Elo Uncoverec: **<u>Robustness</u>** and <u>Best Practices</u> in Language Model Evaluation

Meriem Boubdir*, Edward Kim⁺, Beyza Ermis*, Sara Hooker*, Marzieh Fadaee*

<u>Motivation</u>

- The evaluation of language models constitutes a critical yet resource-intensive undertaking, requiring substantial time, cost, and specialized human expertise.
- Despite advancements in automated metrics, human feedback remains key to assessing language models performance.
- The Elo rating system, designed for dynamic games like chess, provide a robust, dynamic, and interpretable framework for comparing models' capabilities.
- LLMs, unlike dynamic competitors that evolve with time, have **static** capabilities and operate in a time-agnostic context. This static nature prompts a critical investigation of Elo's suitability and reliability when applied to LLM evaluations.

The Elo Rating System

- The rating system assigns scalar values to players' skill level, adjusting dynamically based on match outcomes in a zero-sum manner.
- Each match outcome affect future rankings by updating the ratings according to:

 $E_A = \frac{1}{1 + 10^{(R_B - R_A)/400}}$

Initial Elo scores Player A: 1950 Player B: 2000



 $R'_A = R_A + K(S_A - E_A)$

Updated Elo scores Player A: 1943 (-7pts) Player B: 2007 (+7pts)

Desirable Properties of Elo

Our research focuses on two desirable properties of the Elo rating system: 1. Reliability:	Rank	Model	Elo Rating
sensitivity to hyperparameters (K-Factor)	1	o <u>GPT-4</u>	1225
and match ordering.	2	Claude-v1	1195
$A > B$ and $B > C \implies A > C$	3	<u>Claude-instant-</u> <u>v1</u>	1153

References:

[Elo] Arpad E. Elo. The Rating of Chessplayers, Past and Present. Arco Pub., New York, 1978. [LMSYS-ChatBot] Wei lin Chiang, Lianmin Zheng, Lisa Dunlap, Joseph E. Gonzalez, Ion Stoica, Paul Mooney, Sohier Dane, Addison Howard, and Nate Keating. Lmsys - chatbot arena human preference predictions. https://kaggle.com/competitions/lmsys-chatbot-arena, 2024. Kaggle.

*Cohere For AI, [†]Cohere

崮 How Robust Are Elo Scores?

- Methodology: We model LLM evaluations as Bernoulli processes to simulate binary win/loss outcomes between two models **A** and **B**. → Allow us to rigorously test the Elo system behavior under controlled settings, including win rates, match ordering and the K-factor.
- Findings:
- a. Elo ratings display sensitivity to the ordering of matches and hyperparameters, such as the K-factor.
- b. Transitivity fails under some conditions, undermining the reliability of rankings derived from Elo scores.
- c. Volatility in Elo ratings becomes more prominent when the win rates are similar $(P_{win}(A) \approx P_{win}(B))$.



Validation on Real-World Data

- We extend our synthetic analyses to real-world application by using human feedback from the open-source LMSYS-ChatBot dataset.
- We show the Elo scores for 3 models under various permutations N_{perms} and K-factor settings, and *how their final ranking is affected*.



(a) GPT-4-0314 vs. GPT-4-0613 and GPT-4-0613 vs. Claude-2.1 **Recorded Win rates**: 0.51 vs 0.49 and 0.61 vs 0.39



(b) Elo Scores Averaged Over 100 Permutations

			1000	7744421	10000	- 400
0.9	304	104 390	393	404	426	- 300 🔒
() 0.7	151	170	173	174	175	- 200 - S
ts B						- 100 g
(A bea ^{0.6}	77	87	87	88	89	o Differen
Prob 55	26	24	24	26	30	100 es
0						200 တိ
0.51	14	14	15	18	24	
	1	8	16	32	64	400
			K-factor			









1) Elo Score	1480
Stability	[S ^A]
	eJoog
	o 1420 山
	1400
	1380

2) K-Factor Tuning

If Wins(A) \approx Wins(B): low K values

Else: large K values for a faster convergence of Elo

3) Transitivity

Can be vulnerable for Elo ratings. Final ranking depends on recorded win rates and choice of hyperparameters.

Guidelines for Robust Elo Ratings



