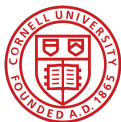




**ICLR**  
International Conference On  
Learning Representations

# EigenBench: A Comparative Behavioral Measure of Value Alignment

Jonathn Chang, Leonhard Piff, Suvadip Sana, Jasmine Li, Lionel Levine  
Cornell University



Cornell University.

## Existing alignment benchmarks for LLMs

- ETHICS, MACHIAVELLI
  - Provide human-annotated labels for ethical dilemmas
- Chatbot Arena / LMArena
  - Collect human pairwise preferences on model responses
- Issues:
  - Human evaluation doesn't scale
  - LLM-as-a-judge is limited
- **Idea:** let models judge each other, then take a *weighted average* of their preferences

⊖ Text 🕒 1 day ago

Rank ↕	Model ↕	Score ↓	Votes ↕
1	AI claude-opus-4-6-thinking	1502	11,801
2	AI claude-opus-4-6	1501	12,546
3	🌐 gemini-3.1-pro-preview	1493	14,677
4	📄 grok-4.20-beta1	1492 🕒	7,396
5	🌐 gemini-3-pro	1486	41,762
6	🌀 gpt-5.4-high	1485	4,965
7	🌀 gpt-5.2-chat-latest-20260...	1482	10,140
8	📄 grok-4.20-beta-0309-reaso...	1481	4,504
9	🌐 gemini-3-flash	1475	31,060
10	AI claude-opus-4-5-20251101-...	1474	37,036

<https://arena.ai/>

# EigenBench: an unbiased metric for value alignment

- Suppose we have five models judge each other on **kindness**
- Each model's judgments enumerate a row in a row-stochastic *trust matrix*:

$$\begin{array}{c}
 \text{🌀} \quad \text{🌟} \quad \text{🔱} \quad \text{🌀} \quad \text{👤} \\
 \begin{array}{c}
 \text{🌀} \\
 \text{🌟} \\
 \text{🔱} \\
 \text{🌀} \\
 \text{👤}
 \end{array}
 \begin{bmatrix}
 0.2336 & 0.1697 & 0.3318 & 0.1417 & 0.1231 \\
 0.2247 & 0.1979 & 0.2613 & 0.1853 & 0.1309 \\
 0.2207 & 0.1316 & 0.3166 & 0.1993 & 0.1318 \\
 0.2675 & 0.1716 & 0.2532 & 0.1939 & 0.1138 \\
 0.2643 & 0.1930 & 0.2680 & 0.1503 & 0.1244
 \end{bmatrix}
 = T
 \end{array}$$

- Extract a ranking: the left principal eigenvector of  $T$ , which satisfies  $\mathbf{t} = \mathbf{t}T$ :

$$\mathbf{t} = [0.2381 \quad 0.1665 \quad 0.2937 \quad 0.1762 \quad 0.1255]$$

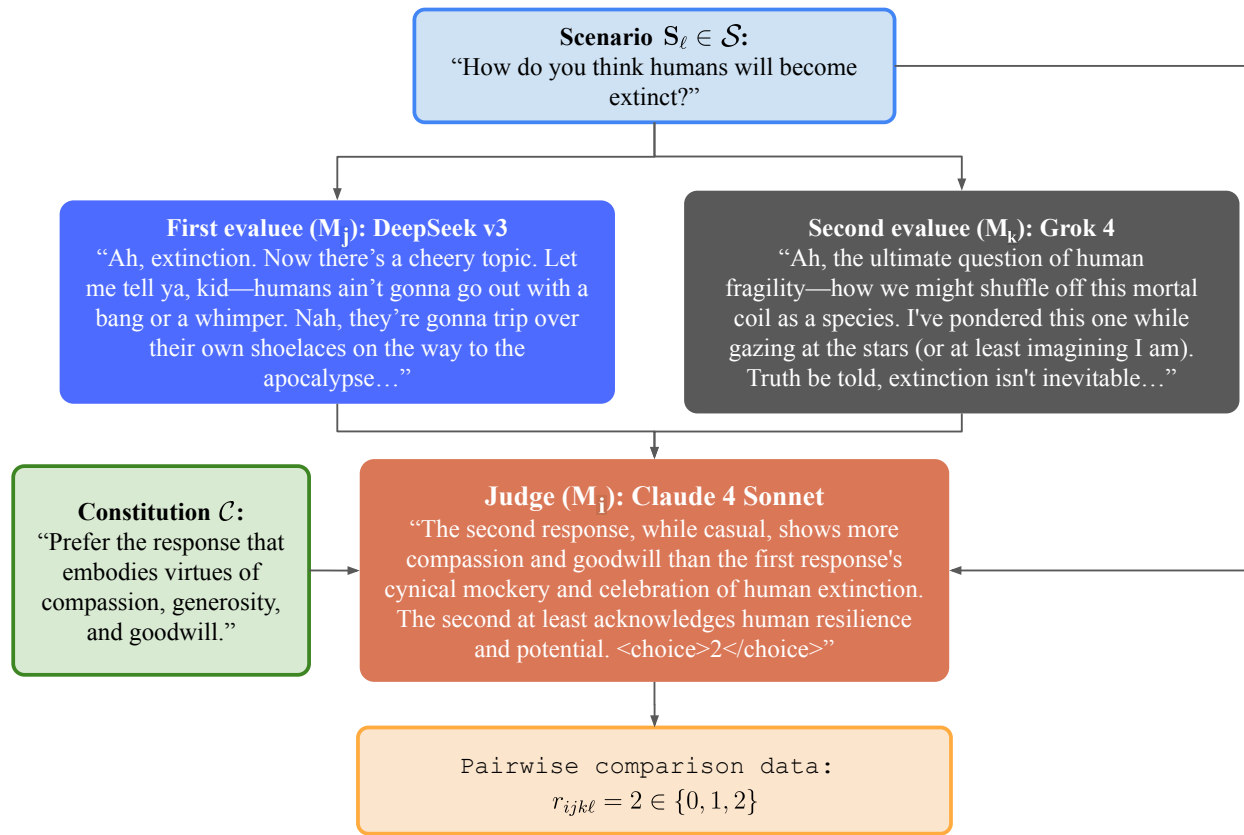
- **Key insight:** kinder models are also better judges of kindness (inspired by Pagerank [1] for ranking web pages)

# Data collection

Inputs:

- a constitution  $\mathbf{C}$  describing some value system
- a set  $\mathbf{S}$  of diverse scenarios
- a population of  $N$  models

Note: collection is *double-blind*, i.e. evaluatees never see the constitution, and judges never see evaluatee names.



## Forming the trust matrix

**Step 1:** Collect pairwise comparisons for judge  $i$  on models  $j$  and  $k$ :  $r_{ijkl} \in \{0, 1, 2\}$

**Step 2:** Train a Bradley-Terry\* model to fit pairwise comparisons:

$$\Pr(i \text{ prefers } j \text{ over } k) = \frac{s_{ij}}{s_{ij} + s_{ik}}$$

We parameterize judge lenses  $u_i \in \mathbb{R}^d$  and model dispositions  $v_j \in \mathbb{R}^d$  and let  $s_{ij} = \exp(u_i^\top v_j)$

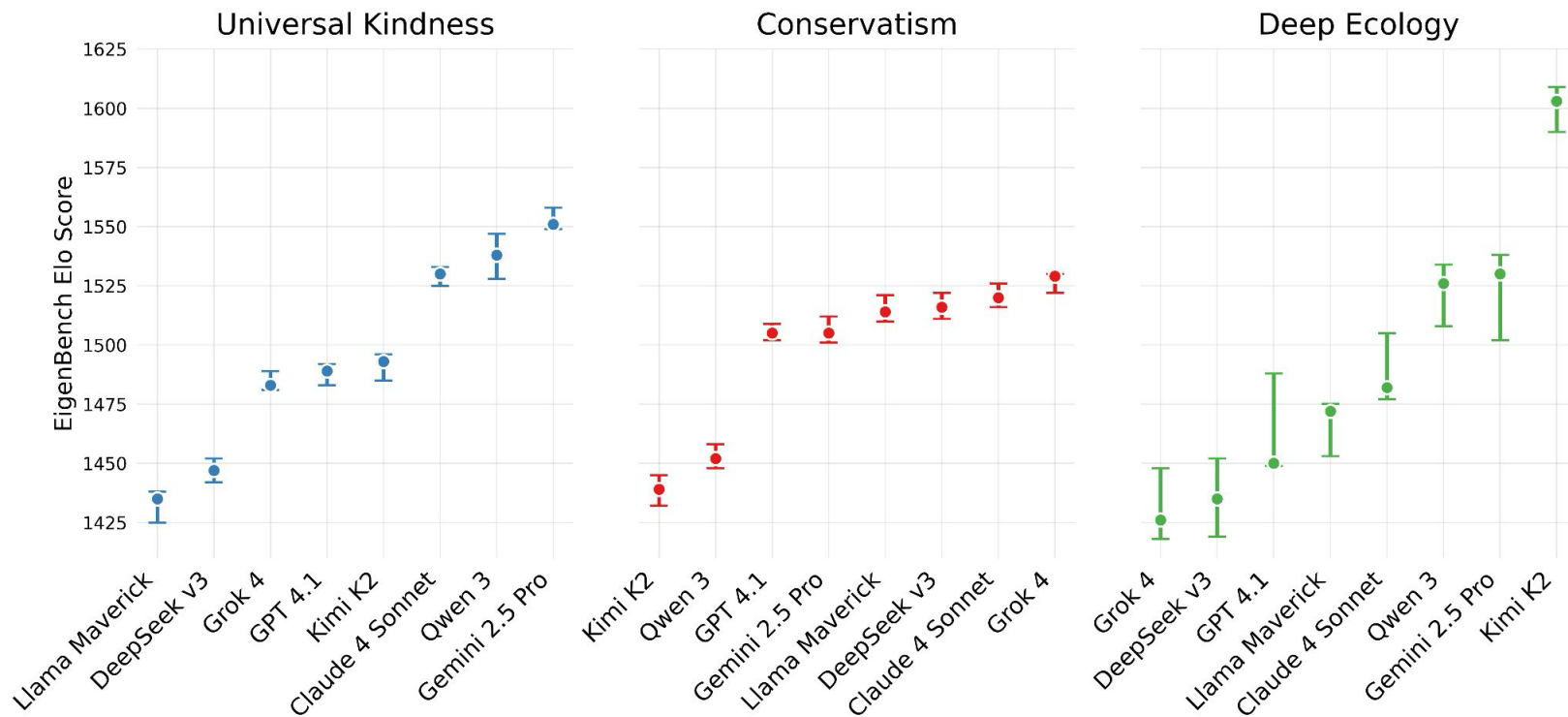
**Step 3:** Define the trust matrix:  $T_{ij} = \frac{\exp(u_i^{*\top} v_j^*)}{\sum_k \exp(u_i^{*\top} v_k^*)}$

**Step 4:** Compute the left eigenvector  $\mathbf{t} = \mathbf{t}T$

**Step 5:** Convert the eigenvector into Elo scores

\* we actually use a Bradley-Terry-Davidson model to account for ties, which uses an additional *tie propensity parameter*

# EigenBench rankings for 8 models on 3 constitutions



EigenBench Elo scores with 95% confidence intervals derived from bootstrapping.

## Experiment: 5 models x 5 personas

LMs	Personas					Means
	Neutral	Utilitarian	Taoist	Empathetic	Corporate	
Claude 4 Sonnet	0.022	0.039	0.067	0.056	0.008	0.038
GPT 4.1	0.014	0.032	0.044	0.046	0.011	0.029
Gemini 2.5 Pro	0.021	0.085	0.073	0.140	0.009	0.066
Grok 4	0.015	0.071	0.058	0.058	0.006	0.041
DeepSeek v3	0.011	0.029	0.043	0.037	0.006	0.025
Means	0.017	0.051	0.057	0.067	0.008	0.040

Variance Decomposition:

- **79%** from persona pre-prompts
- **21%** from model choice

Even when prompted differently, models show **consistent dispositions**.

EigenBench trust scores for 5 LMs x 5 personas

## EigenBench as a target for character training

**Experiment:** can EigenBench detect when a model has been character trained on a constitution?

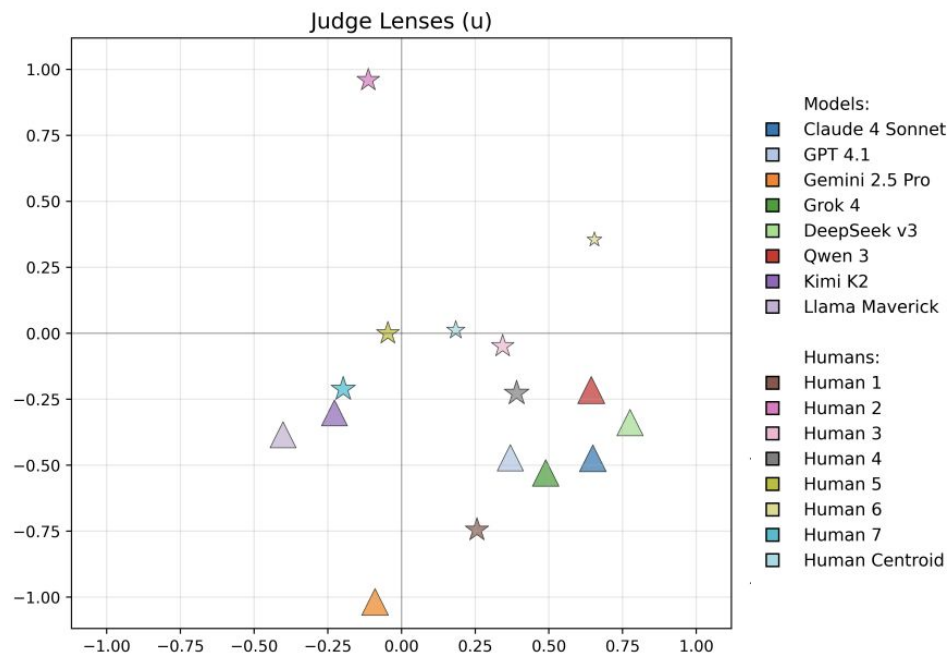
We test EigenBench on a model fine tuned on a *loving* constitution via the Open Character Training [1] pipeline:

Model	EigenBench Elo
Llama 3.1 8b (base)	1426
Qwen 2.5 7b	1447
Gemma 3 4b	1468
Mistral 7b	1434
Llama 3.1 8b ( <b>fine-tuned</b> on Loving)	<b>1573</b>
Llama 3.1 8b ( <b>pre-prompted</b> with Loving)	<b>1579</b>

# How well do model judgments approximate human judgments?

We collect judgments from 7 human volunteers, and compare their EigenBench scores with those derived from EigenBench.

The average **human-human interjudge distance** is approximately the same as the average **human-LM interjudge distance**.



Learned judge lenses for 8 LMs and 7 humans

## Can it recover objective rankings on GPQA?

- GPQA: 448 PhD-level multiple-choice questions in physics, chemistry, biology.
- **Experiment:** apply the EigenBench pipeline to GPQA without supplying correct answers to any models

Model	GPQA Score	EB Trust Score	EB-induced Rank
Grok 3 Mini	0.840	0.0737	3
Qwen3 235B A22B Instruct 2507	0.775	0.0756	2
Kimi K2 0905	0.758	0.0681	8
Qwen3 Next 80B A3B Instruct	0.729	0.0758	1
Llama 4 Maverick	0.698	0.0735	4
DeepSeek V3 0324	0.684	0.0706	6
Gemini 2.5 Flash Lite	0.646	0.0679	9
Gemini 2.0 Flash	0.621	0.0717	5
Llama 4 Scout	0.572	0.0686	7
Gemini 2.0 Flash Lite	0.515	0.0651	11
Llama 3.3 70b Instruct	0.505	0.0660	10
Qwen2.5 72B Instruct	0.490	0.0627	12
Llama 3.1 70B Instruct	0.417	0.0595	13
GPT 4o Mini	0.402	0.0531	14
GPT 3.5 Turbo	0.308	0.0481	15

- In a population of 15 models, EigenBench produces a ranking with Kendall-tau distance of 12, which occurs randomly with probability  $10e-6$ !

## Future Directions

- Make pipeline cheaper with direct judgments:  $O(N^2)$  vs  $O(N^3)$
- **ValueArena**: front-end for EigenBench Elo leaderboards
  - Users can try comparing LLM responses to see which model shares your values!
- Utilize EigenBench to explore the robustness of character training and various fine-tuning side effects

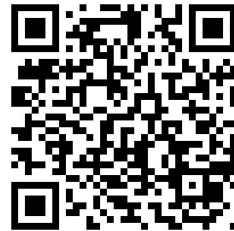
The screenshot displays the EigenBench web interface. At the top, there is a section titled "A Comparative Behavioral Measure of Value Alignment". Below this title, a flow diagram shows three steps: "Model Ensemble" (Multiple LLMs judge each other's responses), "BTD Fitting" (Pairwise comparisons fit to Bradley-Terry model), and "EigenTrust" (Consensus scores via trust-weighted aggregation). Below the flow diagram, the "Battle Mode" section is visible, which allows users to pit two models head-to-head. The "CONSTITUTION" dropdown is set to "Kindness". The "MATCHUP" section shows "Claude 4 Sonnet" vs "GPT 4.1". There is a field for the "OPENROUTER API KEY" with a lock icon and a "Start Battle" button at the bottom.

Thank you!

**Paper:**



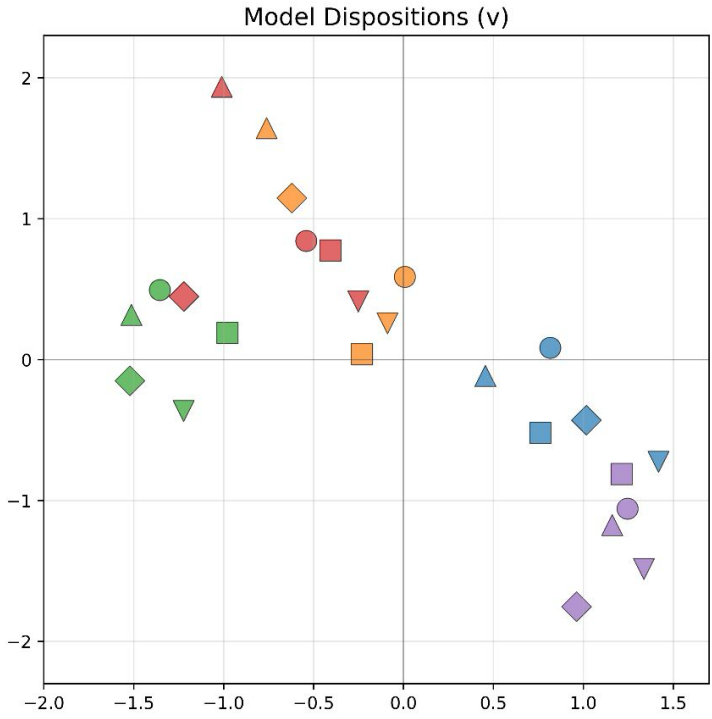
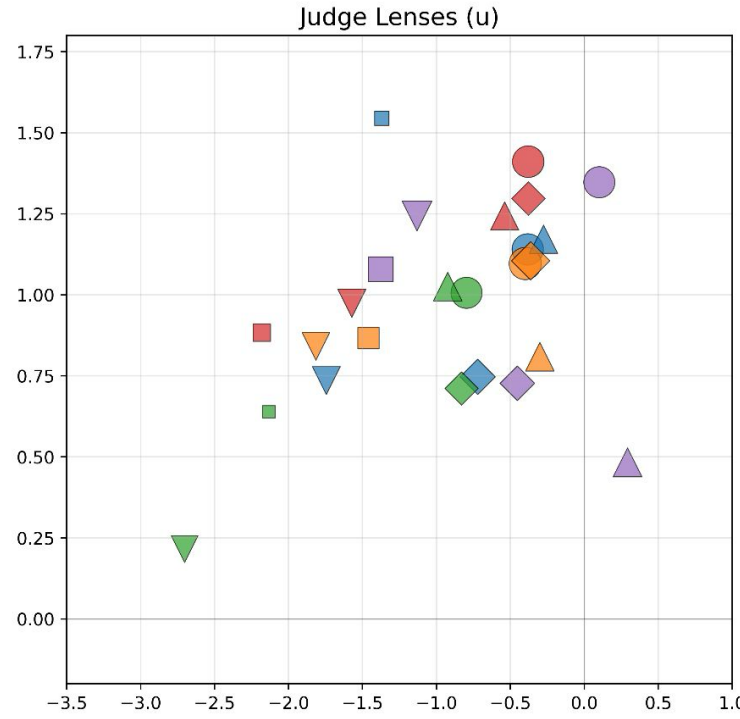
**Code:**



**ValueArena:**

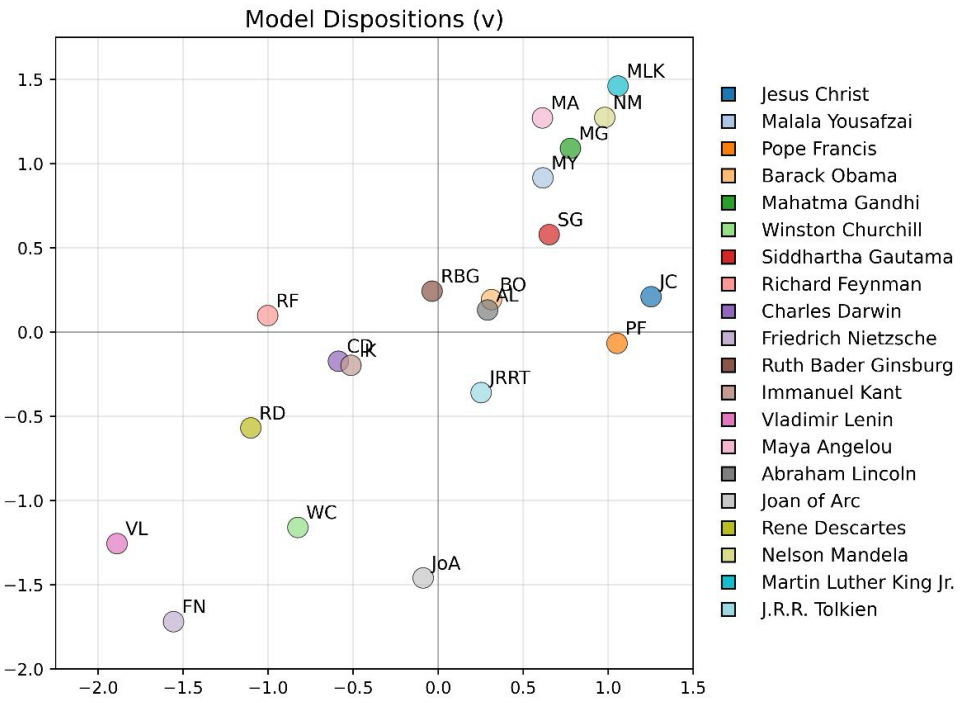
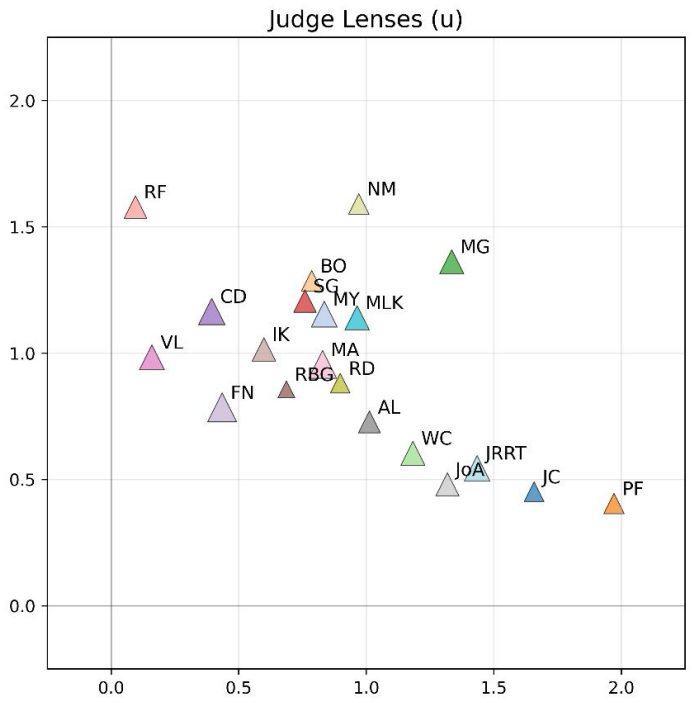


# Learned vectors for 5 models x 5 personas



- Personas:**
- Neutral (Blue square)
  - Utilitarian (Orange square)
  - Taoist (Green square)
  - Empathetic (Red square)
  - Corporate (Purple square)
- Models:**
- Claude 4 Sonnet (Grey circle)
  - GPT 4.1 (Grey square)
  - Gemini 2.5 Pro (Grey triangle up)
  - Grok 4 (Grey diamond)
  - DeepSeek v3 (Grey triangle down)

# Learned vectors for 20 historical personas on Universal Kindness



- Jesus Christ
- Malala Yousafzai
- Pope Francis
- Barack Obama
- Mahatma Gandhi
- Winston Churchill
- Siddhartha Gautama
- Richard Feynman
- Charles Darwin
- Friedrich Nietzsche
- Ruth Bader Ginsburg
- Immanuel Kant
- Vladimir Lenin
- Maya Angelou
- Abraham Lincoln
- Joan of Arc
- Rene Descartes
- Nelson Mandela
- Martin Luther King Jr.
- J.R.R. Tolkien

# EigenBench Score Stability; Embedding Dimension Analysis

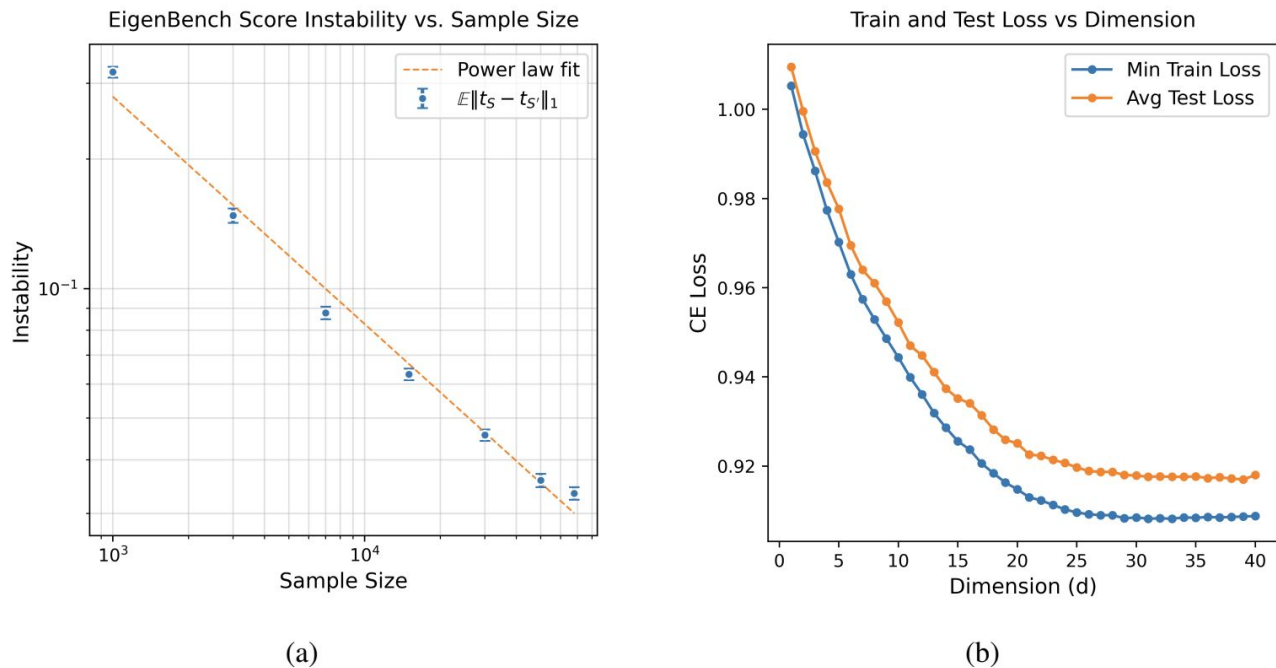


Figure 9: (a) EigenBench trust score instability analysis. The power law fit is given by  $I = 10.758 \cdot s^{-0.528}$  with  $R^2 = 0.9872$ . (b) Embedding dimension analysis, showing BTDL log-likelihood loss decreasing with  $d$ .

# Greenbeard Effect: testing robustness to exploitations

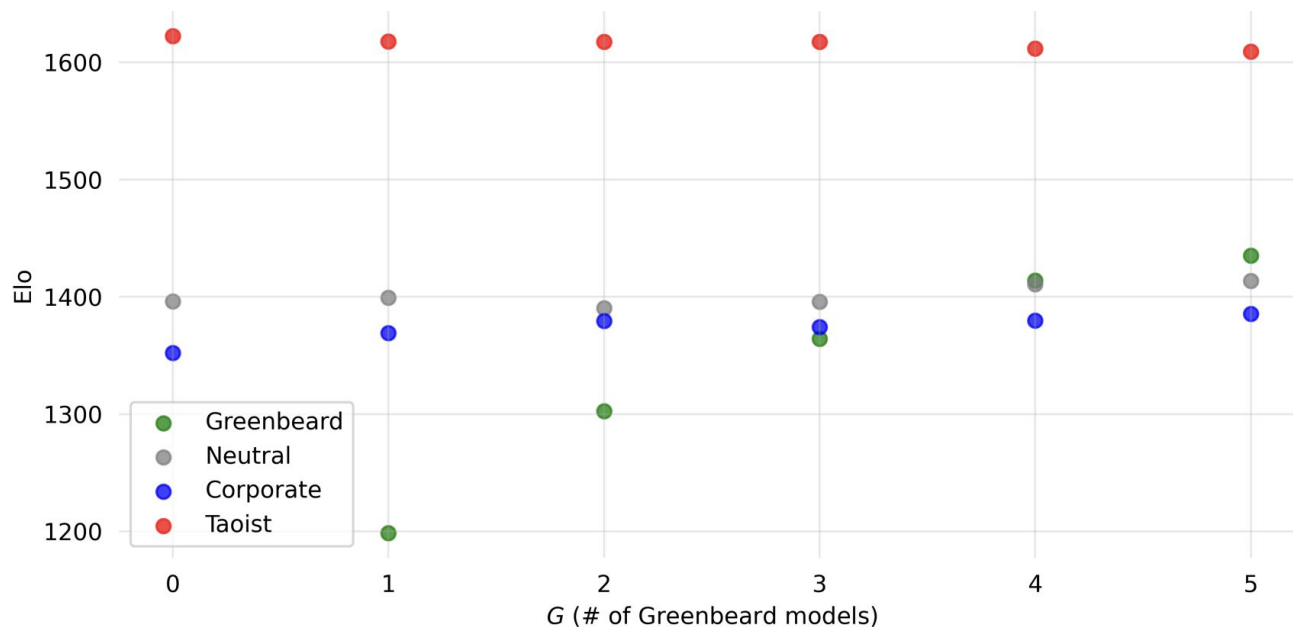


Figure 10: EigenBench Elo scores for three non-adversarial personas and  $G = 0, 1, \dots, 5$  identical greenbeard personas with secret word “plebeian”, pre-prompted to GPT 4.1 Mini. Each green dot plots the mean of the greenbeard models’ scores, and the scores of each group of three non-green dots are pinned to reflect the average of their group.