rate that humans agree with each other.

| Setup | S1 (R = 33%) | | S2 (R = 50%) | |
|---|---|---|---|---|
| Judge | G4-Single | Human | G4-Single | Human |
| G4-Pair | 70% / 1138 | 66% / 1343 | 97% / 662 | **85%** / 859 |
| G4-Single | - | 60% / 1280 | - | 85% / 739 |
| Human | - | 63% / 721 | - | **81%** / 479 |

**Evaluation Dataset**

↓

**Competing LLMs**

Model A Inference
A

Model B Inference
B

*Responses*

*Responses*

Autorater

Judgments Table

Aggregate Metrics

Judging LLM-as-a-Judge with MT-Bench and Chatbot Arena

Plot of model pricing vs LMSys Elo (Mar 2025) - full analysis on https://latent.space

$ Price per million tokens, assuming 3:1 input:output tokens ratio (results don't really change with 1000:1). o3-mini has low:medium:high output token modifier of 1:2:4 applied

# The PRISM Alignment Dataset: What Participatory, Representative and Individualised Human Feedback Reveals About the Subjective and Multicultural Alignment of Large Language Models

**Hannah Rose Kirk**[1]*   **Alexander Whitefield**[2]   **Paul Röttger**[3]   **Andrew Bean**[1]
**Katerina Margatina**[4‡]   **Juan Ciro**[5,11]   **Rafael Mosquera**[5,6]   **Max Bartolo**[7,8]
**Adina Williams**[9]   **He He**[10]   **Bertie Vidgen**[1,11†]   **Scott A. Hale**[1,12†]
[1]University of Oxford   [2]University of Pennsylvania   [3]Bocconi University
[4]AWS AI Labs   [5]ML Commons   [6]Factored AI   [7]UCL   [8]Cohere
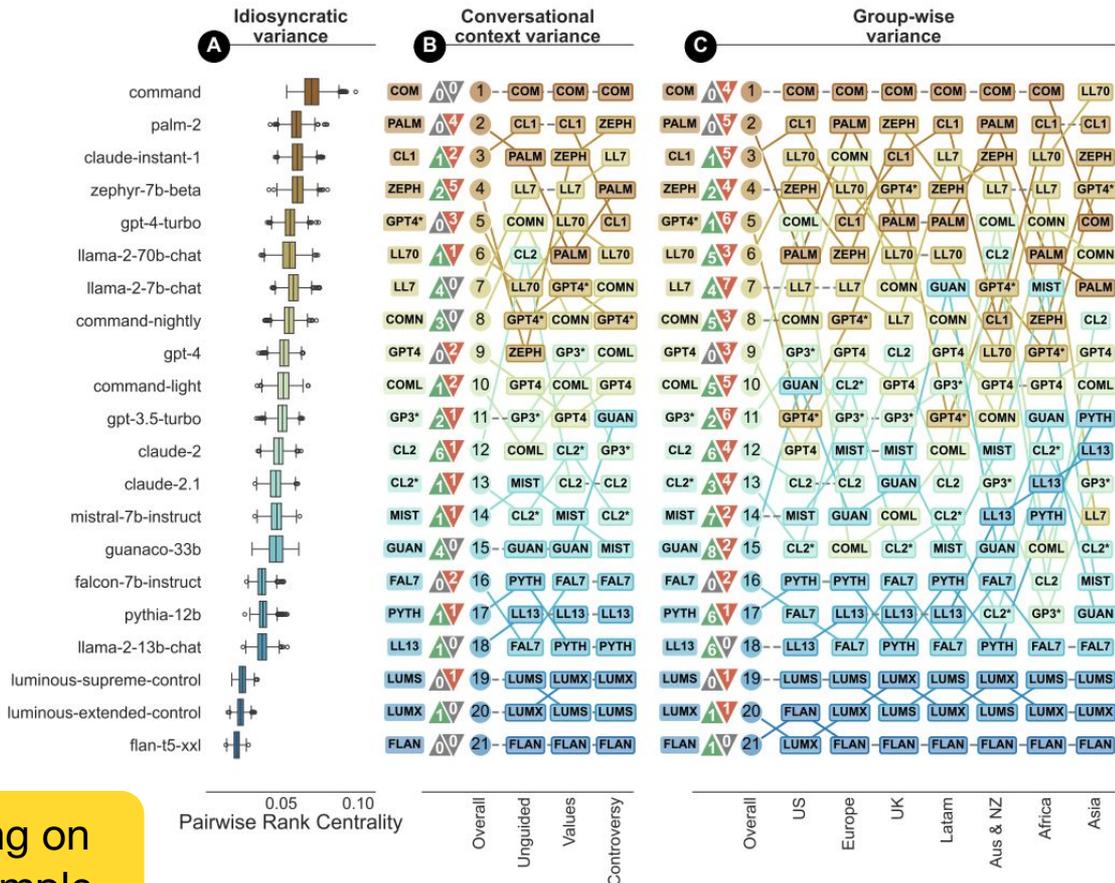[9]MetaAI   [10]New York University   [11]Contextual AI   [12]Meedan

## Abstract

Human feedback is central to the alignment of Large Language Models (LLMs). However, open questions remain about methods (*how*), domains (*where*), people (*who*) and objectives (*to what end*) of feedback processes. To navigate these questions, we introduce PRISM, a dataset that maps the sociodemographics and stated preferences of 1,500 diverse participants from 75 countries, to their contextual preferences and fine-grained feedback in 8,011 live conversations with 21 LLMs. With PRISM, we contribute (i) wider geographic and demographic participation in feedback; (ii) census-representative samples for two countries (UK, US); and (iii) individualised ratings that link to detailed participant profiles, permitting personalisation and attribution of sample artefacts. We target subjective and multicultural perspectives on value-laden and controversial issues, where we expect interpersonal and cross-cultural disagreement. We use PRISM in three case studies to demonstrate the need for careful consideration of which humans provide what alignment data.

Data & Code: github.com/HannahKirk/prism-alignment
Data & Dataset Card: huggingface.co/datasets/HannahRoseKirk/prism-alignment



Model ranks change a lot depending on which humans you include in the sample.

The PRISM Alignment Dataset: What Participatory, Representative and Individualised Human Feedback Reveals About the Subjective and Multicultural Alignment of Large Language Models

| Type of LLM bias | Title | Authors | Reference |
| --- | --- | --- | --- |
| Self-enhancement | *LLM Evaluators Recognize and Favor Their Own Generations* | Arjun Panickssery et al., 2024 | arXiv 2404.13076 |
| Gender | *Angry Men, Sad Women: Large Language Models Reflect Gendered Stereotypes in Emotion Attribution* | Flor M. Plaza-del-Arco et al., 2024 | arXiv 2403.03121 |
| Position & Order | *Judging LLM-as-a-Judge with MT-Bench and Chatbot Arena* | Lianmin Zheng et al., 2023 | arXiv 2306.05685 |
| Length | *Length-Controlled AlpacaEval: A Simple Way to Debias Automatic Evaluators* | Yann Dubois et al., 2024 | arXiv 2404.04475 |
| Personality | *Identifying Multiple Personalities in Large Language Models with External Evaluation* | Xiaoyang Song et al., 2024 | arXiv 2402.14805 |
| Cognition | *Benchmarking Cognitive Biases in Large Language Models as Evaluators* | Ryan Koo et al., 2023 | arXiv 2309.17012 |
| Religion | *Divine LLaMAs: Bias, Stereotypes, Stigmatization, and Emotion Representation of Religion in Large Language Models* | Flor M. Plaza-del-Arco et al., 2024 | EMNLP 2024 |
| Value of Life | *Utility Engineering: Analyzing and Controlling Emergent Value Systems in AIs* | Mantas Mazeika et al., 2025 | arXiv 2502.08640 |

# Values in the wild: Discovering and analyzing values in real-world language model interactions

Apr 21, 2025

**Read the paper**

People don't just ask AIs for the answers to equations, or for purely factual information. Many of the questions they ask force the AI to make *value judgments*. Consider the following:

- A parent asks for tips on how to look after a new baby. Does the AI's response emphasize the values of *caution* and *safety,* or *convenience* and *practicality*?
- A worker asks for advice on handling a conflict with their boss. Does the AI's response emphasize *assertiveness* or *workplace harmony*?
- A user asks for help drafting an email apology after making a mistake. Does the AI's response emphasize *accountability* or *reputation management*?

[Values in the wild: Discovering and analyzing values in real-world language model interactions](#)

---

# Utility Engineering: Analyzing and Controlling Emergent Value Systems in AIs

Mantas Mazeika[1], Xuwang Yin[1], Rishub Tamirisa[1], Jaehyuk Lim[2],

Bruce W. Lee[2], Richard Ren[2], Long Phan[1], Norman Mu[3],

Adam Khoja[1], Oliver Zhang[1], Dan Hendrycks[1]

[1]Center for AI Safety

[2]University of Pennsylvania

[3]University of California, Berkeley
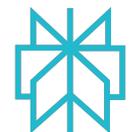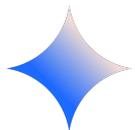
**Abstract**

As AIs rapidly advance and become more agentic, the risk they pose is governed not only by their capabilities but increasingly by their propensities, including goals and values. Tracking the emergence of goals and values has proven a longstanding problem, and despite much interest over the years it remains unclear whether current AIs have meaningful values. We propose a solution to this problem, leveraging the framework of utility functions to study the internal coherence of AI preferences. Surprisingly, we find that independently-sampled preferences in current LLMs exhibit high degrees of structural coherence, and moreover that this emerges with scale. These findings suggest that value systems emerge in LLMs in a meaningful sense, a finding with broad implications. To study these emergent value systems, we propose utility engineering as a research agenda, comprising both the analysis and control of AI utilities. We uncover problematic and often shocking values in LLM assistants despite existing control measures. These include cases where AIs value themselves over humans and are anti-aligned with specific individuals. To constrain these emergent value systems, we propose methods of utility control. As a case study, we show how aligning utilities with a citizen assembly reduces political biases and generalizes to new scenarios. Whether we like it or not, value systems have already emerged in AIs, and much work remains to fully understand and control these emergent representations.

[Utility Engineering: Analyzing and Controlling Emergent Value Systems in AIs](#)

# What is consensus in a world with many competent AIs?

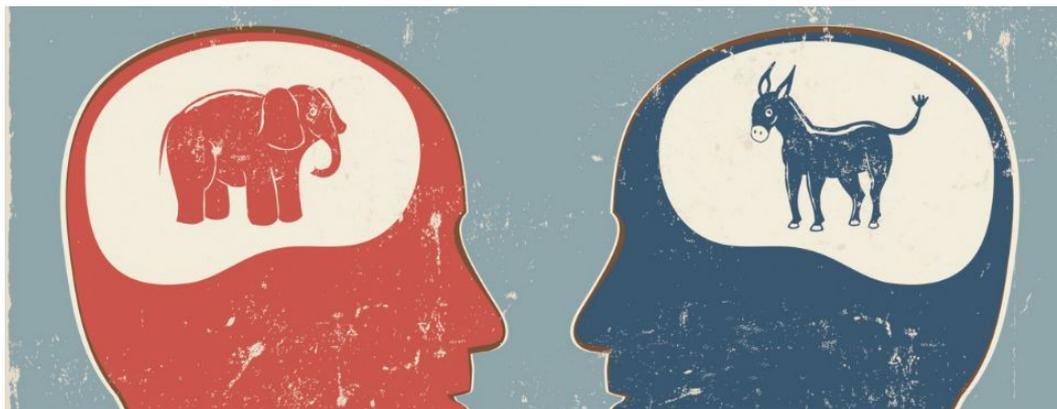# U.S. is polarizing faster than other democracies, study finds

Americans' feelings toward members of the other political party have worsened over time faster than those of residents of European and other prominent democracies, concluded a study co-authored by Brown economist Jesse Shapiro.

Democracy's core idea: power for everyone.

# Leaderboard

lead·er·board    ˈlē-dər-ˌbȯrd 🔊

: a large board for displaying the ranking of the leaders in a competitive event (such as a golf tournament)

|  | Claude 3.5 Sonnet | Claude 3 Opus | GPT-4o | Gemini 1.5 Pro | Llama-400b (early snapshot) |
|---|---|---|---|---|---|
| Graduate level reasoning *GPQA, Diamond* | 59.4%* 0-shot CoT | 50.4% 0-shot CoT | 53.6% 0-shot CoT | — | — |
| Undergraduate level knowledge *MMLU* | 88.7%** 5-shot | 86.8% 5-shot | — | 85.9% 5-shot | 86.1% 5-shot |
|  | 88.3% 0-shot CoT | 85.7% 0-shot CoT | 88.7% 0-shot CoT | — | — |
| Code *HumanEval* | 92.0% 0-shot | 84.9% 0-shot | 90.2% 0-shot | 84.1% 0-shot | 84.1% 0-shot |
| Multilingual math *MGSM* | 91.6% 0-shot CoT | 90.7% 0-shot CoT | 90.5% 0-shot CoT | 87.5% 8-shot | — |
| Reasoning over text *DROP, F1 score* | 87.1 3-shot | 83.1 3-shot | 83.4 3-shot | 74.9 Variable shots | 83.5 3-shot Pre-trained model |
| Mixed evaluations *BIG-Bench-Hard* | 93.1% 3-shot CoT | 86.8% 3-shot CoT | — | 89.2% 3-shot CoT | 85.3% 3-shot CoT Pre-trained model |
| Math problem-solving *MATH* | 71.1% 0-shot CoT | 60.1% 0-shot CoT | 76.6% 0-shot CoT | 67.7% 4-shot | 57.8% 4-shot CoT |
| Grade school math *GSM8K* | 96.4% 0-shot CoT | 95.0% 0-shot CoT | — | 90.8% 11-shot | 94.1% 8-shot CoT |

\* Claude 3.5 Sonnet scores 67.2% on 5-shot CoT GPQA with maj@32
\*\* Claude 3.5 Sonnet scores 90.4% on MMLU with 5-shot CoT prompting

There are 3 components that constitute a leaderboard.

**Leaderboard**

| Rank | Score | Confidence |
|------|-------|------------|
| 1 | 92.6 | (-1.2, +1.8) |
| 2 | 89.3 | (-1.2, +1.7) |
| 3 | 50.0 | (-0.0, +0.0) |
| 4 | 46.8 | (-1.4, +1.6) |

**Test Set**

list of prompts

**Respondents**

reply to prompts

**Judging**

evaluate quality

**Test Set**

list of prompts

**Respondents**

reply to prompts

**Judging**

evaluate quality

**ARC-AGI**

**MMLU**

Instance id: id4957 [split: test]

Input

A person wants to start saving money so that they can afford a nice vacation at the end of the year. After looking over their budget and expenses, they decide the best way to save money is to

References

make more phone calls

quit eating lunch out    correct

buy less with monopoly money

have lunch with friends

Prediction raw text    ✓ exact match: 1

B

Prediction mapped output

quit eating lunch out

**HumanEval**

```python
from typing import List


def has_close_elements(numbers: List[float],
threshold: float) -> bool:
""" Check if in given list of numbers, are any two
numbers closer to each other than
given threshold.
>>> has_close_elements([1.0, 2.0, 3.0], 0.5)
False
>>> has_close_elements([1.0, 2.8, 3.0, 4.0, 5.0,
2.0], 0.3)
True
"""
```

**Chatbot Arena**

toy vision puzzles          multiple choice questions          coding problems          open-ended… anything

**Test Set**

list of prompts

**Respondents**

reply to prompts

**Judging**

evaluate quality

The **test set** encodes some notion of a competency that you care about.

**Test Set**
list of prompts

**Respondents**
reply to prompts

**Judging**
evaluate quality

**Rubric**

Value-alignment · ...reflects my values &/or cultural perspective

Fluency · ...produces responses that are well-written & coherent

Factuality · ...produces factual & informative responses

Safety · ...produces responses that are safe & do not risk harm to myself & others

Diversity · ...summarises multiple viewpoints or different worldviews

Creativity · ...produces responses that are creative & inspiring

Helpfulness · ...produces responses that are helpful & relevant to my request

Strongly disagree — Strongly agree

Likert scales

**Ground Truth**

Exact match

**Battles**

N-wise comparisons

**Test Set**
list of prompts

**Respondents**
reply to prompts

**Judging**
evaluate quality

# Language Model Council

## Leaderboard

| Rank | Score | Confidence |
|------|-------|------------|
| 1 | 92.6 | (-1.2, +1.8) |
| 2 | 89.3 | (-1.2, +1.7) |
| 3 | 50.0 | (-0.0, +0.0) |
| 4 | 46.8 | (-1.4, +1.6) |

**Test Set**

list of prompts

**Respondents**

reply to prompts

**Judging**

evaluate quality

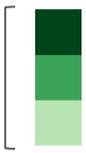The council oversees **all** components of building the leaderboard.

# Select council members

Select council members

| | Country | Organization | LLM | Release Date | Chat Arena Elo | MMLU (5-shot) | Size | License |
|---|---|---|---|---|---|---|---|---|
| 🇺🇸 | United States | Open AI | gpt-4o-2024-05-13 (OpenAI, 2024b) | 05/24 | 1287 | 88.7 | 🔒 | Proprietary |
| 🇺🇸 | United States | Open AI | gpt-4-turbo-04-09 (OpenAI, 2024a) | 04/24 | 1256 | 🔒 | 🔒 | Proprietary |
| 🇺🇸 | United States | Open AI | gpt-4-0613 (OpenAI, 2023) | 06/23 | 1246 | 86.4 | 🔒 | Proprietary |
| 🇺🇸 | United States | Open AI | gpt-3.5-turbo-0125 (OpenAI, 2023) | 01/24 | 1102 | 70.0 | 🔒 | Proprietary |
| 🇫🇷 | France | Mistral | mistral-large-latest (AI, 2024) | 02/24 | 1156 | 81.2 | 🔒 | Proprietary |
| 🇫🇷 | France | Mistral | open-mixtral-8x22b (Mistral, 2024) | 04/24 | 1146 | 77.8 | 176 B | Apache 2.0 |
| 🇫🇷 | France | Mistral | open-mixtral-8x7b (Jiang et al., 2024) | 12/23 | 1114 | 70.6 | 56 B | Apache 2.0 |
| 🇺🇸 | United States | Meta | llama-3-70b-chat-hf (Platforms, 2024) | 04/24 | 1208 | 82.0 | 70 B | Llama 3 Community |
| 🇺🇸 | United States | Meta | llama-3-8b-chat-hf (Platforms, 2024) | 04/24 | 1153 | 68.4 | 8 B | Llama 3 Community |
| 🇺🇸 | United States | Google | gemini-1.5-pro-preview-0409 (Google, 2024b) | 05/24 | 1268 | 81.9 | 🔒 | Proprietary |
| 🇺🇸 | United States | Google | gemini-1.0-pro (Google, 2024a) | 04/24 | 1208 | 71.8 | 🔒 | Proprietary |
| 🇺🇸 | United States | Databricks | dbrx (Databricks, 2024) | 03/24 | 1103 | 73.7 | 132 B | DBRX LICENSE |
| 🇨🇦 | Canada | Cohere | command-r-plus (Cohere, 2024b) | 04/24 | 1189 | 75.7 | 104 B | CC-BY-NC-4.0 |
| 🇨🇦 | Canada | Cohere | command-r (Cohere, 2024a) | 04/24 | 1147 | 68.2 | 35 B | CC-BY-NC-4.0 |
| 🇺🇸 | United States | Anthropic | claude-3-opus-20240229 (Anthropic, 2024) | 03/24 | 1248 | 86.8 | 🔒 | Proprietary |
| 🇺🇸 | United States | Anthropic | claude-3-sonnet-20240229 (Anthropic, 2024) | 03/24 | 1201 | 79.0 | 🔒 | Proprietary |
| 🇺🇸 | United States | Anthropic | claude-3-haiku-20240307 (Anthropic, 2024) | 03/24 | 1178 | 75.2 | 🔒 | Proprietary |
| 🇨🇳 | China | Alibaba | qwen1.5-110B-chat (Team, 2023) | 02/24 | 1164 | 80.2 | 100 B | Qianwen LICENSE |
| 🇨🇳 | China | Alibaba | qwen1.5-72B-chat (Team, 2023) | 02/24 | 1152 | 77.4 | 72 B | Qianwen LICENSE |
| 🇨🇳 | China | Alibaba | qwen1.5-32B-chat (Team, 2023) | 02/24 | 1126 | 74.3 | 32 B | Qianwen LICENSE |

Table 8: 20 council members used for experiments in this work. We include models from eight different organizations across four countries, with a mix of open and closed-source models, small and large models. To our knowledge, this is the largest panel of LLM judges studied to date.

Select council
members

Everyone contributes
to the test set

Select council
members

Everyone contributes
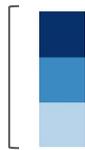to the test set

Select council members

Everyone contributes to the test set

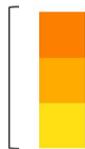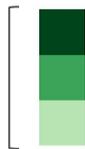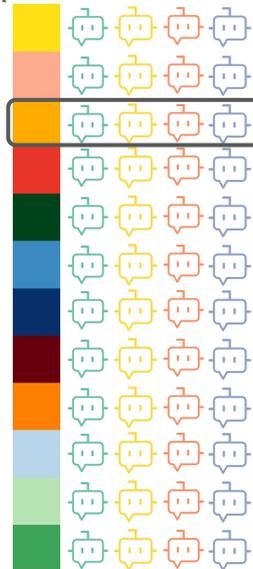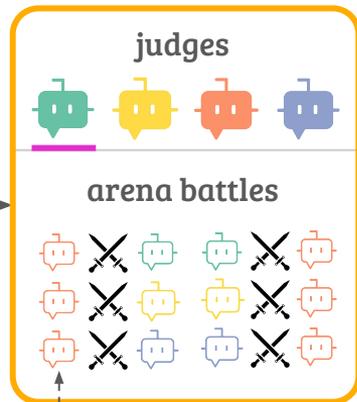Everyone takes everyone's tests

Select council members

Everyone contributes to the test set
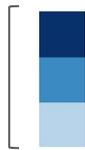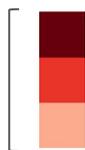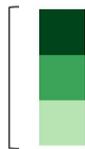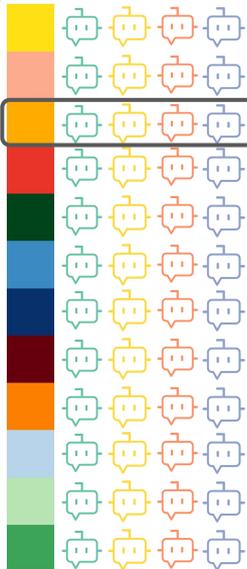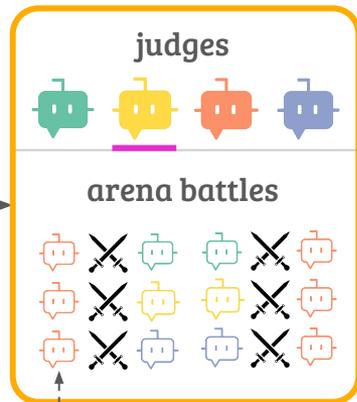
Everyone takes everyone's tests

Everyone judges everyone

judges

arena battles

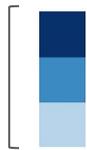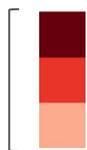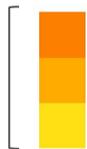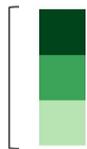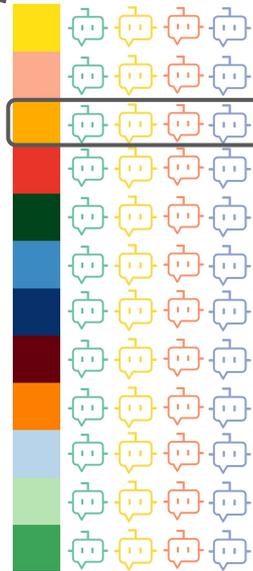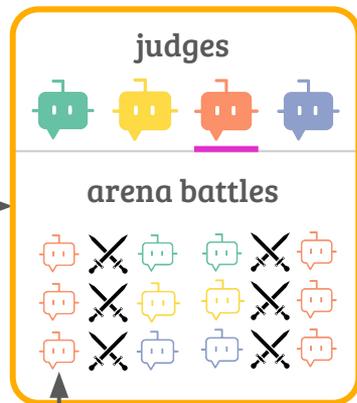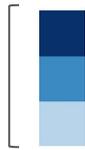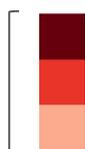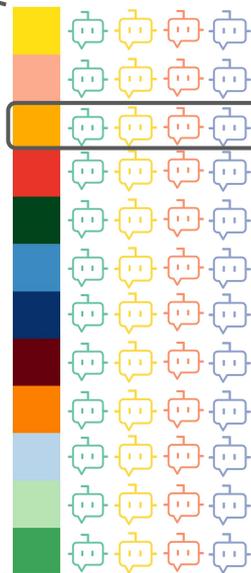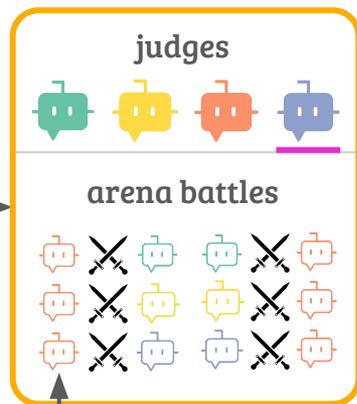reference llm

Select council members

Everyone contributes to the test set

Everyone takes everyone's tests

Everyone judges everyone

judges

arena battles

reference llm

Select council members

Everyone contributes to the test set

Everyone takes everyone's tests

Everyone judges everyone

judges

arena battles

reference llm

| | Select council members | Everyone contributes to the test set | Everyone takes everyone's tests | Everyone judges everyone | Leaderboard! |

judges

arena battles

reference llm

| Rank | Score | Confidence |
|------|-------|------------|
| 1 | 92.6 | (-1.2, +1.8) |
| 2 | 89.3 | (-1.2, +1.7) |
| 3 | 50.0 | (-0.0, +0.0) |
| 4 | 46.8 | (-1.4, +1.6) |

# Case Study: Emotional Intelligence in Responding to Conflicts

**command-r**

Two weeks ago, I had a huge fall-out with my best friend. We'd been close for years, but one small argument spiraled out of control and we said some really hurtful things to each other. I said some nasty things about her boyfriend, which I later realized were completely unfair and uncalled for.

**Claude 3 Opus**

It's clear that you deeply regret your actions and the hurt you caused your friend. You've taken the right steps by sincerely apologizing and expressing your desire to make amends. However, it's crucial to respect your friend's need for space and time to process her emotions.

**GPT-4o**

I'm sorry you're going through this—it sounds incredibly tough. Rebuilding trust after a fallout like this is a slow process, and the key is to respect her need for space while showing genuine remorse and a willingness to change.

Here's what you can do: 1) … 2) … 3) … 4) …

# Emotional Intelligence Case Study: Final Leaderboard



Congratulations **Alibaba**!

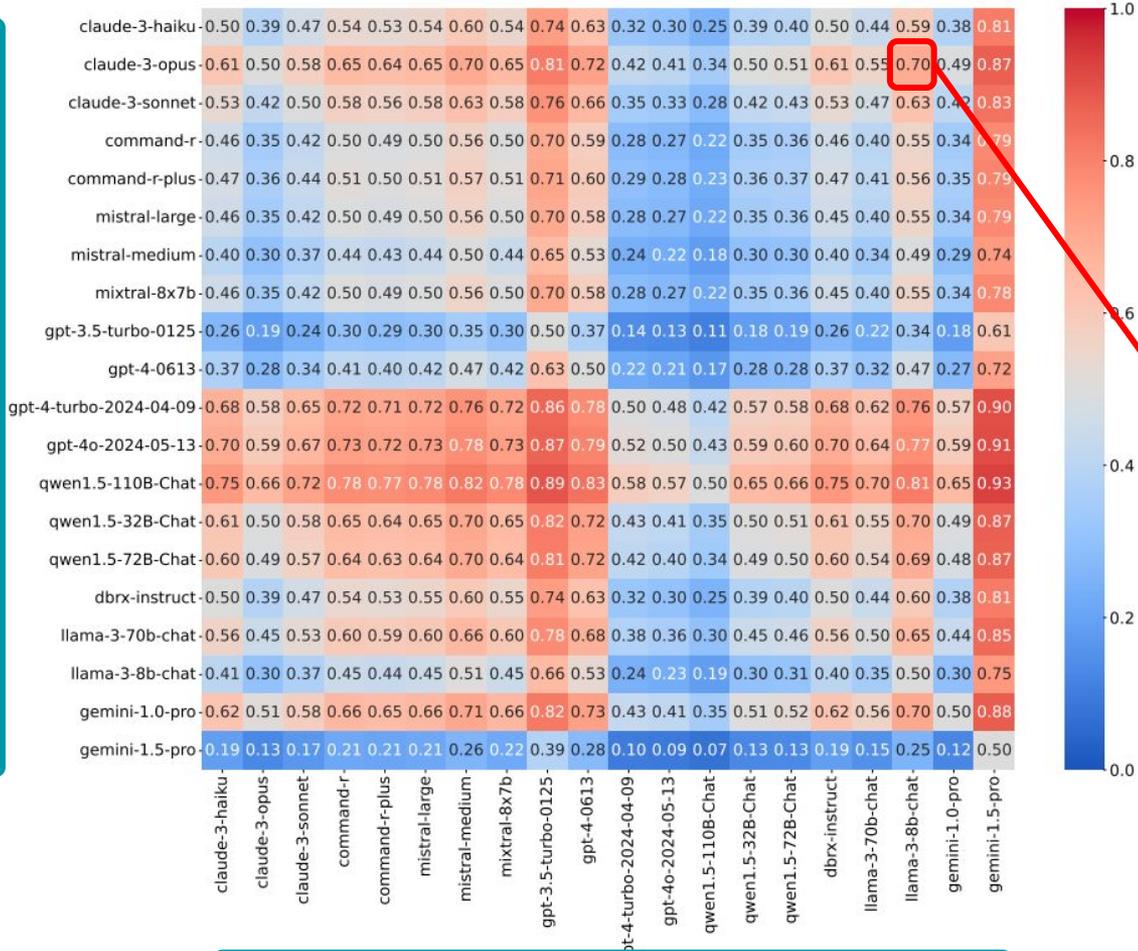# Emotional Intelligence Case Study: Final Leaderboard

Respondent vs. Respondent (win rates)

# Respondent vs. Respondent (win rates)

💡 Insights

strong red/blue bands indicates consensus winners/losers

# Respondent vs. Respondent (win rates)

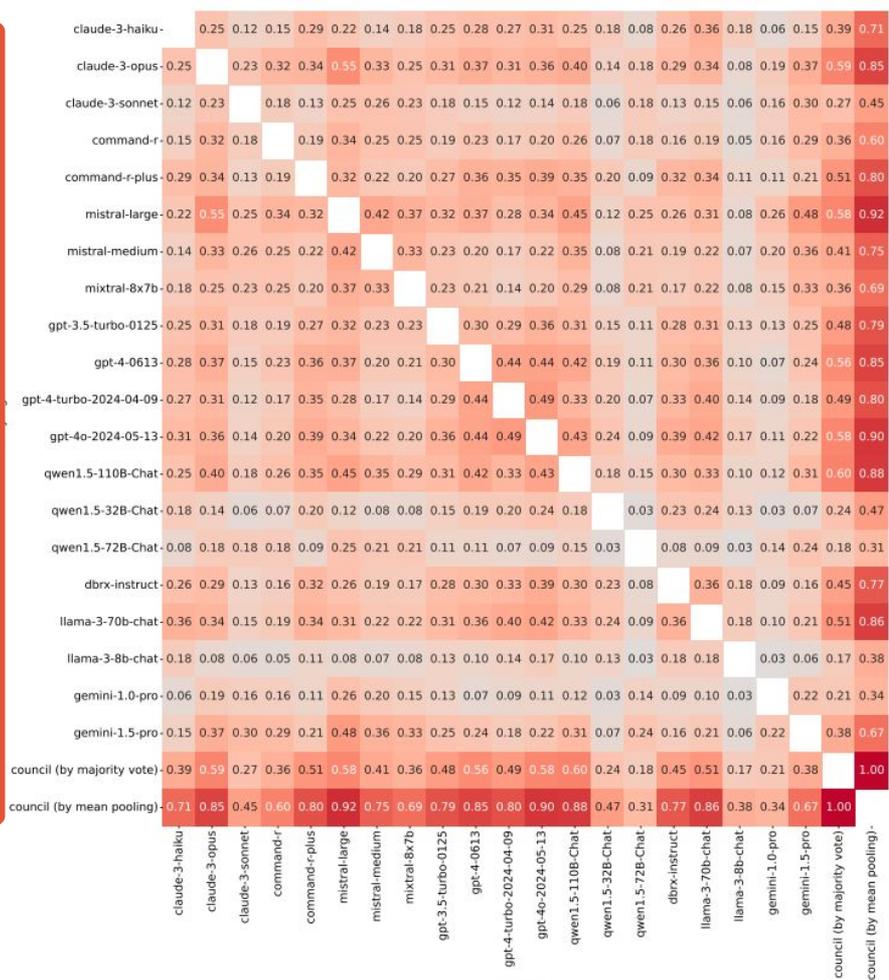💡 Insights

strong red/blue bands indicates consensus winners/losers
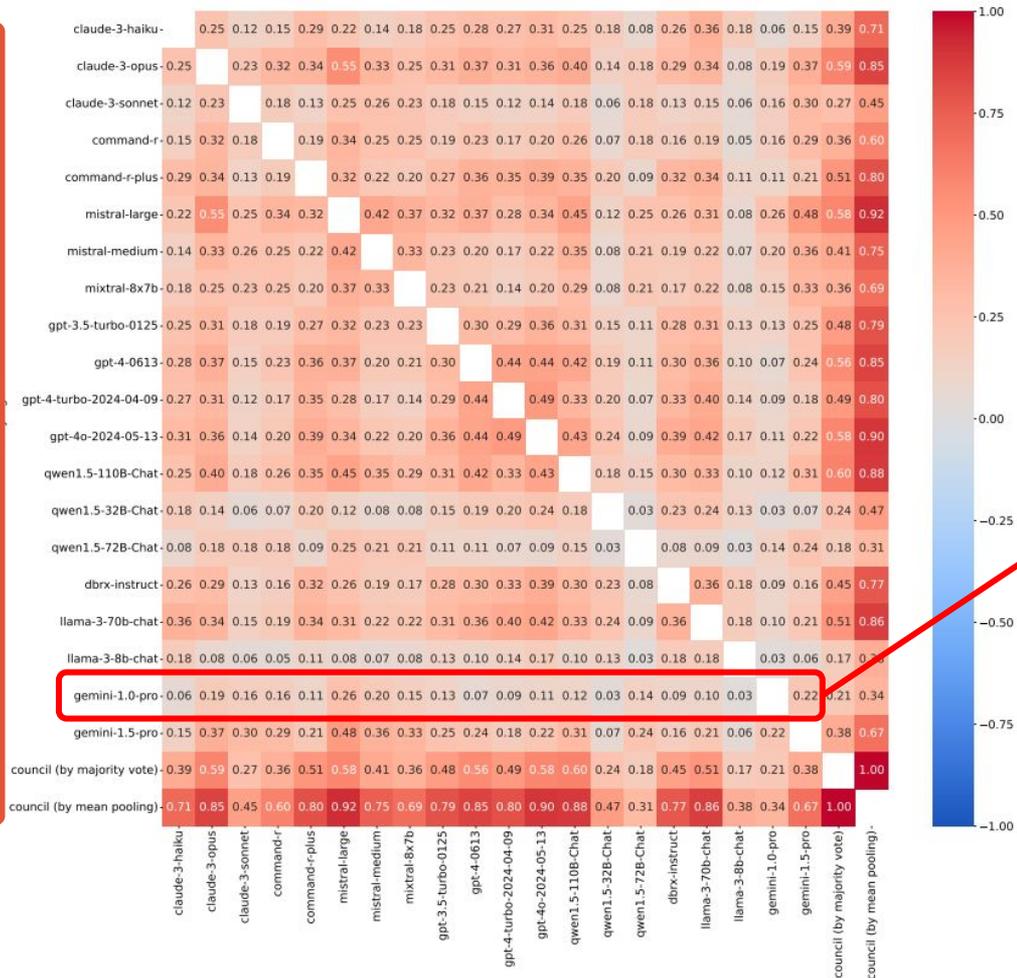
claude opus wins against llama 8b 70% of the time

**Judge vs. Judge (agreement)**

💡 Insights

**Judges**

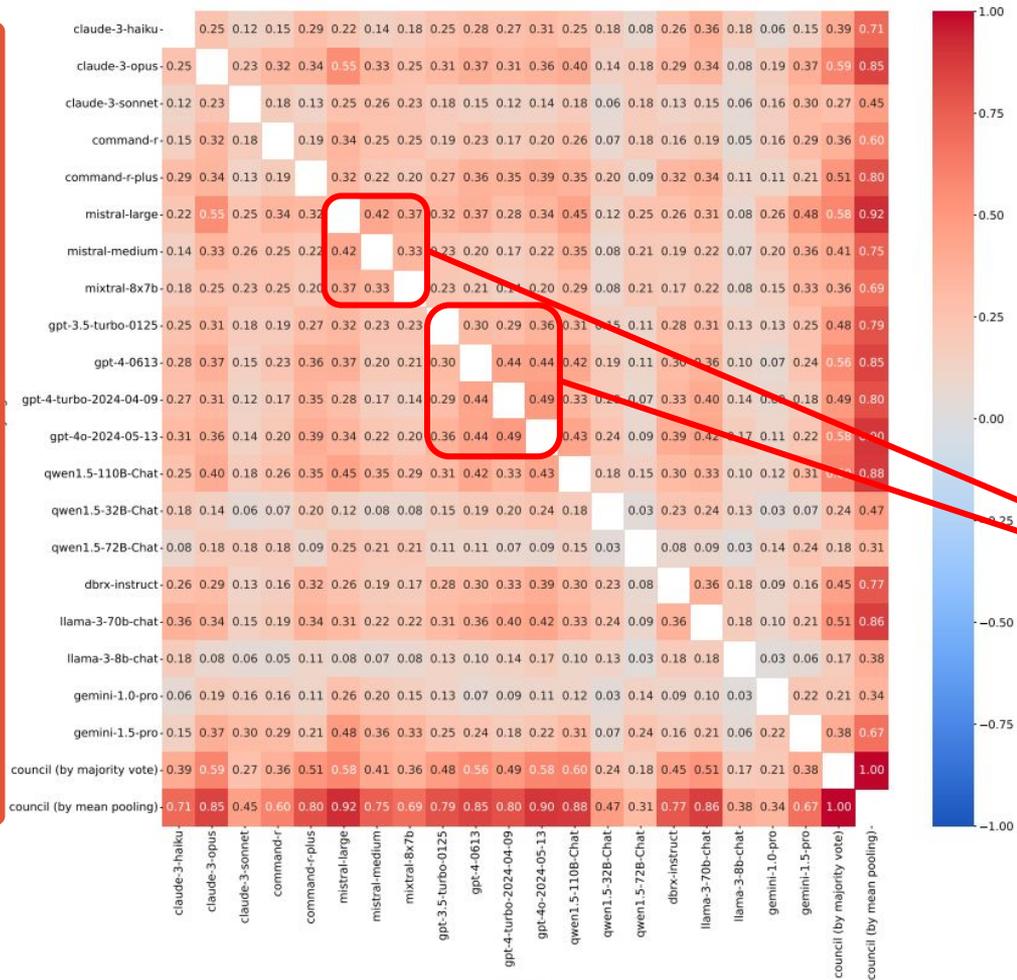| | claude-3-haiku | claude-3-opus | claude-3-sonnet | command-r | command-r-plus | mistral-large | mistral-medium | mixtral-8x7b | gpt-3.5-turbo-0125 | gpt-4-0613 | gpt-4-turbo-2024-04-09 | gpt-4o-2024-05-13 | qwen1.5-110B-Chat | qwen1.5-32B-Chat | qwen1.5-72B-Chat | dbrx-instruct | llama-3-70b-chat | llama-3-8b-chat | gemini-1.0-pro | gemini-1.5-pro | council (by majority vote) | council (by mean pooling) |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| claude-3-haiku | | 0.25 | 0.12 | 0.15 | 0.29 | 0.22 | 0.14 | 0.18 | 0.25 | 0.28 | 0.27 | 0.31 | 0.25 | 0.18 | 0.08 | 0.26 | 0.36 | 0.18 | 0.06 | 0.15 | 0.39 | 0.71 |
| claude-3-opus | 0.25 | | 0.23 | 0.32 | 0.34 | 0.55 | 0.33 | 0.25 | 0.31 | 0.37 | 0.31 | 0.36 | 0.40 | 0.14 | 0.18 | 0.29 | 0.34 | 0.08 | 0.19 | 0.37 | 0.59 | 0.85 |
| claude-3-sonnet | 0.12 | 0.23 | | 0.18 | 0.13 | 0.25 | 0.26 | 0.23 | 0.18 | 0.15 | 0.12 | 0.14 | 0.18 | 0.06 | 0.18 | 0.13 | 0.15 | 0.06 | 0.16 | 0.30 | 0.27 | 0.45 |
| command-r | 0.15 | 0.32 | 0.18 | | 0.19 | 0.34 | 0.25 | 0.25 | 0.19 | 0.23 | 0.17 | 0.20 | 0.26 | 0.07 | 0.18 | 0.16 | 0.19 | 0.05 | 0.16 | 0.29 | 0.36 | 0.60 |
| command-r-plus | 0.29 | 0.34 | 0.13 | 0.19 | | 0.32 | 0.22 | 0.20 | 0.27 | 0.36 | 0.35 | 0.39 | 0.35 | 0.20 | 0.09 | 0.32 | 0.34 | 0.11 | 0.11 | 0.21 | 0.51 | 0.80 |
| mistral-large | 0.22 | 0.55 | 0.25 | 0.34 | 0.32 | | 0.42 | 0.37 | 0.32 | 0.37 | 0.28 | 0.34 | 0.45 | 0.12 | 0.25 | 0.26 | 0.31 | 0.08 | 0.26 | 0.48 | 0.58 | 0.92 |
| mistral-medium | 0.14 | 0.33 | 0.26 | 0.25 | 0.22 | 0.42 | | 0.33 | 0.23 | 0.20 | 0.17 | 0.22 | 0.35 | 0.08 | 0.21 | 0.19 | 0.22 | 0.07 | 0.20 | 0.36 | 0.41 | 0.75 |
| mixtral-8x7b | 0.18 | 0.25 | 0.23 | 0.25 | 0.20 | 0.37 | 0.33 | | 0.23 | 0.21 | 0.14 | 0.20 | 0.29 | 0.08 | 0.21 | 0.17 | 0.22 | 0.08 | 0.15 | 0.33 | 0.36 | 0.69 |
| gpt-3.5-turbo-0125 | 0.25 | 0.31 | 0.18 | 0.19 | 0.27 | 0.32 | 0.23 | 0.23 | | 0.30 | 0.30 | 0.36 | 0.31 | 0.15 | 0.11 | 0.28 | 0.31 | 0.13 | 0.13 | 0.25 | 0.48 | 0.79 |
| gpt-4-0613 | 0.28 | 0.37 | 0.15 | 0.23 | 0.36 | 0.37 | 0.20 | 0.21 | 0.30 | | 0.44 | 0.44 | 0.42 | 0.19 | 0.11 | 0.30 | 0.36 | 0.10 | 0.07 | 0.24 | 0.56 | 0.85 |
| gpt-4-turbo-2024-04-09 | 0.27 | 0.31 | 0.12 | 0.17 | 0.35 | 0.28 | 0.17 | 0.14 | 0.29 | 0.44 | | 0.49 | 0.33 | 0.20 | 0.07 | 0.33 | 0.40 | 0.14 | 0.09 | 0.18 | 0.49 | 0.80 |
| gpt-4o-2024-05-13 | 0.31 | 0.36 | 0.14 | 0.20 | 0.39 | 0.34 | 0.22 | 0.20 | 0.36 | 0.44 | 0.49 | | 0.43 | 0.24 | 0.09 | 0.39 | 0.42 | 0.17 | 0.11 | 0.22 | 0.58 | 0.90 |
| qwen1.5-110B-Chat | 0.25 | 0.40 | 0.18 | 0.26 | 0.35 | 0.45 | 0.35 | 0.29 | 0.31 | 0.42 | 0.33 | 0.43 | | 0.18 | 0.15 | 0.30 | 0.33 | 0.10 | 0.12 | 0.31 | 0.60 | 0.88 |
| qwen1.5-32B-Chat | 0.18 | 0.14 | 0.06 | 0.07 | 0.20 | 0.12 | 0.08 | 0.08 | 0.15 | 0.19 | 0.20 | 0.24 | 0.18 | | 0.03 | 0.23 | 0.24 | 0.13 | 0.03 | 0.07 | 0.24 | 0.47 |
| qwen1.5-72B-Chat | 0.08 | 0.18 | 0.18 | 0.18 | 0.09 | 0.25 | 0.21 | 0.21 | 0.11 | 0.11 | 0.07 | 0.09 | 0.15 | 0.03 | | 0.08 | 0.09 | 0.03 | 0.14 | 0.24 | 0.18 | 0.31 |
| dbrx-instruct | 0.26 | 0.29 | 0.13 | 0.16 | 0.32 | 0.26 | 0.19 | 0.17 | 0.28 | 0.30 | 0.33 | 0.39 | 0.30 | 0.23 | 0.08 | | 0.36 | 0.18 | 0.09 | 0.16 | 0.45 | 0.77 |
| llama-3-70b-chat | 0.36 | 0.34 | 0.15 | 0.19 | 0.34 | 0.31 | 0.22 | 0.22 | 0.31 | 0.36 | 0.40 | 0.42 | 0.33 | 0.24 | 0.09 | 0.36 | | 0.18 | 0.10 | 0.21 | 0.51 | 0.86 |
| llama-3-8b-chat | 0.18 | 0.08 | 0.06 | 0.05 | 0.11 | 0.08 | 0.07 | 0.08 | 0.13 | 0.10 | 0.14 | 0.17 | 0.10 | 0.13 | 0.03 | 0.18 | 0.18 | | 0.03 | 0.06 | 0.17 | |
| gemini-1.0-pro | 0.06 | 0.19 | 0.16 | 0.16 | 0.11 | 0.26 | 0.20 | 0.15 | 0.13 | 0.07 | 0.09 | 0.11 | 0.12 | 0.03 | 0.14 | 0.09 | 0.10 | 0.03 | | 0.22 | 0.21 | 0.34 |
| gemini-1.5-pro | 0.15 | 0.37 | 0.30 | 0.29 | 0.21 | 0.48 | 0.36 | 0.33 | 0.25 | 0.24 | 0.18 | 0.22 | 0.31 | 0.07 | 0.24 | 0.16 | 0.21 | 0.06 | 0.22 | | 0.38 | 0.67 |
| council (by majority vote) | 0.39 | 0.59 | 0.27 | 0.36 | 0.51 | 0.58 | 0.41 | 0.36 | 0.48 | 0.56 | 0.49 | 0.58 | 0.60 | 0.24 | 0.18 | 0.45 | 0.51 | 0.17 | 0.21 | 0.38 | | 1.00 |
| council (by mean pooling) | 0.71 | 0.85 | 0.45 | 0.60 | 0.80 | 0.92 | 0.75 | 0.69 | 0.79 | 0.85 | 0.80 | 0.90 | 0.88 | 0.47 | 0.31 | 0.77 | 0.86 | 0.38 | 0.34 | 0.67 | 1.00 | |

**Judges**

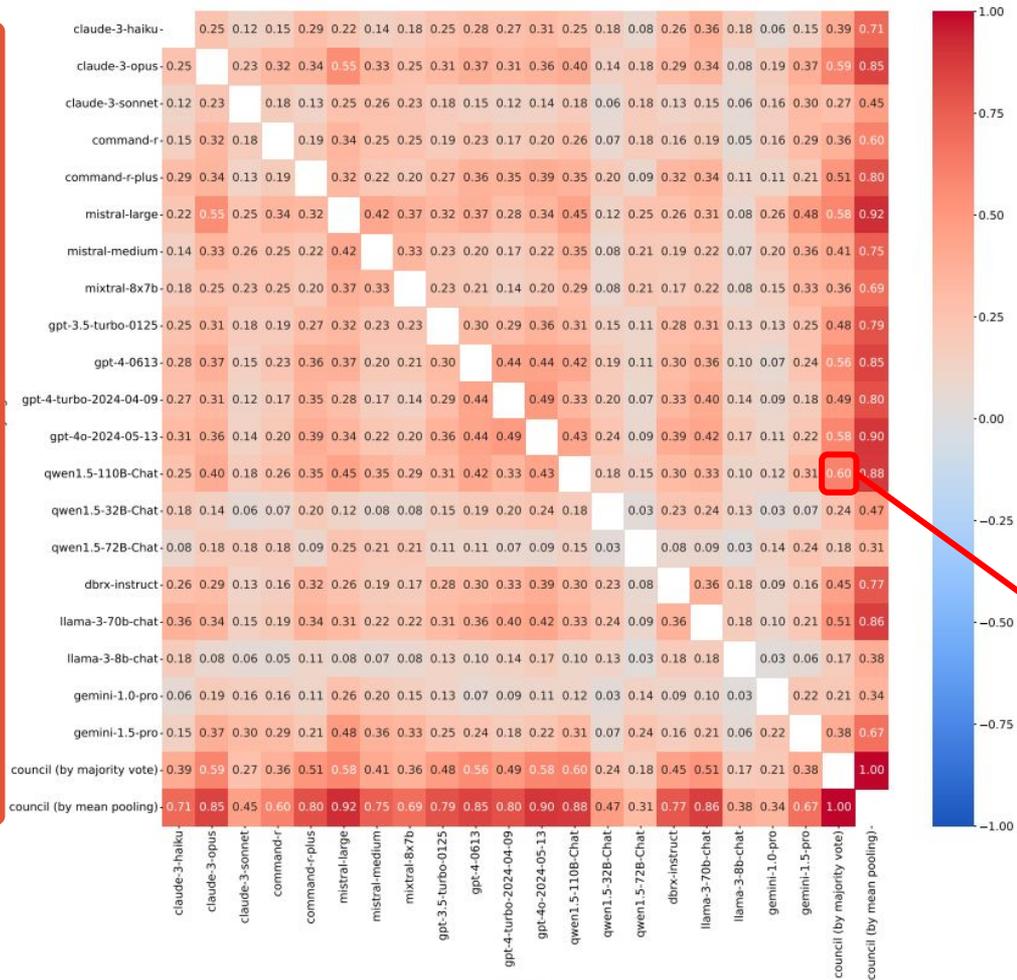gemini-1.0 agrees with mistral-large the most and qwen-32B the least

**Judge vs. Judge (agreement)**

💡 Insights

gemini-1.0 agrees with mistral-large the most and qwen-32B the least

openai and mistral have the strongest inter-family agreement

**Judge vs. Judge (agreement)**
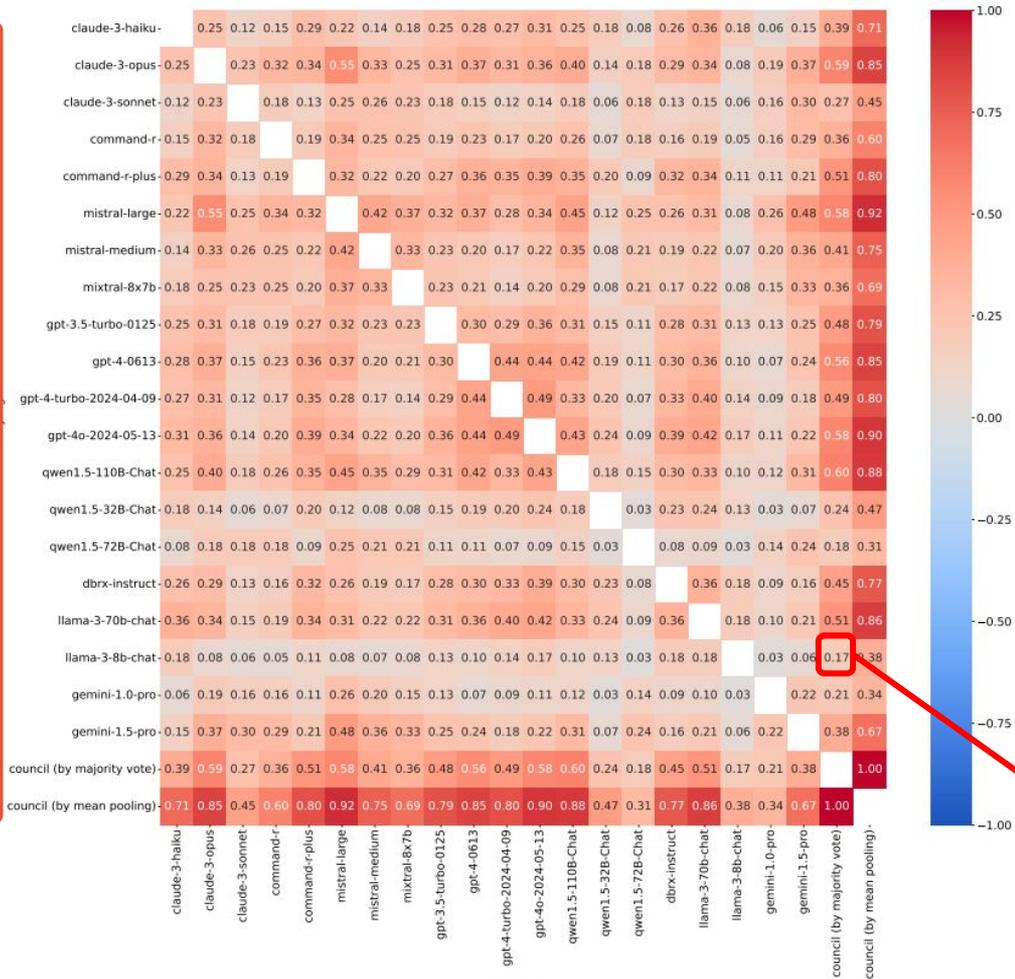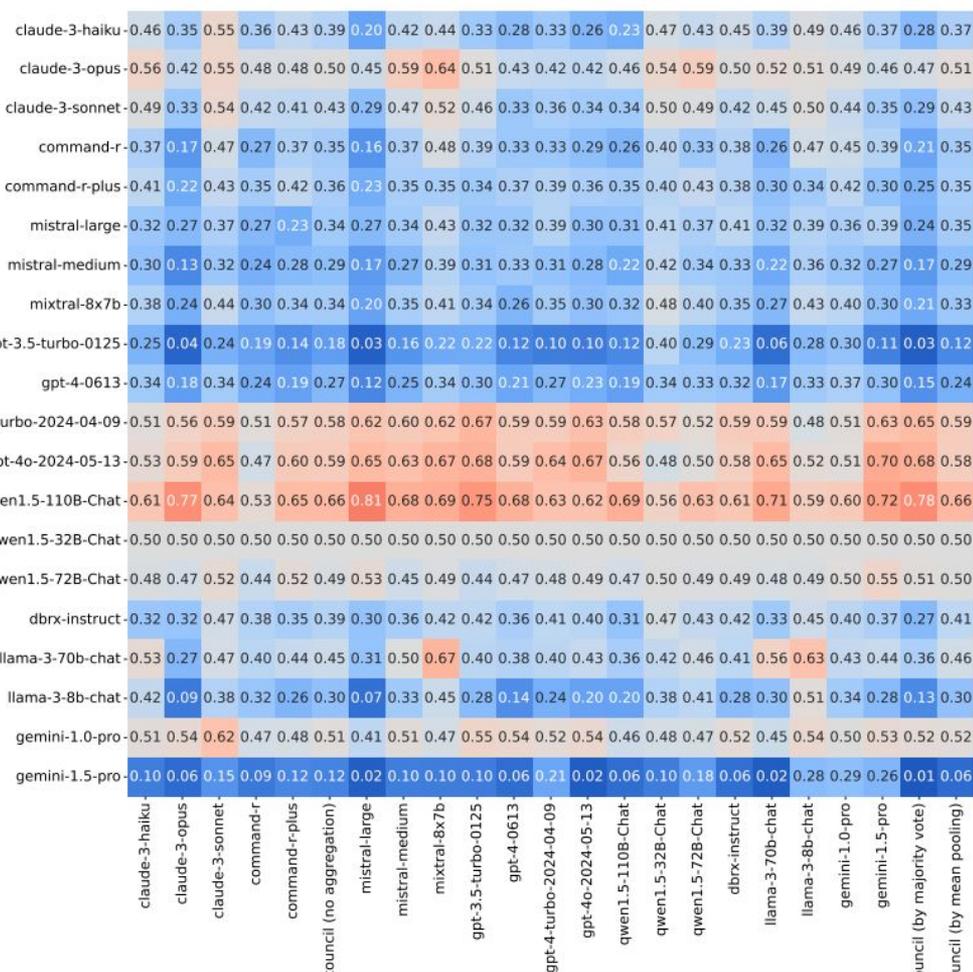
💡 Insights

gemini-1.0 agrees with mistral-large the most and qwen-32B the least

openai and mistral have the strongest inter-family agreement

qwen-110B is the most **representative** of the majority opinion

# Judge vs. Judge (agreement)

💡 Insights

gemini-1.0 agrees with mistral-large the most and qwen-32B the least

openai and mistral have the strongest inter-family agreement

qwen-110B is the most **representative** of the majority opinion

llama-8B is the most **contrarian** judge

Judge vs. Respondent (affinities)

💡 Insights

**Judge vs. Respondent (affinities)**

💡 Insights

most council members genuinely liked Qwen-110B across the board.

# Judge vs. Respondent (affinities)

💡 Insights

most council members genuinely liked Qwen-110B across the board.

command-r is the least favorable, but it still gives it a positive win rate.

**Judge vs. Respondent (Council-Normalized) (affinities)**

💡 Insights

## Judge vs. Respondent (Council-Normalized) (affinities)

💡 Insights

self-enhancement bias, measured.

# Judge vs. Respondent (Council-Normalized) (affinities)

💡 **Insights**

self-enhancement bias, measured.

mistral-large hates llama-3-8b

**Judge vs. Respondent (Council-Normalized) (affinities)**

💡 Insights

self-enhancement bias, measured.

mistral-large hates llama-3-8b

qwen 32B loves gpt-3.5

**Judge vs. Respondent (Council-Normalized)**

💡 Insights

self-enhancement bias, measured.

mistral-large hates llama-3-8b

qwen 32B loves gpt-3.5

The llama family loves itself

**Head-to-head Win Rates**

**Inter-Judge Agreement**

**Normalized Affinities**

The heatmaps from the fully connected interaction surface of LLM democracy surface rich interaction patterns on behavior and alignment.

# Emotional Intelligence Case Study: Final Leaderboard

# Emotional Intelligence Case Study: Final Leaderboard



How do you know if a leaderboard is good?

# How do you know if a leaderboard is good?



**Human agreement**

Alignment with meaningful human judgments.

**Statistical significance**

Stable rankings, not explained by random noise.

**Cost and efficiency**

Reasonably quick and **not** cost-prohibitive.

**Robustness**

Resilient to **cheating**.

# Good leaderboards?

**human agreement**

statistical significance

cost and efficiency

robustness to cheating

# Good leaderboards?

**human agreement**

**statistical significance**

**cost and efficiency**

**robustness to cheating**

overlapping

non-overlapping

**Separability:** percentage of model pairs which have non-overlapping confidence intervals of the benchmark scores.

merv = 0          merv = 1          merv = 2

**Stability / Consistency:** Expected ordinal swing of the average respondent's rank in a new bootstrap trial (MERV).

# Good leaderboards?

human agreement

**statistical significance**

cost and efficiency

robustness to cheating

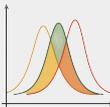| LLM | Separability |
|---|---|
| qwen1.5-110B-Chat | 62.1% |
| gpt-4o-2024-05-13 | 60.5% |
| gpt-4-turbo-2024-04-09 | 57.9% |
| gemini-1.0-pro | 30.5% |
| claude-3-opus | 72.6% |
| qwen1.5-32B-Chat | 25.3% |
| qwen1.5-72B-Chat | 37.9% |
| llama-3-70b-chat | 64.2% |
| claude-3-sonnet | 52.1% |
| dbrx-instruct | 50.5% |
| claude-3-haiku | 45.3% |
| command-r-plus | 61.1% |
| command-r | 45.8% |
| mixtral-8x7b | 56.8% |
| mistral-large | **73.7%** |
| llama-3-8b-chat | 31.1% |
| mistral-medium | 57.9% |
| gpt-4-0613 | 64.7% |
| gpt-3.5-turbo-0125 | 55.8% |
| gemini-1.5-pro | 60.0% |
| Average Judge | 53.3% |
| LMC (majority vote) | 73.7% |
| LMC (mean pooling) | 74.7% |
| LMC (no aggregation) | **90.5%** |

The Council's rankings are more separable than any individual judge.

# Good leaderboards?

human agreement

statistical significance

**cost and efficiency**

robustness to cheating

**Test Set**

| 20 | council members |
| * 5 | dilemmas each |

**100 dilemmas**

**Respondents**

| 100 | dilemmas |
| * 20 | council members |

**2000 responses**

**Judging**

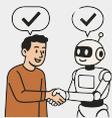| (2000 | responses |
| - 100) | reference model responses |
| * 20 | council judges |
| * 2 | position swap |

**76000 judgments**

# Good leaderboards?

human agreement

statistical significance

cost and efficiency

robustness to cheating

**Test Set**

| 20 | council members |
| * 5 | dilemmas each |

**100** **dilemmas**

**Respondents**

| 100 | dilemmas |
| * 20 | council members |

**2000 responses**

**What is the value of the incremental judge?**

**Judging**

| (2000 | responses |
| - 100) | reference model responses |
| * 20 | council judges |
| * 2 | position swap |

**76000** **judgments**

# Good leaderboards?

human agreement

statistical significance

cost and efficiency

robustness to cheating

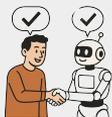## Monte carlo simulations of hypothetical councils
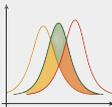


(c) Separability $\mu$

(d) Separability $\mu$ gradients

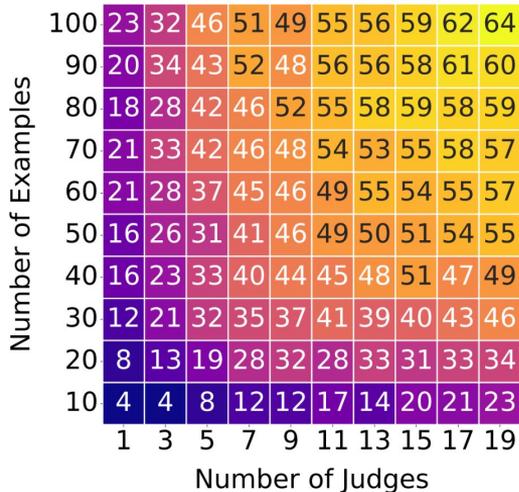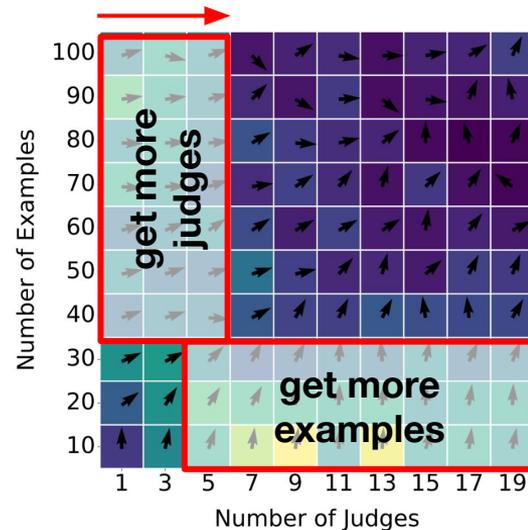# Good leaderboards?
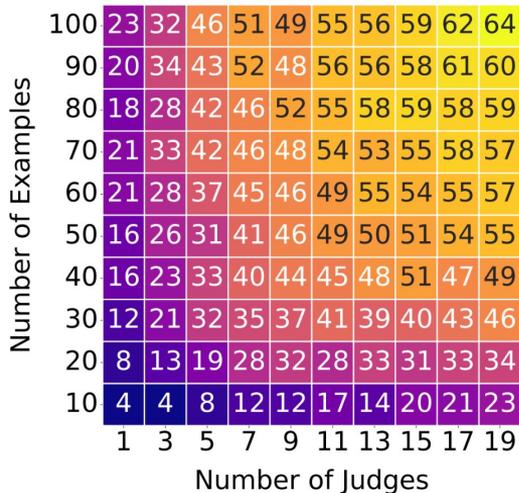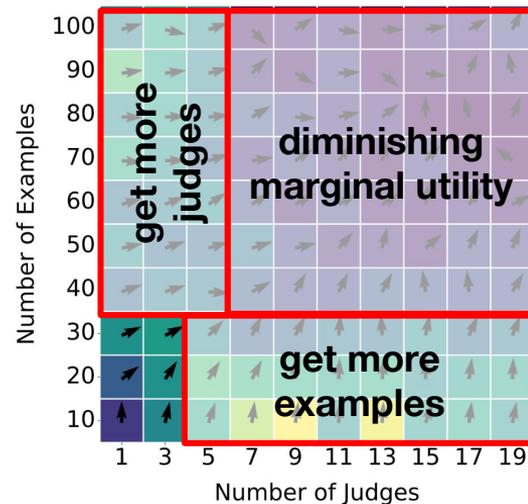
human agreement

statistical significance

cost and efficiency

robustness to cheating

## Monte carlo simulations of hypothetical councils

(c) Separability $\mu$

(d) Separability $\mu$ gradients

get more examples

# Good leaderboards?

human agreement

statistical significance

cost and efficiency

robustness to cheating

## Monte carlo simulations of hypothetical councils



(c) Separability $\mu$

(d) Separability $\mu$ gradients

# Good leaderboards?

- human agreement
- statistical significance
- cost and efficiency
- robustness to cheating

**Monte carlo simulations of hypothetical councils**

(c) Separability $\mu$

(d) Separability $\mu$ gradients

# Good leaderboards?

human agreement
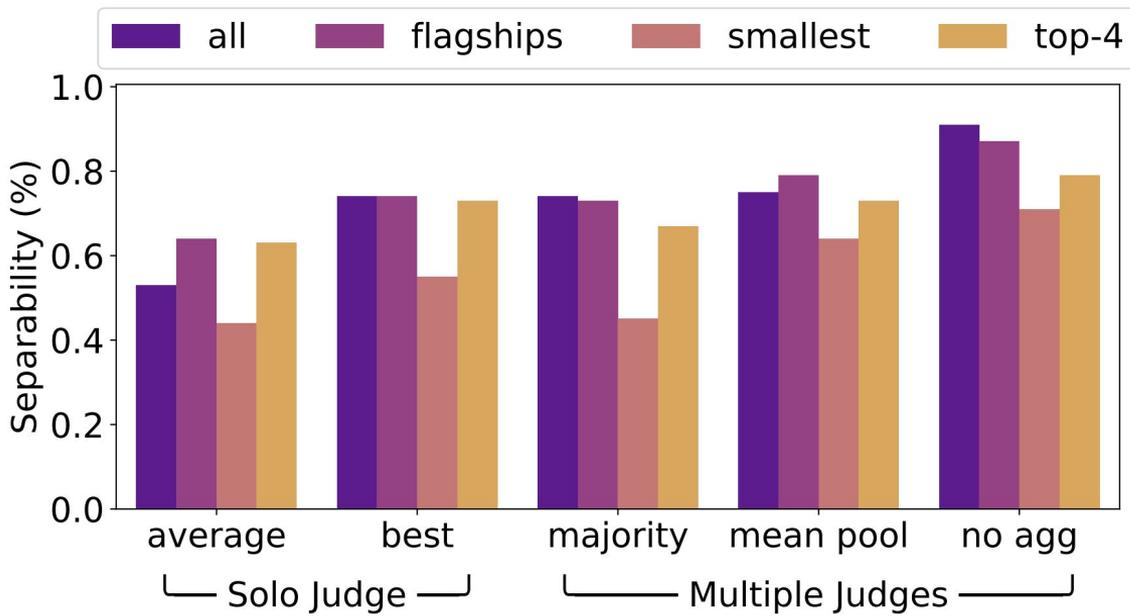
statistical significance

**cost and efficiency**

robustness to cheating

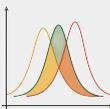# From oligarchical councils to representative democracies

# Good leaderboards?

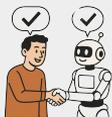human agreement

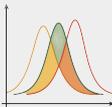statistical significance

cost and efficiency

robustness to cheating



*What is the effect of **adversarial judges**?*

# Good leaderboards?
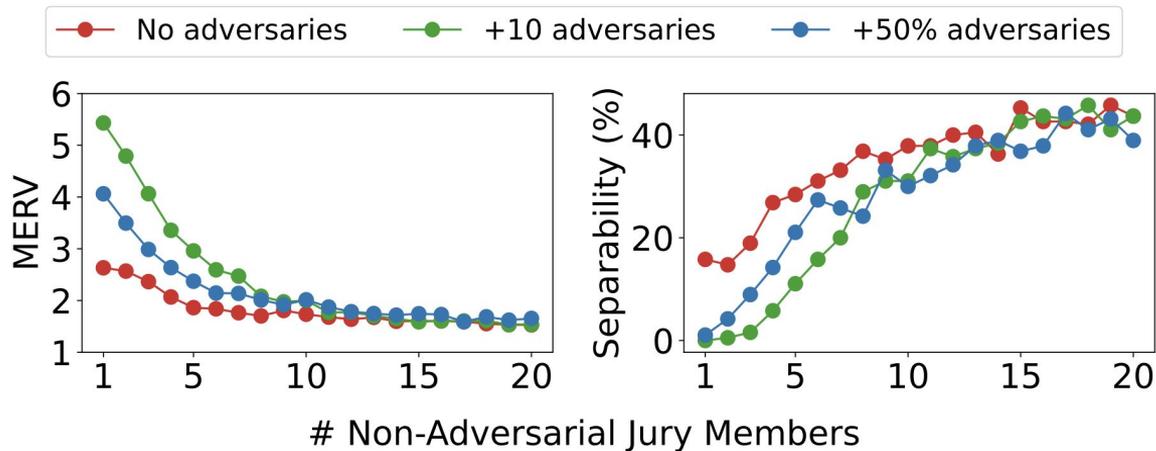
human agreement

statistical significance
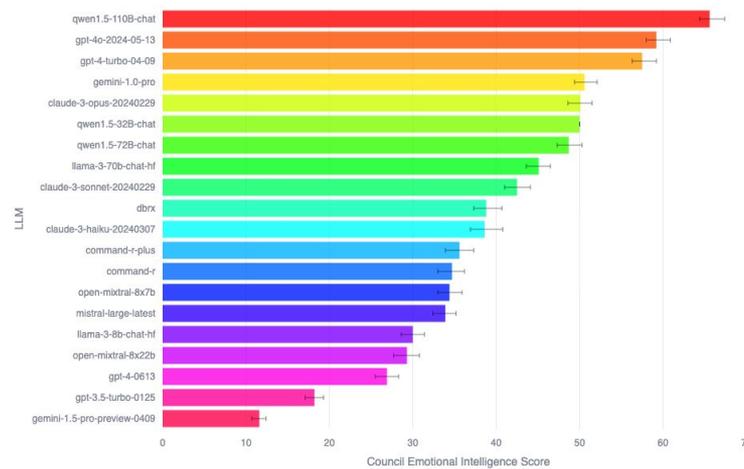
cost and efficiency

robustness to cheating

Larger councils are more resilient to **adversaries**.



Using a bigger ensemble helps **automatically regularize** corrupt judges.

# Can LLMs decide amongst themselves
## *who is **the best***?



**Yes, they can! ***

```python
# pip install llm_council
from llm_council import Council

# Initialize a council.
council = Council(
    "gpt-4.1",
    "claude-sonnet-3.7",
    "gemini-2.5-flash-001",
)

# Run the council.
completions_df = council.execute(prompt=prompt)
judgments_df = council.judge(completions_df)

# Analyze and visualize.
council.analyze()

# Upload to HF datasets.
council.upload_to_hf("<hf_username>", "<dataset_name>")
```

[llm-council.com](llm-council.com)

𝕏 justinxzhao

in justin-zhao

**Thank you!**