

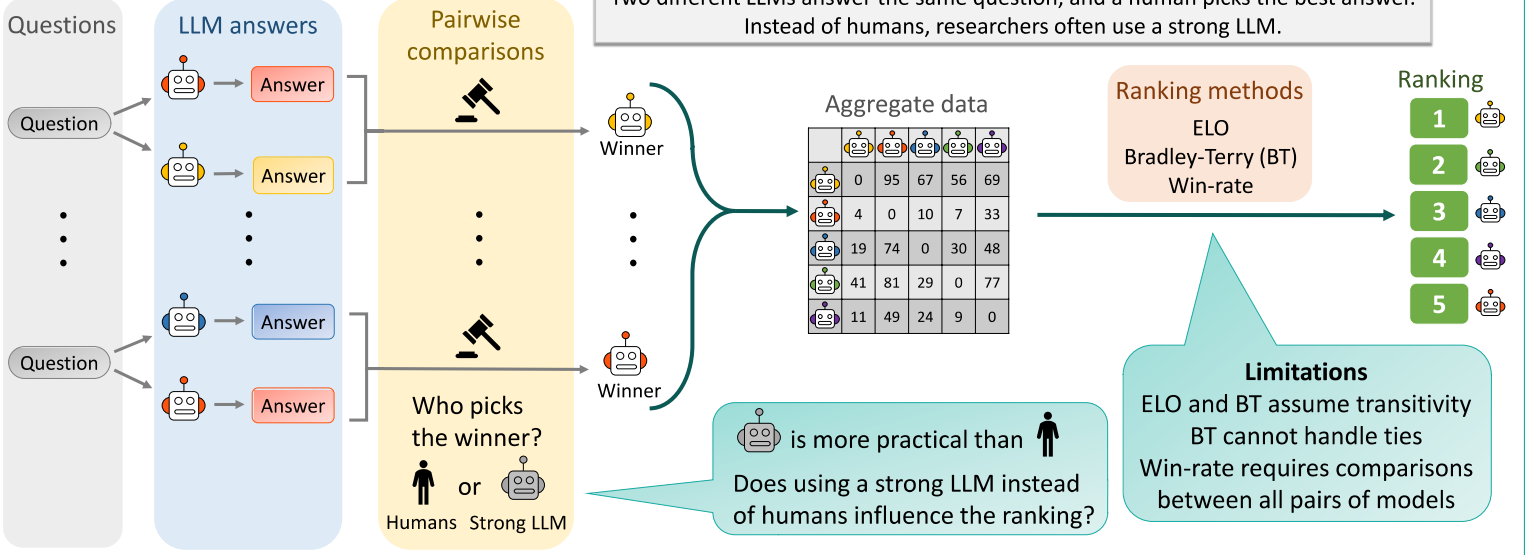
Prediction-Powered Ranking of Large Language Models

Ivi Chatzi Eleni Straitouri Suhas Thejaswi Manuel Gomez-Rodriguez

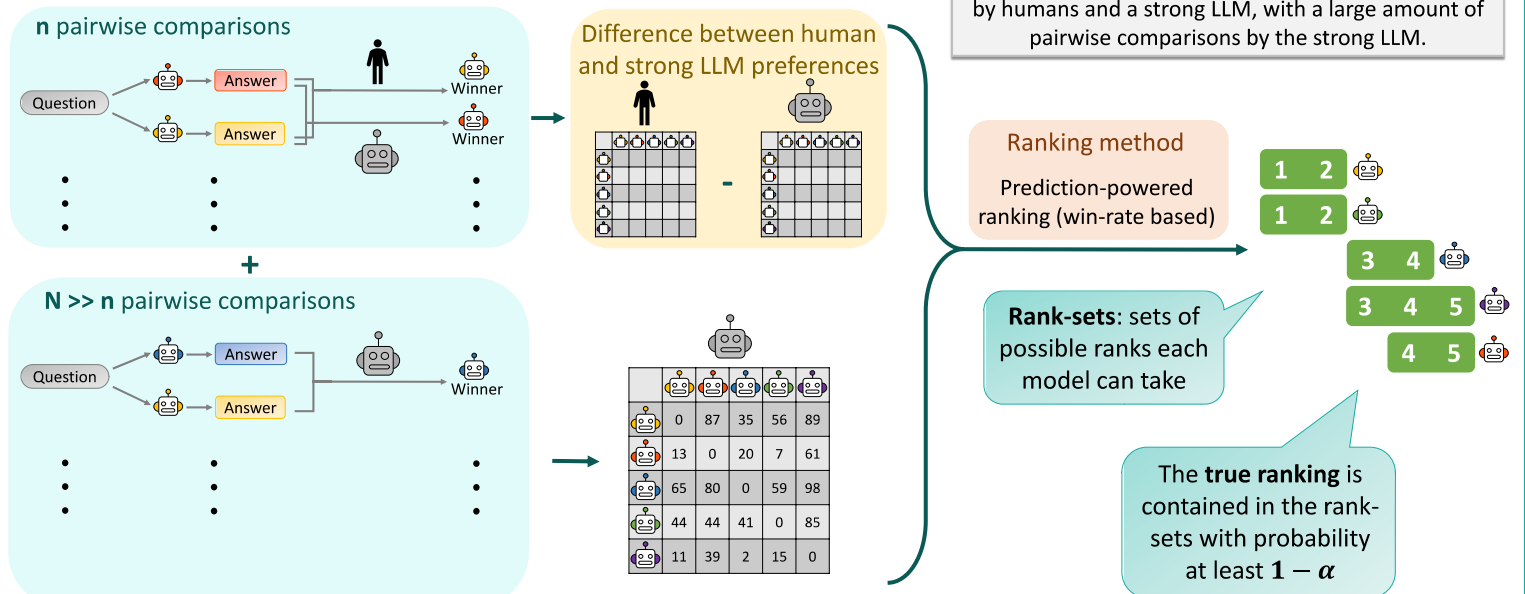


Max Planck Institute for Software Systems

Existing methods

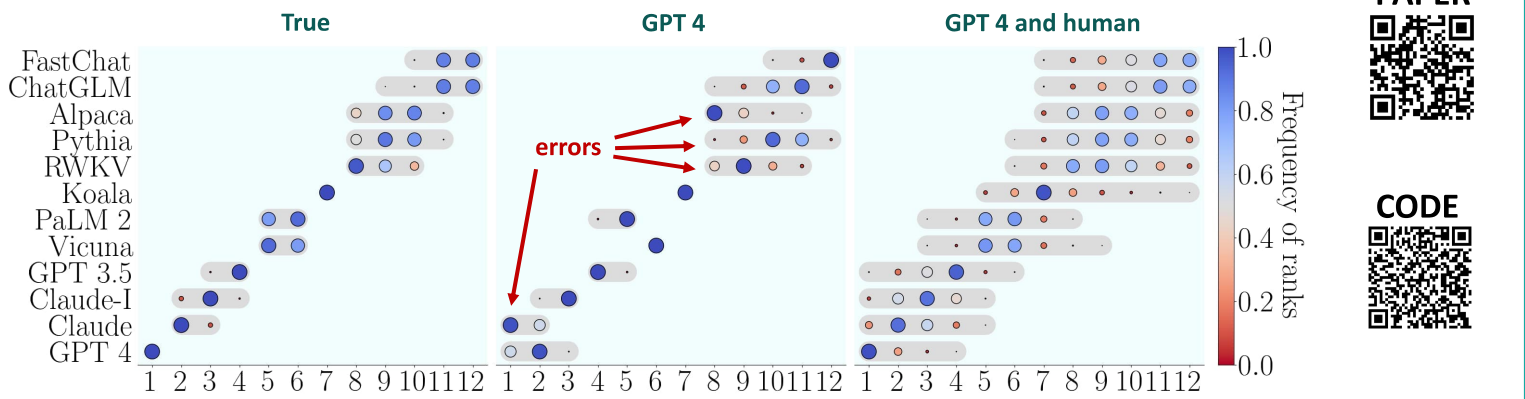


Our framework: Prediction-powered ranking



Rank-sets in practice

How often does the rank-set of each model contain each rank?



Rank-sets using all pairwise comparisons by humans

Rank-sets using only pairwise comparisons by GPT 4 have errors

Rank-sets using many pairwise comparisons by GPT 4 and a few by humans are correct but larger

PAPER



CODE

