





40%

- 35% 🕤

- 25% -

<u>≞</u>: %20

<sup>1,2,3</sup>Michael W. Mahoney @<sup>1</sup>UC Berkeley, <sup>2</sup>ICSI, <sup>3</sup>LBL https://leaderbot.org  $\triangleleft$  CHATBOTS

<sup>1,2</sup>Siavash Ameli
<sup>1</sup>Siyuan Zhuang
<sup>1</sup>Ion Stoica
hael W. Mahoney
keley, <sup>2</sup>ICSI, <sup>3</sup>LBL
FRAMEWORK for
RANKING
LLM-BASED



importance of balanced scaling.LM Characteristics and Scores.



predictionsarginal Probabilities Across Models

Kernel PCA projection of LLMs' dissimilarities into 3D space. Distances reflecting pairwise dissimilarities. Circle size and color indicate scores, highlighting clusters of similar performance.



MDS projection of dissimilarities into 2D space. LLMs are spatially arranged based on pairwise differences, with size and color showing scores. The plot aligns relative scores with dissimilarities, uncovering meaningful patterns.

2024-05-13

gemini-1.5-pro-exp-0801

chatgpt-4o-latest