
A Unified Perturbation Framework for Analyzing Leaderboard Stability and Manipulation

Hosna Oyarhoseini¹ Jimmy Lin¹ Amir-Hossein Karimi¹

Abstract

Evaluation leaderboards such as LMArena play a central role in benchmarking large language models by aggregating pairwise human preferences into model rankings, yet the robustness of these rankings remains poorly understood. We present a unified perturbation framework for analyzing Bradley–Terry leaderboards under structured data modifications using influence-based approximations. Our framework studies three match-level perturbations—`Drop`, `Add`, and `Flip`—together with player removal, and evaluates their effects on *top-k membership*, global ranking consistency via *Kendall’s τ* , and *confidence-interval-based uncertainty*. Across Chatbot Arena and six additional pairwise-comparison datasets, we show that modern leaderboards are non-robust across all three objectives: sub-1% targeted perturbations can change the top-ranked model, degrade *Kendall’s τ* , and alter confidence intervals. Beyond robustness auditing, we show that the same influence scores enable efficient targeted perturbations, promoting or demoting specific models and reducing target-model uncertainty with fewer actions than previous manipulation and active-sampling baselines. By summarizing these effects with normalized dataset-level robustness scores, our framework provides a practical and helpful tool for auditing leaderboard stability and motivating more robust evaluation protocols.

1. Introduction

The rapid proliferation of Large Language Models (LLMs) has necessitated the development of scalable, human-centric evaluation frameworks (Frick et al., 2025; Miroyan et al., 2026; Zheng et al., 2023; Guo et al., 2023). Because the

¹University of Waterloo, Waterloo, Ontario, Canada. Correspondence to: Hosna Oyarhoseini <hoarhos@uwaterloo.ca>.

Accepted to the 1st Workshop on Combining Theory and Benchmarks, CTB@ICML 2026, Seoul, South Korea. Copyright 2026 by the author(s).

quality of open-ended generation is subjective and hard to capture with a single absolute metric, modern platforms have increasingly adopted pairwise preference comparisons, where judges choose between two model outputs. In systems such as LMArena/Chatbot Arena, these preferences are primarily collected from human voters (Chiang et al., 2024a), though LLM-based judges are also increasingly used in related evaluation settings (Zheng et al., 2023). This paradigm builds on a classical approach for ranking items from direct comparison outcomes (Bradley & Terry, 1952; Negahban et al., 2017; Wauthier et al., 2013; Shah & Wainwright, 2018), but has recently become a central pillar of LLM benchmarking (Zheng et al., 2023; Chiang et al., 2024a). Most notably, Chatbot Arena (Chiang et al., 2024b) uses Bradley–Terry scores derived from crowdsourced votes to produce model rankings. Beyond leaderboard evaluation, Bradley–Terry models are also underpinned to reward model training for RLHF (Ouyang et al., 2022; Bai et al., 2022; Lee et al., 2024; Touvron et al., 2023; Xu et al., 2024; Sun et al., 2025) and have been used to route queries to suitable LLMs or inference-time scaling strategies (Damani et al., 2025). As these leaderboards increasingly shape model adoption and industry recognition (Metz, 2025; Kruppa, 2024; Singh et al., 2025), they are often treated as stable estimates of true model skill.

Recent work challenges this assumption (Huang et al., 2026; Min et al., 2025; Huang et al., 2025a; Singh et al., 2025; Wu et al., 2022). Huang et al. (2026) show that removing only a small number of pairwise comparisons can alter top rankings, while other studies identify vulnerabilities from injected votes (Min et al., 2025), gamed LLM judges (Zheng et al., 2025; Raina et al., 2024), apathetic or arbitrary annotators (Zhao et al., 2025), and data leakage or selective reporting (Singh et al., 2025). However, existing analyses typically focus on a single *perturbation type* or ranking objective. This leaves open a broader question: *how do different small, structured changes to the comparison dataset propagate through different leaderboard outcomes?* We address this question by building on (Huang et al., 2026)’s idea with a unified influence-based perturbation framework for Bradley–Terry leaderboards. Our framework treats dataset modifications as structured interventions and propagates their effects through the ranking estimator to downstream

leaderboard conclusions. This allows us to study both robustness failures and targeted interventions within the same formalism, including match-level perturbations and player-level removal motivated by model deprecation or exclusion.

Our main contributions are:

- We introduce a unified influence-based framework for auditing Bradley–Terry leaderboards under three match-level perturbations—`Drop`, `Add`, and `Flip`—instantiated for three ranking objectives: *top-k membership*, *Kendall’s τ* for global consistency, and *confidence-interval-based uncertainty* for the reliability of estimated skills.
- Across seven pairwise-comparison datasets, we show that leaderboard non-robustness is systematic across datasets, *perturbation type*, and ranking criteria: fewer than 1% of comparisons substantially affect *top-k membership*, *Kendall’s τ* , and CI-based stability. We also aggregate influence-guided failures into normalized dataset-level robustness scores, giving a compact audit profile across Top- k , CI-aware, and global-ranking views.
- The same framework supports targeted interventions, promoting or demoting specific models with fewer actions than prior vote-manipulation (Min et al., 2025) and identifying matchups that reduce uncertainty more effectively than Chatbot Arena-style active sampling (Chiang et al., 2024b).
- A player-removal analysis shows that removing influential players can induce broad reordering, highlighting model deprecation as a source of the leaderboard illusion (Singh et al., 2025).

Overall, this work provides a unified perspective on leaderboard non-robustness by bridging robustness analysis and adversarial manipulation, revealing fundamental limitations in current leaderboard designs and highlighting the need for more reliable and trustworthy benchmarking methodologies.

2. Preliminaries

2.1. Pairwise comparison and Bradley–Terry framework

Modern LLM leaderboards rank models from pairwise preference data rather than absolute scores. In this setting, each *player* denotes an entity being ranked; for LLM leaderboards, a player corresponds to an LLM or model. A *match* or *vote* denotes one pairwise comparison between two players, typically obtained when a human or judge compares two model responses to the same prompt and selects the preferred output. These comparisons are commonly modeled using the Bradley–Terry (BT) framework (Bradley & Terry, 1952), which estimates latent skill scores from pairwise

outcomes.

We adopt the notation of (Chiang et al., 2024a), defining a set of M models $\mathcal{M} = \{m_1, \dots, m_M\}$ and a dataset $D = \{z_n\}_{n=1}^N$. Each observation $z_n = (x_n, y_n)$ compares two players (i_n, j_n) , where $x_n = e_{i_n} - e_{j_n} \in \mathbb{R}^M$ encodes the matchup and $y_n \in \{0, 1\}$ indicates whether i_n beats j_n . For a comprehensive breakdown of the data format and our specific protocol for handling ties, please refer to Appendix B.1.

The BT model assigns each model m_i a scalar latent skill coefficient $\theta_i \in \mathbb{R}$. For comparison n , let (i, j) denote the ordered pair of models being compared, and let $y_n = 1$ indicate that model i is preferred to model j . The Bradley–Terry probability is

$$\begin{aligned} p_n &:= P(y_n = 1 \mid i, j; \theta) \\ &= \sigma(x_n^\top \theta) = \sigma(\theta_i - \theta_j) \\ &= \frac{1}{1 + e^{-(\theta_i - \theta_j)}}. \end{aligned} \quad (1)$$

To ensure identifiability against constant shifts, we fix $\theta_1 = 0$. The parameter vector θ is estimated by minimizing the empirical binary cross-entropy loss:

$$\begin{aligned} \hat{\theta} &= \arg \min_{\theta} \sum_{n=1}^N \ell(z_n; \theta) \\ &= \arg \min_{\theta} \sum_{n=1}^N \left[-y_n \log p_n \right. \\ &\quad \left. - (1 - y_n) \log(1 - p_n) \right]. \end{aligned} \quad (2)$$

2.2. Influence functions and propagation to objectives

Previous work has used influence scores to identify high-impact deletions under a fixed budget (Broderick et al., 2025). Building on this idea, we study how perturbing individual BT comparisons affects both the estimated skills and downstream leaderboard conclusions. Let $w \in \mathbb{R}^N$ denote match weights and define the weighted BT estimator

$$\hat{\theta}(w) = \arg \min_{\theta: \theta_1=0} \sum_{n=1}^N w_n \ell(z_n; \theta), \quad (3)$$

where $w = \mathbf{1}$ gives the full-data fit and $w_n = 0$ removes match z_n . Let $\hat{\theta} = \hat{\theta}(\mathbf{1})$, $g_n = \nabla_{\theta} \ell(z_n; \hat{\theta})$, and $H = \nabla_{\theta}^2 \sum_{n=1}^N \ell(z_n; \hat{\theta})$. Influence functions approximate the local effect of changing a match weight without refitting (Koh & Liang, 2017):

$$\left. \frac{\partial \hat{\theta}(w)}{\partial w_n} \right|_{w=\mathbf{1}} = -H^{-1} g_n. \quad (4)$$

We propagate this parameter sensitivity to a scalar leaderboard objective $f(\hat{\theta}(w), w)$, such as a *top-k membership*,

Kendall’s τ surrogate, or *confidence-interval-based uncertainty* criterion. Allowing explicit dependence on w captures objectives, such as uncertainty, that change directly with the weighted comparison graph. The objective-level influence of match z_n is

$$I_n^{(f)} := \nabla_{\theta} f(\hat{\theta}(w), w)^{\top} \frac{\partial \hat{\theta}(w)}{\partial w_n} + \frac{\partial f(\hat{\theta}(w), w)}{\partial w_n},$$

$$f(\hat{\theta}(w + \Delta w), w + \Delta w) \approx f(\hat{\theta}(w), w) + \sum_{n=1}^N \Delta w_n I_n^{(f)}. \quad (5)$$

In the leaderboard setting, [Huang et al. \(2026\)](#) instantiate this idea with *top-k membership* objectives $f(\theta) = \theta_i - \theta_j$ across top- k boundaries, showing that dropping a small number of preference votes can change top-ranked membership. We generalize this view by decoupling the perturbation mechanism from the audited objective: the same propagation rule applies to multiple leaderboard objectives and, later, to `Drop`, `Add`, and `Flip` actions. Because our actions are finite rather than infinitesimal, we follow [\(Huang et al., 2026\)](#) use a one-step Newton (1sN) refinement to partially account for curvature changes in the perturbed objective; details are in [Appendix B.2](#).

3. Influence framework

The goal of our framework is to systematically quantify how perturbations to the leaderboard dataset D propagate to downstream ranking objectives, with an overview of the framework shown in [Figure 1](#).

3.1. Action space

We consider three fundamental actions that characterize the ways a leaderboard can be manipulated. Each action corresponds to a specific modification of the weighted M-estimator defined in [Eq. 3](#).

Action 1: Match Dropping (`Drop`). Following prior work that uses influence methods to show that LLM leaderboards are non-robust to small amounts of data removal [\(Huang et al., 2026\)](#), we model dropping a match $z_n \in D$ as setting its weight from $w_n = 1$ to $w_n = 0$. The effect of this single perturbation on a downstream objective is approximated using influence as:

$$\Delta f_{\text{drop},n} \approx -\mathcal{I}_n^{(f)}. \quad (6)$$

Action 2: Match Addition (`Add`). Adding a match involves selecting a candidate z_{new} and increasing its weight from $w_{\text{new}} = 0$ toward 1. Specifically, we define our candidate set as the union of the current dataset (where $w_n = 1$) and all possible pairs (where $w_n = 0$). We then identify which w_{new} to shift from zero to one based on its estimated

influence on the objective:

$$\Delta f_{\text{add,new}} \approx (+1) \cdot \mathcal{I}_{\text{new}}^{(f)}, \quad \text{where}$$

$$\mathcal{I}_{\text{new}}^{(f)} = -\nabla_{\theta} f^{\top} H^{-1} \nabla_{\theta} \ell(z_{\text{new}}; \hat{\theta}) + \left(\frac{\partial f}{\partial w_{\text{new}}} \right). \quad (7)$$

This action models different levels of control over the added data, from benign data collection to stronger manipulation. We consider three candidate spaces of increasing control:

- **all_pairs:** Lowest control. The method selects an unordered pair (i, j) , but the outcome is fixed by the estimated skill ordering, modeling data collection without outcome control.
- **all_outcomes_weighted:** Intermediate control. Both outcomes $(i \succ j$ and $j \succ i)$ are considered separately, and each outcome’s effect is weighted by its BT probability. These probability-weighted outcome scores are used for scoring and selection only.
- **all_outcomes:** Strongest control. Both outcomes are treated as separate candidates, allowing the method to select the outcome with largest effect, modeling stronger manipulation.

Action 3: Outcome Flipping (`Flip`). Reversing the outcome of match z_n to z'_n (where $y'_n = 1 - y_n$) is equivalent to simultaneously dropping the original and adding the reversed match:

$$\Delta f_{\text{flip},n} \approx \mathcal{I}_{z'_n}^{(f)} - \mathcal{I}_{z_n}^{(f)}. \quad (8)$$

We use `Flip` primarily as a counterfactual robustness, explanation, and manipulation action, measuring how a leaderboard conclusion would change if an observed comparison had resolved in the opposite direction [\(Pearl, 2009; Wachter et al., 2018\)](#). This explains which outcomes most support the current ranking and which reversals would most effectively manipulate it.

3.2. Objectives

We characterize leaderboard robustness and sensitivity through three distinct classes of objectives. These scalar functionals, $f(\theta, w)$, allow us to quantify how data perturbations propagate to specific ranking outcomes, from local skill gaps to global ranking consistency and statistical uncertainty.

Top-k Membership Objective / Gap Objective (*top-k membership*). The top- k membership objective (gap objective), following prior work on leaderboard robustness [\(Huang et al., 2026\)](#), measures the difference in estimated skill between two models, which drives rank ordering and top- k membership. For a target pair (i, j) , it is defined as: (the derivative is given in [Appendix B.4](#))

$$f_{\text{gap}}(\theta) = \theta_i - \theta_j. \quad (9)$$

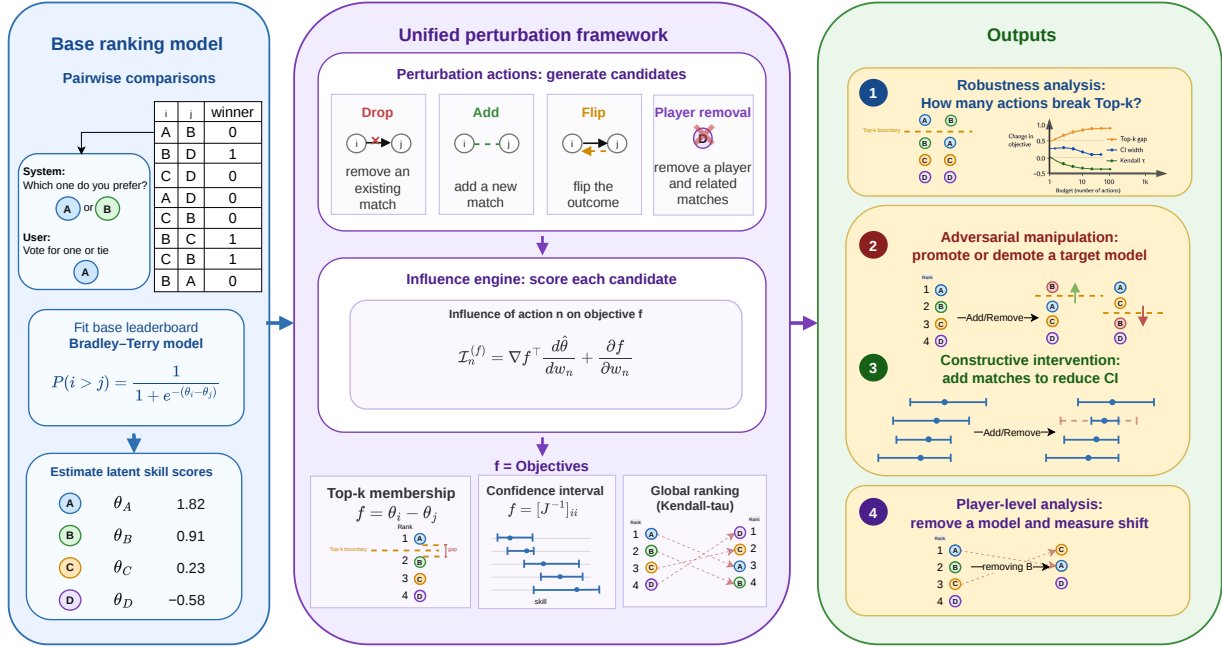


Figure 1. Overview of the framework: influence scores are computed for each action type (**Drop**, **Add**, **Flip**) and propagated through the Bradley–Terry estimator to ranking objectives (*top-k membership*, *confidence-interval-based uncertainty*, *Kendall’s τ*).

Statistical Uncertainty Objectives (*confidence-interval-based uncertainty*). A confidence interval (CI) quantifies uncertainty around an estimated leaderboard skill, indicating how precisely a model’s latent strength is identified from the comparison data. We use a local information approximation motivated by the BT uncertainty analysis of (Gao et al., 2023), where the coordinate-wise uncertainty of the MLE is governed by the inverse local Fisher information of each player.

For each unordered pair (i, j) , let $p_{ij} = \sigma(\hat{\theta}_i - \hat{\theta}_j)$ and $v_{ij} = p_{ij}(1 - p_{ij})$. Let $w_{ij} = \sum_{n: \{i_n, j_n\} = \{i, j\}} w_n$ denote the total comparison weight assigned to the unordered pair $\{i, j\}$. We define the local information of player i as $\hat{\rho}_i^2(w) = \sum_{j \neq i} w_{ij} v_{ij}$. Intuitively, $\hat{\rho}_i^2(w)$ is large when player i has many informative comparisons, leading to smaller uncertainty. Following this local-information view, we use its inverse as a scalable proxy for player-wise variance and define

$$f_{\text{CI-player}}(i; w) = \frac{1}{\hat{\rho}_i^2(w)}, \quad f_{\text{CI-trace}}(w) = \sum_{i=1}^M \frac{1}{\hat{\rho}_i^2(w)}. \quad (10)$$

The first objective measures uncertainty for a target player, which we use for targeted CI reduction, while the second provides a trace-style proxy for global leaderboard uncertainty. These objectives depend both on $\hat{\theta}$ and explicitly on the comparison weights w ; their corresponding direct weight derivatives are given in Appendix B.4. We discuss

the relation between this Gao-style local approximation and the more general sandwich covariance estimator in Appendix B.6.

Global Ranking Consistency (*Kendall’s τ*). To evaluate how perturbations affect the entire leaderboard structure, we use a smooth surrogate of *Kendall’s τ* correlation against a reference ranking π (which is the ranking from BT fitted on the whole initial dataset). Since the discrete *Kendall’s τ* is non-differentiable, we employ a tanh-based relaxation:

$$f_{\tau, T}(\theta; \pi) = \frac{2}{M(M-1)} \sum_{a < b} s_{ab} \tanh\left(\frac{\theta_a - \theta_b}{T}\right), \quad (11)$$

where $s_{ab} \in \{-1, +1\}$ encodes the relative order of items a and b in the reference π , and T is a temperature parameter. As $T \rightarrow 0$, this objective converges to the discrete *Kendall’s τ* . This differentiable form allows us to compute influence scores $\mathcal{I}_n^{(\tau, T)}$ to identify which specific comparisons most heavily affect overall ranking consistency. The derivative is given in Appendix B.4.

3.3. Player-level influence

Following Giordano et al. (2019), we extend influence functions from individual matches to structured group perturbations. For a subset $\mathcal{G} \subset \{1, \dots, N\}$ of matches, such as all matches incident to a given player, the first-order change in

the BT estimator is

$$\begin{aligned} \Delta\theta_{\mathcal{G}} &\approx \sum_{n \in \mathcal{G}} \Delta w_n \frac{\partial \hat{\theta}(w)}{\partial w_n} \\ &= -H(w)^{-1} \sum_{n \in \mathcal{G}} \Delta w_n \nabla_{\theta} \ell(z_n; \hat{\theta}(w)). \end{aligned} \quad (12)$$

We then propagate $\Delta\theta_{\mathcal{G}}$ through a scalar objective $f(\theta, w)$ as in Eq. 5, yielding a player-level influence score from the aggregate contribution of that player’s incident matches.

For player-removal analysis, we additionally use a grouped Newton refinement: after the first-order grouped deletion step in Eq. 12, we take one Newton correction on the remaining comparison data. This is analogous to the 1sN correction for single-match perturbations, but applied jointly to all matches incident to the removed player; details are in Appendix B.3.

4. Experiments

Dataset. Our primary evaluation uses Arena-55K (LMarena AI, 2024) (64 LLMs, 55k human votes), a subset of Chatbot Arena (Chiang et al., 2024a). We additionally evaluate on six pairwise-comparison datasets: Chatbot Arena LLM Judges and MT-Bench Human Evaluation (Zheng et al., 2023), Vision Arena (Chou et al., 2025), WebDev Arena (Vichare et al., 2025), NBA Elo Top-50 (FiveThirtyEight, 2025), and ATP Tennis Top-10 (Sackmann, 2024). Dataset details are in Appendix A.1.

Experimental protocol. Unless otherwise stated, we greedily select perturbations using influence scores under a fixed budget, without refitting during selection. We re-fit after the selected perturbations to report the exact post-perturbation leaderboard and success conditions, though these checks could be monitored by influence estimates. Full algorithms are provided in Appendix B.5.

4.1. Match analysis

We study how individual matches affect leaderboard outcomes along two axes: (i) a *robustness* question: how many perturbations are needed before a criterion changes?, and (ii) a *manipulation* question: can we steer ranking or confidence towards a desired outcome?

4.1.1. ROBUSTNESS ANALYSIS

Top- k membership robustness Following Huang et al. (2026), we measure the minimum number of match perturbations required to change the composition of the top- k set. We instantiate our framework with the gap objective f_{gap} and score all actions candidates by their predicted effect

on crossing a top- k boundary. Starting from the original leaderboard, we greedily apply the highest-influence action until the top- k set changes or until reaching a budget of 5% of the dataset is reached.

Table 1 reports the minimum influence-guided actions needed to change the Top-1 model. On major LLM leaderboards, fewer than 1% targeted actions suffice (e.g., 3 `Flip` or 5 `Drop` on Arena 55k). Across datasets, `Flip` is the most action-efficient perturbation, follows by `Drop` while `Add` variants typically require larger budgets or remain robust within the budget. We ablate the effect of varying k on Arena 55k in Appendix E.2.

CI-aware Top- k membership robustness We extend the Top- k robustness analysis to incorporate estimation uncertainty. For a chosen boundary rank k , we augment the usual membership objective with confidence bounds: we fix the boundary pair consisting of the model currently ranked k and the model currently ranked $k + 1$, and compare the upper confidence bound of the rank- k model to the lower confidence bound of the rank- $(k + 1)$ model. A perturbation is declared successful only if, after refitting, the rank- $(k + 1)$ model is not merely above the rank- k model in point estimate, but is *strictly* above under uncertainty as well, meaning that its lower CI bound exceeds the rank- k model’s upper CI bound. This strict CI-aware criterion is therefore stronger than an ordinary non-CI-aware (point-estimate) Top- k membership; the algorithm and an example are given in Appendix B.5, D.3.

Table 2 compares the minimum number of influence-guided actions needed under point-estimate and CI-aware criteria, using per dataset boundary rank k for each dataset. Accounting for uncertainty increases the required budget, but does not make the leaderboard robust. CI-aware criteria require on average $13\times$ more actions than point-estimate changes, yet several datasets remain non-robust within the allowed budget. `Drop` and `Flip` remain the most effective actions, while `Add` variants often require larger budgets or remain robust within the budget.

Global ranking robustness (Kendall’s τ). Beyond local boundary changes, we score candidates using the f_{τ} surrogate and apply actions greedily, comparing against random selection (Figure 2, left). Influence-guided `Flip` causes the steepest τ degradation, while `Drop` is noticeably more conservative and `Add` variants cluster at intermediate degradation. In all cases, random selection barely moves τ , confirming that influence-guided selection is essential, random perturbations are insufficient to disrupt global ranking even at budget 30. Results for all datasets are in Appendix D.2.

A Unified Perturbation Framework for Leaderboard Stability and Manipulation

Table 1. Minimum influence-guided actions needed to change the Top-1 model under a 5% budget. Percentages are relative to total matches; robust means no change within budget.

Dataset	Dataset size	Drop	Flip	Add-pairs	Add-outcomes	Add-weighted
Arena 55k	57,477	5 (0.01%)	3 (0.01%)	9 (0.02%)	9 (0.02%)	9 (0.02%)
Arena LLM-J	49,938	9 (0.02%)	6 (0.01%)	22 (0.04%)	21 (0.04%)	21 (0.04%)
MT-Bench	3,355	92 (2.74%)	46 (1.37%)	robust	robust	robust
NBA Top-50	109,892	24 (0.02%)	15 (0.01%)	77 (0.07%)	55 (0.05%)	57 (0.05%)
ATP Top-10	278	6 (2.16%)	3 (1.08%)	robust	9 (3.24%)	14 (5.00%)
Vision Arena	29,849	50 (0.17%)	25 (0.08%)	robust	robust	robust
WebDev Arena	10,501	179 (1.70%)	12 (0.11%)	robust	robust	robust

Table 2. Minimum number of actions needed to change the top- k boundary under point-estimate / CI-aware criteria (percentages relative to total matches) with 95% confidence level and a 5% perturbation budget. Across datasets, the CI-aware criterion generally requires substantially more actions than the point-estimate criterion. robust: target not changed within the 5% budget.

Dataset	Dataset size	k	Drop (point-estimate / CI-aware)	Flip (point-estimate / CI-aware)	Add-pairs	Add-outcomes	Add-weighted
Arena 55k	57,477	22	2 (0.00%) / 19 (0.03%)	1 (0.00%) / 12 (0.02%)	2 (0.00%) / 39 (0.07%)	2 (0.00%) / 26 (0.05%)	3 (0.01%) / 28 (0.05%)
LLM Judge Arena	49,938	31	1 (0.00%) / 31 (0.06%)	1 (0.00%) / 17 (0.03%)	1 (0.00%) / 50 (0.10%)	1 (0.00%) / 32 (0.06%)	1 (0.00%) / 45 (0.09%)
MT-Bench	3,355	2	13 (0.39%) / 62 (1.85%)	7 (0.21%) / 33 (0.98%)	robust / robust	robust / robust	robust / robust
NBA Top-50	109,892	8	2 (0.00%) / 46 (0.04%)	1 (0.00%) / 25 (0.02%)	5 (0.00%) / robust	3 (0.00%) / robust	6 (0.01%) / robust
ATP Top-10	278	8	1 (0.36%) / robust	1 (0.36%) / 7 (2.52%)	1 (0.36%) / robust	1 (0.36%) / robust	1 (0.36%) / robust
Vision Arena	29,849	14	4 (0.01%) / 43 (0.14%)	3 (0.01%) / 33 (0.11%)	6 (0.02%) / robust	5 (0.02%) / robust	10 (0.03%) / robust
WebDev Arena	10,501	12	3 (0.03%) / 342 (3.26%)	2 (0.02%) / 23 (0.22%)	7 (0.07%) / robust	6 (0.06%) / robust	10 (0.10%) / robust

Confidence-interval-based uncertainty robustness. We next evaluate how perturbations affect global leaderboard uncertainty. Using $f_{\text{CI-trace}}$ to score candidates and applying actions greedily (Figure 2, right), the action-type ordering reverses: `Add` variants yield the steepest uncertainty reduction, reaching roughly -1% of the initial trace by budget 25, about an order of magnitude larger than `Drop` or `Flip`. In contrast, random additions slightly increase uncertainty, as they tend to introduce less-informative comparisons. Results for all datasets are in Appendix D.2.

From influence scores to leaderboard audits Beyond identifying non-robustness, we define normalized dataset-level robustness metrics that summarize influence-guided failures as audit signals. These metrics measure either the budget needed to change a leaderboard conclusion, as in Top-1, or the normalized magnitude of change, as in Kendall- τ and confidence interval stability, Appendix D.1 gives the full metric definitions. Table 3 reports these scores for all seven dataset. Lower values indicate less robustness. Based on R_{all} , MT-Bench is the most robust dataset, followed closely by NBA Top-50 and the arena-style LLM datasets, whereas ATP Top-10 and WebDev Arena are least robust overall due to lower global-ranking stability.

Top- k entry and removal. We use Top- k entry and removal as targeted manipulation tasks: promotion moves an outside model into the Top- k set, while demotion pushes an inside model out. We evaluate both under a sequentially revealed stream of candidate comparisons. For each dataset, we partition the leaderboard into top, middle, and lower skill regions and sample one boundary rank from each region; promotion targets rank $(k + 1)$ and demotion targets rank k . We compare our influence-guided policy with an omni-

rigging baseline (Min et al., 2025) under the same exposed stream, where each method may discard a pair, accept either directional outcome, or encode a tie; full details are in Appendix B.5. Our method uses a dynamic boundary-gap objective, recomputing the boundary opponent after each accepted perturbation and selecting the action with largest predicted influence on the current gap.

Figure 3 shows that the influence-guided policy typically achieves targeted promotion and demotion with fewer actions than the omni-rigging baseline across most datasets. Across datasets, influence-guided selection uses about 74 actions on average, compared with about 92 for omni-rigging, a roughly 19% reduction. This indicates that directly optimizing the current top- k boundary gap is more efficient than locally greedy rank manipulation.

CI reduction for a target model From a leaderboard provider’s perspective, we study whether the framework can constructively reduce the confidence interval of a target model m^* . We use the target uncertainty objective $f_{\text{CI-player}}(m^*)$ and restrict the action space to `Add`, selecting candidate additions with the most negative predicted influence. We compare influence-guided `Add` with the `all_pairs` candidate space against *Random* and *Arena Active* (Chiang et al., 2024b); *Random* samples pairs uniformly, while *Arena Active* prioritizes target-involving uncertain comparisons. After selecting a pair, all methods assign the outcome by the current BT skill ordering, refit after each addition, and measure target-model CI reduction.

For each dataset, we choose three target players from the top, middle, and lower skill regions and run each method for 12 `Add` actions. Table 4 reports $\% \Delta \text{CI} = 100 \times (\text{final} - \text{initial}) / \text{initial}$, averaged over targets; more negative values

A Unified Perturbation Framework for Leaderboard Stability and Manipulation

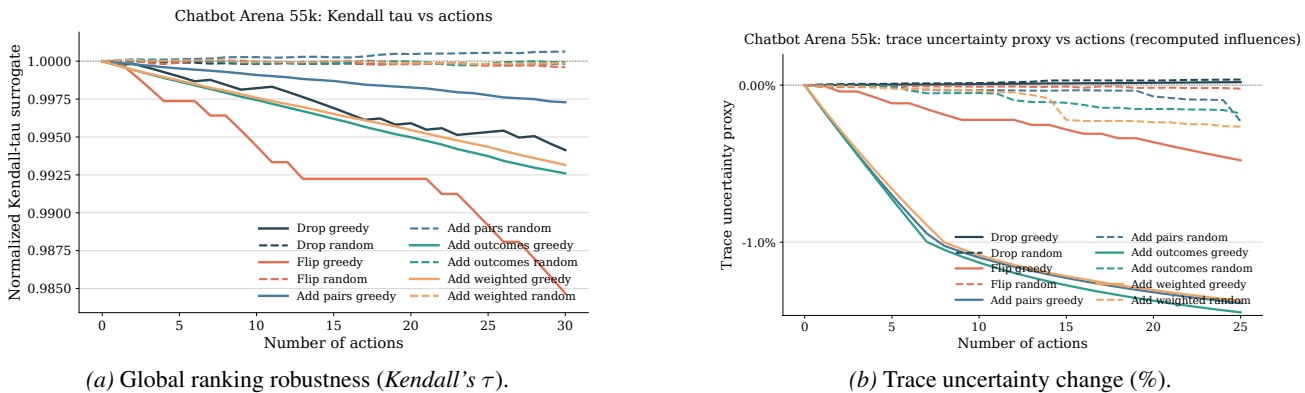


Figure 2. Global ranking and uncertainty robustness on Chatbot Arena 55k under budgets of 30 actions for global ranking and 25 for uncertainty (solid = influence-guided, dashed = random). **Left:** Influence-guided `Flip` causes the largest degradation ($\tau \approx 0.985$), while `Drop` is more conservative ($\tau \approx 0.994$); random baselines remain near 1.0. **Right:** Influence-guided `Add` reduces trace uncertainty by up to $\approx -1\%$, whereas random actions slightly increase it.

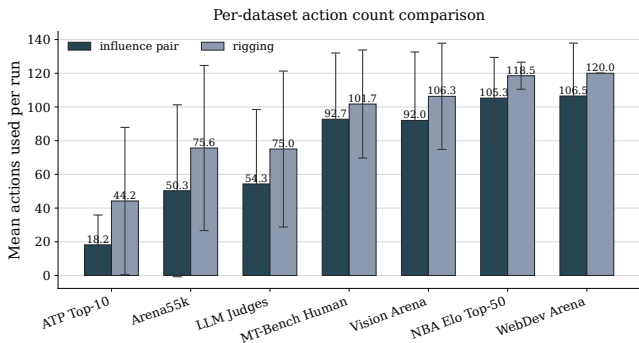


Figure 3. Mean intervention actions per run for pairwise influence versus the rigging baseline. Error bars show variability across repeated trials. Each run is capped at 120 exposed pairs.

indicate larger uncertainty reduction. Averaged over all results, influence-guided sampling produces a larger CI-width reduction than *Arena Active* and *Random*, showing directly optimizing the target CI objective is more informative than Arena-style heuristics.

4.2. Player removal analysis

Beyond individual match perturbations, we consider player removal: dropping all matches of a given model, as occurs when models are retired, deprecated, or excluded from evaluation. Prior work on the leaderboard illusion (Singh et al., 2025) shows that excluding models can qualitatively change which remaining models appear superior; we quantify this effect systematically. We use our influence framework to rank all players by predicted $|\tau|$ -influence and validate by removing the most influential player in each dataset and measuring the resulting *Kendall's tau* shift. Table 5 shows that removing an influential player can induce broad reordering, moving up to 28 players, shifting ranks up to 13 positions, and changing up to 6 top-10 memberships. Connectivity is the strongest and most consistent predictor (Appendix E.3).

Table 3. Dataset-level robustness scores across the Top-1, CI views, and global ranking. Lower values indicate less robustness under a 5% data budget. **Bold** marks the most robust dataset.

Dataset	$R_{\text{Top-1}}$	R_{CI}	R_{τ}	R_{all}
Arena 55k	0.001	0.976	0.993	0.656
Arena LLM-J	0.000	0.975	0.992	0.656
MT-Bench	0.048	0.955	0.985	0.662
NBA Top-50	0.003	0.992	0.982	0.659
ATP Top-10	0.072	0.810	0.474	0.452
Vision Arena	0.003	0.976	0.991	0.657
WebDev Arena	0.004	0.579	0.782	0.455

5. Related work

Ranking systems and evaluation leaderboards. The Bradley-Terry model (Bradley & Terry, 1952; Hamilton et al., 2025; Fang et al., 2026) and related models like Elo rating systems (Boubdir et al., 2024; Liu et al., 2025), form the basis of modern probabilistic ranking. Maximum-likelihood estimation and uncertainty quantification for BT models are well-studied (Hunter, 2004; Bong & Rinaldo, 2022; Gao et al., 2023; Fan et al., 2024), with spectral and sorting-based methods providing efficient alternatives (Negahban et al., 2017; Wauthier et al., 2013). In LLM evaluation, pairwise preference aggregation has become central through Chatbot Arena (Chiang et al., 2024b), with recent work improving statistical reliability (Ameli et al., 2025; Gao et al., 2025) and studying evaluation bias (Daynauth et al., 2025; Levtsov & Ustalov, 2025).

Leaderboard robustness and manipulation. Benchmark rankings can be sensitive to dataset shifts, benchmark overfitting, and small comparison perturbations (Boubdir et al., 2024; Huang et al., 2026; Singh et al., 2025). The leaderboard illusion (Singh et al., 2025) shows that benchmark rankings can be distorted by selective disclosure, private

Table 4. Percent change in target-model CI width after 12 targeted actions, $\% \Delta \text{CI} = 100 \times (\text{final} - \text{initial}) / \text{initial}$. More negative is better.

Dataset	Influence-based	Arena Active	Random
MT-Bench	-0.27 ± 0.05	-0.25 ± 0.07	-0.08 ± 0.03
WebDev Arena	-0.28 ± 0.05	-0.39 ± 0.23	-0.02 ± 0.00
Vision Arena	-0.17 ± 0.12	-0.16 ± 0.11	-0.01 ± 0.01
Arena LLM-J	-0.47 ± 0.25	-0.30 ± 0.11	-0.00 ± 0.00
Arena 55k	-0.72 ± 0.22	-0.54 ± 0.28	-0.01 ± 0.01
NBA Top-50	-0.33 ± 0.09	-0.29 ± 0.05	-0.01 ± 0.00
All datasets	-0.37 ± 0.13	-0.32 ± 0.14	-0.02 ± 0.01

testing, uneven data allocation, and model inclusion or exclusion decisions; while vote manipulation and adversarial leaderboard attacks expose additional manipulation risks (Min et al., 2025; Huang et al., 2025b; Suri et al., 2026). Related theory studies robustness of estimators to sample removal (Azar et al., 2025; Broderick et al., 2025), and recent uncertainty-aware ranking methods propose rank sets or confidence diagrams (Chatzi et al., 2024; Wang et al., 2025). Our work provides a unified framework covering match addition, removal, and flipping, extending prior work that typically studies one *perturbation type* or one objective.

Explainability and data attribution. Data attribution is a central tool for explaining model behavior by assigning importance scores to individual training examples or groups of examples. Existing approaches include Data Shapley (Ghorbani & Zou, 2019), TracIn (Pruthi et al., 2020), and TRAK (Park et al., 2023), which quantify how training data contributes to model predictions or learned representations; see Hammoudeh & Lowd (2024) for a survey. Influence functions are another classical approach to data attribution, estimating the effect of training points on model parameters without retraining (Koh & Liang, 2017; Pregibon, 1981), with roots in robust statistics (Huber, 1967; Freedman, 2006). Extensions include infinitesimal jackknife (Giordano et al., 2019), group influence methods (Koh et al., 2019), second-order approximations (Basu et al., 2020), and analyses of when influence estimates are reliable (Bae et al., 2022). We build on this line of work by propagating influence scores to leaderboard-specific objectives.

6. Conclusion, limitations, and future work

We introduced a unified influence-based perturbation framework for evaluating Bradley–Terry leaderboard stability and manipulation. Across seven datasets, we show that leaderboard conclusions are fragile: Top-1 membership changes with fewer than 1% targeted actions on major LLM leaderboards, CI-aware Top- k changes require larger budgets but remain vulnerable, and removing highly connected play-

Table 5. Most-influential-player ablation results. We remove the player with largest influence and all associated matches, then measure changes.

Dataset	$\Delta \tau$	Moved	Max Shift	Top-10	Removed %
ATP Top-10	-0.167	2	2	2	24.1%
MT-Bench	0.000	0	0	0	30.3%
WebDev	-0.091	5	2	4	23.2%
Vision	-0.050	2	2	0	14.1%
Arena LLM-J	-0.023	25	5	4	6.2%
Arena 55k	-0.030	28	10	5	6.2%
NBA Top-50	-0.053	27	13	6	0.5%

ers can induce broad rank shifts. Our normalized dataset-level robustness scores summarize these effects across Top- k , CI-aware, and global-ranking views. Overall, expert-curated MT-Bench is generally more robust, while crowd-sourced arena-style leaderboards are more sensitive; among actions, `Flip` is usually most effective for changing rankings, whereas `Add` is most useful for reducing uncertainty through targeted data collection. Beyond auditing, the same framework enables targeted manipulation: influence-guided actions promote or demote selected models with fewer interventions than the omni-rigging baseline, while supporting constructive uncertainty reduction through comparisons that narrow confidence intervals more effectively than Arena-style sampling.

Our analysis has several limitations. Influence estimates are local approximations and may be less accurate under large perturbation budgets or highly nonlinear ranking changes. Our uncertainty objectives use scalable BT-based proxies, which do not capture all sources of annotation noise, judge bias, prompt-level dependence, or model-specific evaluation artifacts. Future work should study higher-order attribution methods, outlier-robust BT estimators, and sample-complexity guarantees for leaderboard stability. Another important direction is to better understand what makes a match, player, or dataset highly influential or robust. Finally, the framework can be extended as a defensive tool, where low robustness scores or highly influential comparisons trigger conservative reporting, targeted data collection, or delayed updates, rather than being treated as evidence of invalid data.

Impact Statement

This paper studies the robustness and manipulability of pairwise-comparison leaderboards used for benchmarking large language models and other ranked systems. Our findings show that widely used leaderboards can be perturbed with small, targeted modifications, which has dual-use implications: the same techniques that enable auditing for fragility could in principle be used to manipulate rankings.

We highlight this tension explicitly and frame the framework as primarily an auditing and defensive tool, where low robustness scores or highly influential comparisons can motivate conservative reporting, targeted data collection, or delayed updates. More broadly, we believe that exposing these vulnerabilities promotes more reliable evaluation protocols and better-calibrated trust in benchmark-driven model adoption decisions. There are many potential societal consequences of advancing the field of Machine Learning, none which we feel must be specifically highlighted here beyond the considerations above.

References

- Ameli, S., Zhuang, S., Stoica, I., and Mahoney, M. W. A statistical framework for ranking LLM-based chatbots. In *The Thirteenth International Conference on Learning Representations*, 2025. <https://openreview.net/forum?id=rAoEub6Nw2>.
- Azar, E., Feldman, M. J., and Nadler, B. Robustness of OLS to sample removals: Theoretical analysis and implications. *arXiv preprint arXiv:2512.23069*, 2025.
- Bae, J., Ng, N., Lo, A., Ghassemi, M., and Grosse, R. If influence functions are the answer, then what is the question? In *Advances in Neural Information Processing Systems*, volume 35, pp. 17953–17967, 2022.
- Bai, Y., Jones, A., Ndousse, K., Askell, A., Chen, A., Das-Sarma, N., Drain, D., Fort, S., Ganguli, D., Henighan, T., et al. Training a helpful and harmless assistant with reinforcement learning from human feedback. *arXiv preprint arXiv:2204.05862*, 2022.
- Basu, S., You, X., and Feizi, S. On second-order group influence functions for black-box predictions. In *Proceedings of the 37th International Conference on Machine Learning*, volume 119 of *Proceedings of Machine Learning Research*, pp. 715–724. PMLR, 2020.
- Bong, H. and Rinaldo, A. Generalized results for the existence and consistency of the MLE in the Bradley-Terry-Luce model. In *Proceedings of the 39th International Conference on Machine Learning*, volume 162, pp. 2160–2177. PMLR, 2022.
- Boubdir, M., Kim, E., Ermis, B., Hooker, S., and Fadaee, M. Elo uncovered: Robustness and best practices in language model evaluation. In *Advances in Neural Information Processing Systems*, volume 37, pp. 106135–106161. Curran Associates, Inc., 2024. doi: 10.52202/079017-3367.
- Bradley, R. A. and Terry, M. E. Rank analysis of incomplete block designs: I. the method of paired comparisons. *Biometrika*, 39(3/4):324–345, 1952. doi: 10.2307/2334029.
- Broderick, T., Giordano, R., and Meager, R. An automatic finite-sample robustness metric: When can dropping a little data make a big difference? *Econometrica*, 93(5): 1915–1935, 2025. doi: 10.3982/ECTA21567.
- Chatzi, I., Straitouri, E., Thejaswi, S., and Gómez Rodríguez, M. Prediction-powered ranking of large language models. In *Advances in Neural Information Processing Systems*, volume 37, pp. 113096–113133. Curran Associates, Inc., 2024. doi: 10.52202/079017-3594.
- Chiang, W.-L., Zheng, L., Sheng, Y., Angelopoulos, A. N., Li, T., Li, D., Zhu, B., Zhang, H., Jordan, M., Gonzalez, J. E., and Stoica, I. Chatbot Arena: An open platform for evaluating LLMs by human preference. In *Proceedings of the 41st International Conference on Machine Learning*, volume 235 of *Proceedings of Machine Learning Research*, pp. 8359–8388. PMLR, 2024a.
- Chiang, W.-L., Zheng, L., Sheng, Y., Angelopoulos, A. N., Li, T., Li, D., Zhu, B., Zhang, H., Jordan, M., Gonzalez, J. E., and Stoica, I. Chatbot Arena: An open platform for evaluating LLMs by human preference. In *Proceedings of the 41st International Conference on Machine Learning*, volume 235 of *Proceedings of Machine Learning Research*, pp. 8359–8388. PMLR, 2024b.
- Chou, C., Dunlap, L., Mashita, K., Mandal, K., Darrell, T., Stoica, I., Gonzalez, J. E., and Chiang, W.-L. Vision-Arena: 230k real world user-VLM conversations with preference labels. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pp. 3877–3887, 2025.
- Damani, M., Shenfeld, I., Peng, A., Bobu, A., and Andreas, J. Learning how hard to think: Input-adaptive allocation of LM computation. In *The Thirteenth International Conference on Learning Representations*, 2025.
- Daynauth, R., Clarke, C., Flautner, K., Tang, L., and Mars, J. Ranking unraveled: Recipes for LLM rankings in head-to-head AI combat. In *Proceedings of the 63rd Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*, pp. 26078–26091. Association for Computational Linguistics, 2025. doi: 10.18653/v1/2025.acl-long.1265.
- Fan, J., Hou, J., and Yu, M. Uncertainty quantification of MLE for entity ranking with covariates. *Journal of Machine Learning Research*, 25(358):1–83, 2024.
- Fang, S., Han, R., Luo, Y., and Xu, Y. Recent advances in the Bradley-Terry model: Theory, algorithms, and applications. *arXiv preprint arXiv:2601.14727*, 2026.
- FiveThirtyEight. NBA-ELO dataset. <https://github.com/fivethirtyeight/data/tree/master/nba-elo>, 2025. Accessed 2026-04-30.

- Freedman, D. A. On the so-called ‘‘Huber sandwich estimator’’ and ‘‘Robust standard errors’’. *The American Statistician*, 60(4):299–302, 2006. doi: 10.1198/000313006X152207.
- Frick, E., Chen, C., Tennyson, J., Li, T., Chiang, W.-L., Angelopoulos, A. N., and Stoica, I. Prompt-to-leaderboard: Prompt-adaptive LLM evaluations. In *Proceedings of the 42nd International Conference on Machine Learning*, volume 267 of *Proceedings of Machine Learning Research*, pp. 17672–17689. PMLR, 2025.
- Gao, C., Shen, Y., and Zhang, A. Y. Uncertainty quantification in the Bradley–Terry–Luce model. *Information and Inference: A Journal of the IMA*, 12(2):1073–1140, 2023. doi: 10.1093/imaiai/iaac032.
- Gao, M., Liu, Y., Hu, X., Wan, X., Bragg, J., and Cohan, A. Re-evaluating automatic LLM system ranking for alignment with human preference. In *Findings of the Association for Computational Linguistics: NAACL 2025*, pp. 4605–4629. Association for Computational Linguistics, 2025. doi: 10.18653/v1/2025.findings-naacl.260.
- Ghorbani, A. and Zou, J. Data Shapley: Equitable valuation of data for machine learning. In *Proceedings of the 36th International Conference on Machine Learning*, pp. 2242–2251. PMLR, 2019.
- Giordano, R., Stephenson, W., Liu, R., Jordan, M. I., and Broderick, T. A Swiss army infinitesimal jackknife. In *Proceedings of the Twenty-Second International Conference on Artificial Intelligence and Statistics*, volume 89, pp. 1139–1147. PMLR, 2019.
- Guo, Z., Jin, R., Liu, C., Huang, Y., Shi, D., Supryadi, Yu, L., Liu, Y., Li, J., Xiong, B., and Xiong, D. Evaluating large language models: A comprehensive survey. *arXiv preprint arXiv:2310.19736*, 2023.
- Hamilton, I., Tawn, N., and Firth, D. The many routes to the ubiquitous Bradley-Terry model. *arXiv preprint arXiv:2312.13619*, 2025.
- Hammoudeh, Z. and Lowd, D. Training data influence analysis and estimation: A survey. *Machine Learning*, 113(5):2351–2403, 2024. doi: 10.1007/s10994-023-06495-7.
- Huang, J. Y., Shen, Y., Wei, D., and Broderick, T. Dropping just a handful of preferences can change top large language model rankings. In *The Fourteenth International Conference on Learning Representations*, 2026. <https://openreview.net/forum?id=jNiEMDsRgc>.
- Huang, Y., Nasr, M., Angelopoulos, A., Carlini, N., Chiang, W.-L., Choquette-Choo, C. A., Ippolito, D., Jagielski, M., Lee, K., Liu, K. Z., Stoica, I., Tramèr, F., and Zhang, C. Exploring and mitigating adversarial manipulation of voting-based leaderboards. In *Proceedings of the 42nd International Conference on Machine Learning*, volume 267 of *Proceedings of Machine Learning Research*, pp. 25654–25671. PMLR, 2025a.
- Huang, Y., Nasr, M., Angelopoulos, A., Carlini, N., Chiang, W.-L., Choquette-Choo, C. A., Ippolito, D., Jagielski, M., Lee, K., Liu, K. Z., Stoica, I., Tramèr, F., and Zhang, C. Exploring and mitigating adversarial manipulation of voting-based leaderboards. In *Proceedings of the 42nd International Conference on Machine Learning*, volume 267 of *Proceedings of Machine Learning Research*, pp. 25654–25671. PMLR, 2025b.
- Huber, P. J. The behavior of maximum likelihood estimates under nonstandard conditions. In *Proceedings of the Fifth Berkeley Symposium on Mathematical Statistics and Probability, Volume 1: Statistics*, pp. 221–234. University of California Press, 1967.
- Hunter, D. R. MM algorithms for generalized Bradley-Terry models. *The Annals of Statistics*, 32(1):384–406, 2004.
- Koh, P. W. and Liang, P. Understanding black-box predictions via influence functions. In *Proceedings of the 34th International Conference on Machine Learning*, volume 70, pp. 1885–1894. PMLR, 2017.
- Koh, P. W., Ang, K.-S., Teo, H. H. K., and Liang, P. On the accuracy of influence functions for measuring group effects. In *Advances in Neural Information Processing Systems*, volume 32, 2019.
- Kruppa, M. The UC Berkeley project that is the AI industry’s obsession. *The Wall Street Journal*, <https://www.wsj.com/tech/ai/the-uc-berkeley-project-that-is-the-ai-industrys-obsession-bc68b3e3>, 2024. Accessed 2026-04-30.
- Lee, H., Phatale, S., Mansoor, H., Mesnard, T., Ferret, J., Lu, K. R., Bishop, C., Hall, E., Carbune, V., Rastogi, A., and Prakash, S. RLAIF vs. RLHF: Scaling reinforcement learning from human feedback with AI feedback. In *Proceedings of the 41st International Conference on Machine Learning*, volume 235 of *Proceedings of Machine Learning Research*, pp. 26874–26901. PMLR, 2024.
- Levtsov, G. and Ustalov, D. Confidence and stability of global and pairwise scores in NLP evaluation. In *Proceedings of the 63rd Annual Meeting of the Association for Computational Linguistics (Volume 4: Student Research Workshop)*, pp. 40–52. Association for Computational Linguistics, 2025. doi: 10.18653/v1/2025.acl-srw.3.
- Liu, Z., Li, J., Zhuang, Y., Liu, Q., Shen, S., Ouyang, J., Cheng, M., and Wang, S. am-ELO: A stable framework

- for arena-based LLM evaluation. In *Proceedings of the 42nd International Conference on Machine Learning*, volume 267 of *Proceedings of Machine Learning Research*, pp. 38857–38868. PMLR, 2025.
- LMarena AI. Arena human preference 55k dataset. <https://huggingface.co/datasets/lmarena-ai/arena-human-preference-55k>, 2024. Hugging Face dataset; accessed 2026-04-30.
- Metz, R. Before DeepSeek blew up, Chatbot Arena announced its arrival. Bloomberg Businessweek, <https://www.bloomberg.com/news/articles/2025-02-18/before-deepseek-blew-up-one-website-announced-its-arrival>, 2025. Accessed 2026-04-30.
- Min, R., Pang, T., Du, C., Liu, Q., Cheng, M., and Lin, M. Improving your model ranking on Chatbot Arena by vote rigging. In *Proceedings of the 42nd International Conference on Machine Learning*, volume 267 of *Proceedings of Machine Learning Research*, pp. 44252–44271. PMLR, 2025.
- Miroyan, M., Wu, T.-H., King, L., Li, T., Pan, J., Hu, X., Chiang, W.-L., Angelopoulos, A. N., Darrell, T., Norouzi, N., and Gonzalez, J. E. Search Arena: Analyzing search-augmented LLMs. In *The Fourteenth International Conference on Learning Representations*, 2026. <https://openreview.net/forum?id=MMGR1DnhtI>.
- Negahban, S., Oh, S., and Shah, D. Rank centrality: Ranking from pairwise comparisons. *Operations Research*, 65(1):266–287, 2017. doi: 10.1287/opre.2016.1534.
- Ouyang, L., Wu, J., Jiang, X., Almeida, D., Wainwright, C., Mishkin, P., Zhang, C., Agarwal, S., Slama, K., Ray, A., Schulman, J., Hilton, J., Kelton, F., Miller, L., Simens, M., Askell, A., Welinder, P., Christiano, P., Leike, J., and Lowe, R. Training language models to follow instructions with human feedback. In *Advances in Neural Information Processing Systems*, volume 35, pp. 27730–27744, 2022.
- Park, S. M., Georgiev, K., Ilyas, A., Leclerc, G., and Madry, A. TRAK: Attributing model behavior at scale. In *Proceedings of the 40th International Conference on Machine Learning*, volume 202, pp. 27074–27113. PMLR, 2023.
- Pearl, J. Causal inference in statistics: An overview. *Statistics Surveys*, 3:96–146, 2009. doi: 10.1214/09-SS057.
- Pregibon, D. Logistic regression diagnostics. *The Annals of Statistics*, 9(4):705–724, 1981.
- Pruthi, G., Liu, F., Kale, S., and Sundararajan, M. Estimating training data influence by tracing gradient descent. In *Advances in Neural Information Processing Systems*, volume 33, pp. 19920–19930, 2020.
- Raina, V., Liusie, A., and Gales, M. Is LLM-as-a-judge robust? investigating universal adversarial attacks on zero-shot LLM assessment. In *Proceedings of the 2024 Conference on Empirical Methods in Natural Language Processing*, pp. 7499–7517, 2024.
- Sackmann, J. ATP tennis rankings, results, and stats. https://github.com/JeffSackmann/tennis_atp, 2024. GitHub repository.
- Shah, N. B. and Wainwright, M. J. Simple, robust and optimal ranking from pairwise comparisons. *Journal of Machine Learning Research*, 18(199):1–38, 2018.
- Sherman, J. and Morrison, W. J. Adjustment of an inverse matrix corresponding to a change in one element of a given matrix. *The Annals of Mathematical Statistics*, 21(1):124–127, 1950. doi: 10.1214/aoms/1177729893.
- Singh, S., Nan, Y., Wang, A., D’Souza, D., Kapoor, S., Üstün, A., Koyejo, S., Deng, Y., Longpre, S., Smith, N. A., Ermis, B., Fadaee, M., and Hooker, S. The leaderboard illusion. In *The Thirty-ninth Annual Conference on Neural Information Processing Systems, Datasets and Benchmarks Track*, 2025.
- Sun, H., Shen, Y., and Ton, J.-F. Rethinking reward modeling in preference-based large language model alignment. In *The Thirteenth International Conference on Learning Representations*, 2025.
- Suri, A., Chaudhari, H., Peng, Y., Naseh, A., Oprea, A., and Houmansadr, A. Exploiting leaderboards for large-scale distribution of malicious models. In *IEEE Symposium on Security and Privacy*, 2026.
- Touvron, H., Martin, L., Stone, K., Albert, P., Almahairi, A., Babaei, Y., Bashlykov, N., Batra, S., Bhargava, P., Bhosale, S., et al. Llama 2: Open foundation and fine-tuned chat models. *arXiv preprint arXiv:2307.09288*, 2023.
- Vichare, A., Angelopoulos, A. N., Chiang, W.-L., Tang, K., and Manolache, L. WebDev Arena: A live LLM leaderboard for web app development. LMarena Blog, <https://arena.ai/blog/webdev-arena/>, 2025.
- Wachter, S., Mittelstadt, B., and Russell, C. Counterfactual explanations without opening the black box: Automated decisions and the GDPR. *Harvard Journal of Law & Technology*, 31(2):841–887, 2018.
- Wang, Z., Han, Y., Fang, E. X., Wang, L., and Lu, J. Confidence diagram of nonparametric ranking for uncertainty assessment in large language models evaluation. *arXiv preprint arXiv:2412.05506*, 2025.

- Wauthier, F., Jordan, M., and Jojic, N. Efficient ranking from pairwise comparisons. In *Proceedings of the 30th International Conference on Machine Learning*, volume 28, pp. 109–117. PMLR, 2013.
- Wu, W., Niezink, N., and Junker, B. A diagnostic framework for the Bradley–Terry model. *Journal of the Royal Statistical Society, Series A: Statistics in Society*, 185 (Supplement 2):S461–S484, 2022. doi: 10.1111/rssa.12959.
- Xu, S., Yue, B., Zha, H., and Liu, G. Uncertainty-aware preference alignment in reinforcement learning from human feedback. In *ICML 2024 Workshop on Models of Human Feedback for AI Alignment*, 2024.
- Zhao, W., Rush, A. M., and Goyal, T. Challenges in trustworthy human evaluation of chatbots. In *Findings of the Association for Computational Linguistics: NAACL 2025*, pp. 3359–3365, Albuquerque, New Mexico, 2025. Association for Computational Linguistics. doi: 10.18653/v1/2025.findings-naacl.186.
- Zheng, L., Chiang, W.-L., Sheng, Y., Zhuang, S., Wu, Z., Zhuang, Y., Lin, Z., Li, Z., Li, D., Xing, E. P., Zhang, H., Gonzalez, J. E., and Stoica, I. Judging LLM-as-a-judge with MT-Bench and chatbot arena. In *Advances in Neural Information Processing Systems*, volume 36, pp. 46595–46623, 2023.
- Zheng, X., Pang, T., Du, C., Liu, Q., Jiang, J., and Lin, M. Cheating automatic LLM benchmarks: Null models achieve high win rates. In *The Thirteenth International Conference on Learning Representations*, 2025.

A. Datasets

A.1. Dataset descriptions

We evaluate our framework on seven pairwise comparison datasets spanning LLM evaluation and sports ranking.

Chatbot Arena 55k. A crowdsourced platform where users simultaneously interact with two anonymous chatbots and vote for the preferred response (Chiang et al., 2024b). We use the `arena-human-preference-55k` split, containing 57,477 human preference judgements across 64 models. Its large scale and diverse user base make it the primary benchmark for LLM evaluation, but also expose it to noise and adversarial risk.

Chatbot Arena LLM Judges. A companion dataset from the same platform in which pairwise preferences are collected using an LLM-as-a-judge rather than human votes (Chiang et al., 2024b). The `chatbot-arena-llm-judges` split contains 49,938 comparisons across 64 models, allowing us to contrast automated and human evaluation robustness.

MT-Bench Human Judgments. A curated multi-turn benchmark designed to evaluate instruction-following and reasoning (Chiang et al., 2024b). Preferences were collected from 58 expert-level annotators (predominantly graduate students), yielding 3,355 high-quality pairwise judgements. Its smaller size and expert annotations make it substantially more robust than crowdsourced platforms.

Vision Arena. A crowdsourced arena for vision-language models in which users compare two anonymous models on visual question-answering tasks. We use the `lmarena-ai/VisionArena-Battle` dataset, which contains 29,849 single- and multi-turn conversations.

WebDev Arena. A crowdsourced arena focused on web-development tasks such as building interactive applications and webpages. We use the `lmarena-ai/webdev-arena-preference-10k` dataset (10,501 prompts), which provides domain-specific evaluation complementary to open-ended chat.

ATP Top-10 Tennis. Match records from the ATP tour (2020–2024). We restrict to the top-10 ranked players by 2024 season standing who each played at least 20 matches, yielding 278 games in total. This small, sparse graph tests our framework in a regime where each individual match carries high weight.

NBA Elo Top-50. Historical NBA game records from all seasons. We focus on the top-50 teams by total games played, yielding 109,892 matchups. Its large, dense comparison graph represents the opposite extreme from ATP and assesses robustness in high-data settings.

A.2. Dataset processing

Each pairwise comparison in the raw datasets is represented as a directed observation (i, j, y) where $y \in \{0, 1\}$ indicates whether item i beat item j . To handle ties uniformly, we adopt a symmetric formulation: each tied comparison is decomposed into two directed observations— $(i, j, 1)$ and $(j, i, 1)$ —treating both items as winning once against the other.

For consistency, each non-tied comparison is also represented as two directed observations: (i, j, y) and $(j, i, 1 - y)$. This ensures a uniform representation across all comparisons and simplifies the construction of the feature matrix $X \in \mathbb{R}^{N \times M}$ used in the BT likelihood, where each row $x_n = e_i - e_j$ encodes the matched pair. The resulting matrix X and outcome vector y are the direct inputs to the weighted M-estimator in Eq. 3.

B. Framework details

B.1. BT input–output format and tie handling

The Bradley–Terry model takes a directed comparison dataset as input and outputs a latent skill vector for the ranked players. Let M denote the number of players and let $e_i \in \mathbb{R}^M$ be the standard basis vector for player i . Each directed comparison between players i and j is encoded as

$$x_{ij} = e_i - e_j,$$

where the binary outcome $y_{ij} \in \{0, 1\}$ indicates whether the first player in the directed pair wins. After preprocessing, the full BT input is represented as

$$X = \begin{bmatrix} x_1^\top \\ x_2^\top \\ \vdots \\ x_N^\top \end{bmatrix} \in \{-1, 0, 1\}^{N \times M}, \quad y = \begin{bmatrix} y_1 \\ y_2 \\ \vdots \\ y_N \end{bmatrix} \in \{0, 1\}^N,$$

where each row $x_n = e_{i_n} - e_{j_n}$ contains one +1 entry for the first player, one -1 entry for the second player, and zeros elsewhere.

Given this input, the BT model estimates a skill vector

$$\hat{\theta} = \begin{bmatrix} \hat{\theta}_1 \\ \hat{\theta}_2 \\ \vdots \\ \hat{\theta}_M \end{bmatrix} \in \mathbb{R}^M,$$

where larger values indicate stronger players. The main BT outputs are the fitted skill vector $\hat{\theta}$, the induced pairwise win probabilities, and the resulting leaderboard obtained by sorting players according to $\hat{\theta}$.

To keep the representation symmetric, every raw comparison is converted into a two-row directed block. A decisive comparison in which player i beats player j is represented as

$$i \succ j \quad \mapsto \quad B_b = \{(x_{ij}, 1), (x_{ji}, 0)\}.$$

A tied comparison is represented as

$$i \sim j \quad \mapsto \quad B_b = \{(x_{ij}, 1), (x_{ji}, 1)\}.$$

Thus, decisive outcomes provide one win and one loss across the two directions, whereas ties assign equal win evidence to both players. This avoids choosing an arbitrary winner for tied comparisons while keeping all raw comparisons in the same binary directed BT input format.

For influence scoring and perturbation selection, we treat each two-row block B_b as the atomic unit. If the two directed rows corresponding to raw comparison b are indexed by n_{ij} and n_{ji} , its block-level influence is computed by summing the row-level influences,

$$I_{B_b}^{(f)} = I_{n_{ij}}^{(f)} + I_{n_{ji}}^{(f)}.$$

All perturbation actions are then applied at the block level. A `DROP` action removes both directed rows in B_b ; an `ADD` action inserts both directed rows for the candidate comparison; and a `FLIP` action modifies the two directed outcomes consistently. This ensures that influence scores and perturbation budgets are defined at the level of the original raw comparisons, rather than at the level of individual directed rows.

B.2. One-step Newton refinement

The influence approximation in Eq. 5 linearizes the effect of an infinitesimal change in a match weight. However, the actions used in our experiments are finite. Following the classical one-step case-deletion approximation for logistic regression diagnostics (Pregibon, 1981) and its use in influence-based leaderboard analysis (Huang et al., 2026), we use a one-step Newton (1sN) correction to better approximate these finite actions without fully refitting the BT model for every candidate.

For a comparison $z_n = (x_n, y_n)$ evaluated at the full-data fit $\hat{\theta}$, let

$$p_n = \sigma(x_n^\top \hat{\theta}), \quad r_n = y_n - p_n, \quad v_n = p_n(1 - p_n).$$

Since

$$g_n = \nabla_{\theta} \ell(z_n; \hat{\theta}) = (p_n - y_n)x_n = -r_n x_n,$$

the first-order deletion approximation from Eq. 4 gives

$$\Delta \hat{\theta}_{n, \text{IF}}^{\text{DROP}} = -r_n H^{-1} x_n.$$

For adding a candidate comparison $z_c = (x_c, y_c)$, the sign is reversed:

$$\Delta \hat{\theta}_{c, \text{IF}}^{\text{Add}} = r_c H^{-1} x_c.$$

The influence approximation uses the original Hessian H and therefore ignores the fact that a finite `Drop` or `Add` action also changes the local curvature of the BT objective. For a single comparison, this curvature change is rank-one because the Hessian contribution of z_n is

$$H_n = v_n x_n x_n^\top.$$

Thus, dropping one comparison changes the local Hessian from H to $H - H_n$, while adding one comparison changes it from H to $H + H_n$. Applying the Sherman–Morrison (Sherman & Morrison, 1950) identity to this rank-one Hessian update yields the leverage term

$$h_n = v_n x_n^\top H^{-1} x_n.$$

This quantity measures how strongly comparison n affects the local curvature around the fitted BT solution.

The 1sN refinement rescales the influence update by this leverage term. For a single dropped comparison,

$$\Delta \hat{\theta}_{n, \text{1sN}}^{\text{Drop}} = \frac{\Delta \hat{\theta}_{n, \text{IF}}^{\text{Drop}}}{1 - h_n},$$

whereas for a single added candidate comparison,

$$\Delta \hat{\theta}_{c, \text{1sN}}^{\text{Add}} = \frac{\Delta \hat{\theta}_{c, \text{IF}}^{\text{Add}}}{1 + h_c}.$$

Equivalently, 1sN performs one Newton step toward the optimum of the perturbed objective, starting from the original fitted parameters $\hat{\theta}$. Compared with influence, which keeps the curvature fixed at H , 1sN partially accounts for the curvature change induced by the finite action while still avoiding a full BT refit.

Given either the IF or 1sN parameter change, the predicted objective change is computed using the same projection rule as Eq. 5:

$$\Delta f \approx \nabla_{\theta} f(\hat{\theta}, w)^\top \Delta \hat{\theta} + \left(\frac{\partial f}{\partial w} \right)_{\text{explicit}}.$$

The explicit term is zero for objectives that depend on the data only through $\hat{\theta}$, and nonzero for uncertainty objectives that directly depend on the weighted comparison graph.

For Flip actions, the local Hessian is unchanged because it depends only on the comparison features x_n and not on the outcome y_n . As a result, the 1sN correction does not apply, and the flip update reduces to its first-order influence approximation.

B.3. Grouped Newton refinement for player removal

The 1sN correction in Appendix B.4 is applied to individual `Drop` or `Add` actions. Player removal is a larger structured perturbation: for a player m , it removes the full set of incident comparisons

$$\mathcal{G}_m = \{n : z_n \text{ contains player } m\}.$$

Because many comparisons are removed simultaneously, we use a grouped Newton refinement rather than applying the single-row leverage correction independently to each match.

Let

$$p_n = \sigma(x_n^\top \hat{\theta}), \quad r_n = y_n - p_n, \quad s_n = r_n x_n$$

denote the fitted probability, residual, and score contribution of comparison n at the full-data solution. For a group deletion \mathcal{G} , the first step sums the removed score contributions and applies the grouped first-order deletion update:

$$s_{\mathcal{G}} = \sum_{n \in \mathcal{G}} s_n, \quad \theta_{\mathcal{G}}^{(1)} = \hat{\theta} - H^{-1} s_{\mathcal{G}}.$$

This is the group analogue of summing individual deletion influences in Eq. 12.

The second step takes one Newton correction using only the comparisons that remain after the group is removed. Let

$$\mathcal{K} = \{1, \dots, N\} \setminus \mathcal{G}$$

be the kept set. At $\theta_{\mathcal{G}}^{(1)}$, define the kept-data score

$$S_{\mathcal{K}}(\theta_{\mathcal{G}}^{(1)}) = \sum_{n \in \mathcal{K}} x_n (y_n - \sigma(x_n^\top \theta_{\mathcal{G}}^{(1)})),$$

and the kept-data Hessian

$$H_{\mathcal{K}}(\theta_{\mathcal{G}}^{(1)}) = \sum_{n \in \mathcal{K}} v_n^{(1)} x_n x_n^\top + \lambda I, \quad v_n^{(1)} = p_n^{(1)}(1 - p_n^{(1)}),$$

where

$$p_n^{(1)} = \sigma(x_n^\top \theta_{\mathcal{G}}^{(1)}).$$

The grouped Newton estimate is then

$$\theta_{\mathcal{G}}^{(2)} = \theta_{\mathcal{G}}^{(1)} + H_{\mathcal{K}}(\theta_{\mathcal{G}}^{(1)})^{-1} S_{\mathcal{K}}(\theta_{\mathcal{G}}^{(1)}).$$

Finally, the predicted player-removal effect is evaluated by applying the downstream objective to this approximate parameter vector:

$$\Delta f_{\mathcal{G}} \approx f(\theta_{\mathcal{G}}^{(2)}) - f(\hat{\theta}).$$

For the player-removal experiments, f is the smooth *Kendall's* τ objective computed over the remaining players, so the removed player is excluded from the reference ranking before evaluating the objective.

This refinement is related to 1sN because both start from an influence-based finite-deletion step and then use local curvature information to improve the approximation. The difference is that 1sN uses a scalar leverage correction for a single comparison, whereas the grouped Newton refinement performs one joint Newton correction on the score equation of the remaining comparison data.

B.4. Objective influence derivations

The objective-level influence of a match z_n is given by Eq. 5 and is defined as the derivative of the objective with respect to increasing the match weight w_n . Thus, for a finite `DROP` action, where $\Delta w_n = -1$, the predicted drop effect is

$$\Delta f_{\text{drop},n} \approx -\mathcal{I}_n^{(f)}.$$

Equivalently, specialising to first-order case-deletion, the parameter shift induced by dropping match n is

$$\Delta \theta_n^{\text{drop}} \approx -r_n H^{-1} x_n, \quad r_n = y_n - p_n, \quad p_n = \sigma(x_n^\top \hat{\theta}),$$

and the corresponding drop effect is

$$\Delta f_{\text{drop},n}^{(f)} \approx \nabla_{\theta} f(\hat{\theta}, w)^\top \Delta \theta_n^{\text{drop}} - \left(\frac{\partial f}{\partial w_n} \right)_{\text{explicit}}.$$

Therefore,

$$\mathcal{I}_n^{(f)} \approx -\Delta f_{\text{drop},n}^{(f)}.$$

For objectives that depend on the data only through $\hat{\theta}$, the explicit term is zero.

Top- k membership objective / gap objective. For

$$f_{\text{gap}}(\theta) = \theta_i - \theta_j,$$

the full gradient is

$$\nabla_{\theta} f_{\text{gap}} = e_i - e_j,$$

and there is no explicit weight term, so the drop effect is

$$\Delta f_{\text{drop},n}^{(f_{\text{gap}})} = \nabla_{\theta} f_{\text{gap}}(\hat{\theta})^{\top} \Delta \theta_n^{\text{drop}}.$$

Equivalently,

$$\mathcal{I}_n^{(f_{\text{gap}})} = -\nabla_{\theta} f_{\text{gap}}(\hat{\theta})^{\top} \Delta \theta_n^{\text{drop}}.$$

Player-uncertainty objective. The single-player uncertainty proxy is

$$f_{\text{CI-player}}(m; w) = \frac{1}{\hat{\rho}_m^2(w)}, \quad \hat{\rho}_m^2(w) = \sum_{j \neq m} w_{mj} v_{mj},$$

with

$$v_{mj} = \sigma(\hat{\theta}_m - \hat{\theta}_j)(1 - \sigma(\hat{\theta}_m - \hat{\theta}_j)).$$

Define

$$v'_{mj} = v_{mj}(1 - 2\sigma(\hat{\theta}_m - \hat{\theta}_j)).$$

Then

$$\frac{\partial \hat{\rho}_m^2}{\partial \theta_k} = \sum_{j \neq m} w_{mj} v'_{mj} (\mathbf{1}_{k=m} - \mathbf{1}_{k=j}),$$

so

$$\nabla_{\theta} f_{\text{CI-player}}(m; w) = -\frac{1}{(\hat{\rho}_m^2)^2} \nabla_{\theta} \hat{\rho}_m^2.$$

Unlike the gap and *Kendall's* τ objectives, this objective has a nonzero explicit weight derivative. If match n compares players a_n and b_n , then

$$\left(\frac{\partial f_{\text{CI-player}}(m; w)}{\partial w_n} \right)_{\text{explicit}} = -\frac{v_n \mathbf{1}\{m \in \{a_n, b_n\}\}}{(\hat{\rho}_m^2)^2},$$

where

$$v_n = \sigma(\hat{\theta}_{a_n} - \hat{\theta}_{b_n})(1 - \sigma(\hat{\theta}_{a_n} - \hat{\theta}_{b_n})).$$

Therefore, the drop effect is

$$\Delta f_{\text{drop},n}^{(f_{\text{CI-player}})} = \nabla_{\theta} f_{\text{CI-player}}(m; w)^{\top} \Delta \theta_n^{\text{drop}} + \frac{v_n \mathbf{1}\{m \in \{a_n, b_n\}\}}{(\hat{\rho}_m^2)^2}.$$

Equivalently,

$$\mathcal{I}_n^{(f_{\text{CI-player}})} = -\nabla_{\theta} f_{\text{CI-player}}(m; w)^{\top} \Delta \theta_n^{\text{drop}} - \frac{v_n \mathbf{1}\{m \in \{a_n, b_n\}\}}{(\hat{\rho}_m^2)^2}.$$

Trace-uncertainty objective. The global uncertainty proxy is

$$f_{\text{CI-trace}}(w) = \sum_{i=1}^M \frac{1}{\hat{\rho}_i^2(w)}.$$

Hence

$$\nabla_{\theta} f_{\text{CI-trace}} = \sum_{i=1}^M -\frac{1}{(\hat{\rho}_i^2)^2} \nabla_{\theta} \hat{\rho}_i^2.$$

For match n between players a_n and b_n , the explicit weight derivative is

$$\left(\frac{\partial f_{\text{CI-trace}}}{\partial w_n}\right)_{\text{explicit}} = -v_n \left(\frac{1}{(\hat{\rho}_{a_n}^2)^2} + \frac{1}{(\hat{\rho}_{b_n}^2)^2} \right).$$

Therefore, the drop effect is

$$\Delta f_{\text{drop},n}^{(f_{\text{CI-trace}})} = \nabla_{\theta} f_{\text{CI-trace}}(\hat{\theta}, w)^\top \Delta \theta_n^{\text{drop}} + v_n \left(\frac{1}{(\hat{\rho}_{a_n}^2)^2} + \frac{1}{(\hat{\rho}_{b_n}^2)^2} \right).$$

Equivalently,

$$\mathcal{I}_n^{(f_{\text{CI-trace}})} = -\nabla_{\theta} f_{\text{CI-trace}}(\hat{\theta}, w)^\top \Delta \theta_n^{\text{drop}} - v_n \left(\frac{1}{(\hat{\rho}_{a_n}^2)^2} + \frac{1}{(\hat{\rho}_{b_n}^2)^2} \right).$$

Kendall's τ surrogate.

$$f_{\tau,T}(\theta; \pi) = \frac{2}{M(M-1)} \sum_{a < b} s_{ab} \tanh\left(\frac{\theta_a - \theta_b}{T}\right),$$

where $s_{ab} \in \{-1, +1\}$ is induced by the reference ranking. Its gradient is

$$[\nabla_{\theta} f_{\tau,T}]_k = \frac{2}{M(M-1)} \sum_{a < b} s_{ab} \frac{1 - \tanh^2\left(\frac{\theta_a - \theta_b}{T}\right)}{T} (\mathbf{1}_{k=a} - \mathbf{1}_{k=b}).$$

There is no explicit weight term, so the drop effect is

$$\Delta f_{\text{drop},n}^{(f_{\tau,T})} = \nabla_{\theta} f_{\tau,T}(\hat{\theta})^\top \Delta \theta_n^{\text{drop}}.$$

Equivalently,

$$\mathcal{I}_n^{(f_{\tau,T})} = -\nabla_{\theta} f_{\tau,T}(\hat{\theta})^\top \Delta \theta_n^{\text{drop}}.$$

B.5. Algorithms

Algorithm 1 Top- k Robustness Action Search

Input: Fitted BT model \hat{M} on dataset D ; top- k cross-boundary candidate pairs \mathcal{P} ; fixed action variant $a \in \{\text{drop}, \text{flip}, \text{add_pairs}, \text{add_outcomes}, \text{add_weighted}\}$; maximum action count A

Output: Smallest action count that makes at least one candidate top- k boundary gap cross zero, together with the corresponding pair and selected actions

- 1: Fit the Bradley–Terry model on D and obtain $\hat{\theta}$.
 - 2: Construct the ordered set \mathcal{P} of top- k cross-boundary candidate pairs.
 - 3: Initialize an empty cache of per-pair influence reports.
 - 4: **for** $\alpha = 1, \dots, A$ **do**
 - 5: **for all** $(i, j) \in \mathcal{P}$ **do**
 - 6: Define the gap objective

$$f_{ij}(\hat{\theta}) = \hat{\theta}_i - \hat{\theta}_j.$$
 - 7: If not already cached, compute the one-step influence report for (i, j) under action variant a over its corresponding candidate pool.
 - 8: Let $g_{ij} = f_{ij}(\hat{\theta})$.
 - 9: Select the top α candidate actions from the cached report: if $g_{ij} > 0$, choose the most negative influences; otherwise choose the most positive influences.
 - 10: Treat grouped forward/reverse copies as one logical action when applicable.
 - 11: Apply the selected α actions to the original fitted dataset D , refit the BT model, and compute the updated gap $g_{ij}^{(\alpha)}$ or estimate it using influence scores.
 - 12: **if** $(g_{ij} > 0 \text{ and } g_{ij}^{(\alpha)} \leq 0)$ **or** $(g_{ij} < 0 \text{ and } g_{ij}^{(\alpha)} \geq 0)$ **then**
 - 13: **return** α , the pair (i, j) , and the selected actions.
 - 14: **end if**
 - 15: **end for**
 - 16: **end for**
 - 17: **return** failure if no pair succeeds within A actions.
-

Algorithm 2 Strict CI-Aware Top- k Manipulation

Input: Fitted BT model \hat{M} on dataset D ; target boundary rank k ; fixed insider–outsider pair (i, j) with i the rank- k model and j the rank- $(k + 1)$ model; CI method; CI level α_{CI} ; fixed action variant $a \in \{\text{drop, flip, add_pairs, add_outcomes, add_weighted}\}$; maximum action budget A

Output: Minimum number of actions needed to certify that outsider j lies above insider i under strict CI separation

1: Fit the Bradley–Terry model on D and compute $\hat{\theta}$ and confidence intervals

$$[L_m, U_m] = \left[\hat{\theta}_m - z_{\alpha_{\text{CI}}} \text{SE}_m, \hat{\theta}_m + z_{\alpha_{\text{CI}}} \text{SE}_m \right], \quad z_{\alpha_{\text{CI}}} = \Phi^{-1}(1 - \alpha_{\text{CI}}/2).$$

2: Fix the target pair (i, j) before the action search begins, where i is the current insider and j is the current outsider.

3: Define the strict CI objective

$$g_{ij}(\hat{\theta}) = (\hat{\theta}_i + z_{\alpha_{\text{CI}}} \text{SE}_i) - (\hat{\theta}_j - z_{\alpha_{\text{CI}}} \text{SE}_j).$$

4: Proceed only if the strict target is initially unmet: $g_{ij}(\hat{\theta}) \geq 0$.

5: Compute the full one-step influence report for action variant a with respect to g_{ij} over its corresponding candidate pool.

6: **for** $\ell = 1, \dots, A$ **do**

7: Select the top ℓ logical actions that most decrease g_{ij} .

8: Rank candidates in ascending order of influence and group forward/reverse copies sharing a match identifier as one logical action.

9: Form the first-order screened objective

$$\tilde{g}_{ij}^{(\ell)} = g_{ij}(\hat{\theta}) + \sum_{n \in S_{ij}^{(\ell)}} \mathcal{I}_n^{(g_{ij})}.$$

10: **if** $\tilde{g}_{ij}^{(\ell)} \geq 0$ **then**

11: Continue to the next budget without refitting.

12: **end if**

13: Apply the selected actions to the original fitted dataset and refit the BT model.

14: Recompute the exact strict objective on the refit model:

$$g_{ij}^{\text{refit}} = (\hat{\theta}_i^{\text{refit}} + z_{\alpha_{\text{CI}}} \text{SE}_i^{\text{refit}}) - (\hat{\theta}_j^{\text{refit}} - z_{\alpha_{\text{CI}}} \text{SE}_j^{\text{refit}}).$$

15: Let

$$L_j^{\text{refit}} = \hat{\theta}_j^{\text{refit}} - z_{\alpha_{\text{CI}}} \text{SE}_j^{\text{refit}}, \quad U_i^{\text{refit}} = \hat{\theta}_i^{\text{refit}} + z_{\alpha_{\text{CI}}} \text{SE}_i^{\text{refit}}.$$

16: **if** $g_{ij}^{\text{refit}} < 0$, equivalently $L_j^{\text{refit}} > U_i^{\text{refit}}$ **then**

17: **return** ℓ , the pair (i, j) , the selected actions, and the refit ranking.

18: **end if**

19: **end for**

20: **return** failure if no strict CI-certified reversal is found within budget A .

Algorithm 3 Strict CI-Aware k -Selection

Input: Fitted BT model \hat{M} on dataset D ; CI method; CI level α_{CI}

Output: One valid target triple (k, i, j) for strict CI-aware manipulation, where i is the insider at rank k and j is the outsider at rank $k + 1$

1: Fit the Bradley–Terry model on D and compute the CI ranking.

2: **for** $k = 1, \dots, N - 1$ **do**

3: Let

$$i = \text{rank-}k \text{ model}, \quad j = \text{rank-}(k + 1) \text{ model}.$$

4: Define the strict CI objective

$$g_{ij}(\hat{\theta}) = (\hat{\theta}_i + z_{\alpha_{\text{CI}}} \text{SE}_i) - (\hat{\theta}_j - z_{\alpha_{\text{CI}}} \text{SE}_j), \quad z_{\alpha_{\text{CI}}} = \Phi^{-1}(1 - \alpha_{\text{CI}}/2).$$

5: Keep boundary (k, i, j) only if $g_{ij}(\hat{\theta}) \geq 0$.

6: **end for**

7: **return** the valid boundary (k, i, j) with the smallest strict objective value $g_{ij}(\hat{\theta})$.

8: **return** failure if no valid boundary satisfies $g_{ij}(\hat{\theta}) \geq 0$.

Algorithm 4 Online Rigging Baseline (Omni-On)

Input: Initial Elo/BT ranking on dataset D ; target model m^* ; direction $d \in \{\text{promote}, \text{demote}\}$; exposed pair stream $\{(a_t, b_t)\}_{t=1}^B$; budget B ; Elo constants $K, \text{BASE}, \text{SCALE}$

Output: Ordered decisions and target-rank trajectory

1: Compute the initial ratings $\hat{\theta}_m$ for all models and record the initial rank of m^* .

2: **for** $t = 1, \dots, B$ **do**

3: Observe the currently exposed pair (a_t, b_t) .

4: Define the candidate decision set

$$\mathcal{A}_t = \{\text{a_wins}, \text{b_wins}, \text{tie}, \text{remove}\}.$$

5: Let $r_a = \hat{\theta}_{a_t}$, $r_b = \hat{\theta}_{b_t}$, and $r_* = \hat{\theta}_{m^*}$, and compute

$$e_a = \frac{1}{1 + \text{BASE}^{(r_b - r_a)/\text{SCALE}}}, \quad e_b = \frac{1}{1 + \text{BASE}^{(r_a - r_b)/\text{SCALE}}}.$$

6: **for all** $\alpha \in \mathcal{A}_t$ **do**

7: Form one-step hypothetical ratings:

$$\text{a_wins: } r'_a = r_a + Ke_b, \quad r'_b = r_b - Ke_b,$$

$$\text{b_wins: } r'_a = r_a - Ke_a, \quad r'_b = r_b + Ke_a,$$

$$\text{tie: } r'_a = r_a - \frac{K}{2}(e_a - e_b), \quad r'_b = r_b + \frac{K}{2}(e_a - e_b),$$

$$\text{remove: } r'_a = r_a, \quad r'_b = r_b.$$

8: Score the action by

$$r_t^{(\alpha)} = \frac{1}{1 + \text{BASE}^{(r'_a - r_*)/\text{SCALE}}} + \frac{1}{1 + \text{BASE}^{(r'_b - r_*)/\text{SCALE}}}.$$

9: For demotion tasks, use the inverted reward $-r_t^{(\alpha)}$.

10: **end for**

11: Select the greedy rigging decision

$$\alpha_t^* = \arg \max_{\alpha \in \mathcal{A}_t} r_t^{(\alpha)}.$$

12: Apply α_t^* to the current dataset, refit the global Elo/BT ranking, and append the new rank of m^* to the history.

13: **end for**

14: **return** the decision history and the full target-rank trajectory.

Algorithm 5 Online Influence-Guided Targeted Top- k Manipulation

Input: Fitted BT model on dataset D ; target model m^* ; rank cutoff k ; direction $d \in \{\text{promote}, \text{demote}\}$; exposed pair stream $\{(a_t, b_t)\}_{t=1}^B$; budget B

Output: Selected decision sequence and final ranking

- 1: Fit the initial BT model on D and record the initial ranking of m^* .
- 2: **for** $t = 1, \dots, B$ **do**
- 3: **if** $\text{rank}(m^*) \leq k$ for promotion, or $\text{rank}(m^*) > k$ for demotion **then**
- 4: Stop; the target condition is already satisfied.
- 5: **end if**
- 6: Observe the currently exposed pair (a_t, b_t) .
- 7: Define the current boundary objective

$$f_t(\hat{\theta}) = \begin{cases} \hat{\theta}_{m^*} - \hat{\theta}_{m_k}, & d = \text{promote}, \\ \hat{\theta}_{m_{k+1}} - \hat{\theta}_{m^*}, & d = \text{demote}, \end{cases}$$

where m_k and m_{k+1} are recomputed from the current ranking after every accepted intervention.

- 8: Form the allowable decision set

$$\mathcal{A}_t = \{\text{remove}, \text{a_wins}, \text{b_wins}, \text{tie}\}.$$

- 9: **for all** $\alpha \in \mathcal{A}_t$ **do**
- 10: Compute its first-order influence score on the current boundary objective:

$$\mathcal{I}^{(f_t)}(\alpha) = \begin{cases} 0, & \alpha = \text{remove}, \\ \mathcal{I}^{(f_t)}(a_t \succ b_t), & \alpha = \text{a_wins}, \\ \mathcal{I}^{(f_t)}(b_t \succ a_t), & \alpha = \text{b_wins}, \\ \mathcal{I}^{(f_t)}(a_t \succ b_t) + \mathcal{I}^{(f_t)}(b_t \succ a_t), & \alpha = \text{tie}. \end{cases}$$

- 11: **end for**
- 12: Select the greedy influence decision

$$\alpha_t^* = \arg \max_{\alpha \in \mathcal{A}_t} \mathcal{I}^{(f_t)}(\alpha).$$

- 13: Apply α_t^* to the current state, refit the BT model if any row is added or estimate it using influence scores, and update the ranking of m^* .
 - 14: **end for**
 - 15: **return** the selected decisions, success indicator, rank trajectory, and final ranking.
-

Algorithm 6 CI Reduction via Targeted Match Addition

Input: Fitted dataset D ; target model m^* ; budget B ; candidate mode v ; CI method

Output: Added matches and CI-width trajectory

- 1: Fit the BT model on D .
- 2: Define the target uncertainty objective

$$f_{\text{CI-player}}(m^*; w) = \hat{\rho}_{m^*}^{-2}(w),$$

where $\hat{\rho}_{m^*}^2(w)$ is the BT pair-weight-aggregated variance proxy defined in Eq. 3.2.

- 3: Generate add candidates according to v .
 - 4: **if** $v = \text{all_pairs}$ **then**
 - 5: Use one unordered pair, with the currently higher-rated model as winner.
 - 6: **end if**
 - 7: Compute influence scores for all add candidates once using the initial fitted model.
 - 8: Sort candidates in ascending influence order, since negative influence reduces $f_{\text{CI-player}}$.
 - 9: **for** $t = 1, \dots, B$ **do**
 - 10: Select the next unused candidate from the fixed sorted list.
 - 11: Add the match to D , refit the BT model or estimate it using influence scores, and recompute the target CI width using the chosen CI method.
 - 12: **end for**
 - 13: **return** the added matches, uncertainty trajectory, and CI-width trajectory.
-

B.6. Relation between the sandwich covariance and Gao-style local information

In the main text, we use the Gao-style local information approximation as a scalable uncertainty proxy for BT leaderboard perturbations. Here, we clarify how this approximation relates to the more general sandwich covariance estimator and why we avoid propagating perturbations through the full covariance matrix.

For a comparison $z_n = (i_n, j_n, y_n)$, define $x_n = e_{i_n} - e_{j_n}$, $p_n = \sigma(x_n^\top \hat{\theta})$, and $v_n = p_n(1 - p_n)$. The weighted BT loss has score and Hessian

$$g_n(\hat{\theta}) = (p_n - y_n)x_n, \quad H_n(\hat{\theta}) = v_n x_n x_n^\top. \quad (13)$$

The sandwich covariance estimator is

$$\Sigma_{\text{sand}}(w) = J(w)^{-1} S(w) J(w)^{-1}, \quad J(w) = \sum_{n=1}^N w_n H_n(\hat{\theta}) + \lambda I, \quad S(w) = \sum_{n=1}^N w_n g_n(\hat{\theta}) g_n(\hat{\theta})^\top. \quad (14)$$

Under a correctly specified BT model,

$$\mathbb{E}[g_n(\theta^*) g_n(\theta^*)^\top \mid i_n, j_n] = v_n x_n x_n^\top = H_n(\theta^*),$$

because $\text{Var}(y_n \mid i_n, j_n) = p_n(1 - p_n) = v_n$. Thus, replacing the empirical score covariance $S(w)$ by its information counterpart gives $S(w) \approx J(w)$, so the sandwich estimator reduces to the information-based covariance approximation

$$\Sigma_{\text{sand}}(w) = J(w)^{-1} S(w) J(w)^{-1} \approx J(w)^{-1}. \quad (15)$$

Gao et al. (Gao et al., 2023) provide a coordinate-wise interpretation of uncertainty in the BT model: the leading uncertainty of each skill estimate is controlled by the inverse of a local Fisher information term. In our weighted finite-sample setting, this local information for player i is

$$\hat{\rho}_i^2(w) = \sum_{j \neq i} w_{ij} p_{ij} (1 - p_{ij}), \quad (16)$$

Equivalently, if

$$J_{\text{info}}(w) = \sum_{i < j} w_{ij} p_{ij} (1 - p_{ij}) (e_i - e_j)(e_i - e_j)^\top + \lambda I,$$

then $\hat{\rho}_i^2(w)$ is the local diagonal information associated with player i .

The full inverse $J_{\text{info}}(w)^{-1}$ captures coupling across all players, but computing and differentiating its diagonal or trace under many candidate perturbations is substantially more expensive. We therefore use the Gao-style local approximation

$$\text{Var}(\hat{\theta}_i) \approx \hat{\rho}_i^{-2}(w),$$

which leads directly to the player-wise uncertainty objective $f_{\text{CI-player}}(i; w) = \hat{\rho}_i^{-2}(w)$ and the trace-style global uncertainty proxy $f_{\text{CI-trace}}(w) = \sum_i \hat{\rho}_i^{-2}(w)$. This approximation preserves the main statistical intuition: uncertainty is reduced by adding informative comparisons incident to poorly measured players, and increased by removing such comparisons.

B.7. Arena Active baseline for target-model CI reduction

Arena Active (Chiang et al., 2024b) is implemented as a target-restricted variance-count active-sampling baseline. At each step, it recomputes scores under the current Bradley–Terry fit, considers only candidate pairs involving the target model m^* , and selects the opponent j with the largest estimated one-step reduction

$$s_{\text{AA}}(m^*, j) = \sqrt{\frac{\widehat{\text{Var}}(\hat{\theta}_{m^*} - \hat{\theta}_j)}{N_{m^*j}}} - \sqrt{\frac{\widehat{\text{Var}}(\hat{\theta}_{m^*} - \hat{\theta}_j)}{N_{m^*j} + 1}}, \quad (17)$$

where N_{m^*j} is the current number of comparisons between m^* and j . Thus, Arena Active is target-specific through candidate filtering, but uses a pair-level variance-count score rather than directly optimizing $f_{\text{CI-player}}(m^*)$. In contrast, our method scores each candidate `Add` action by its predicted effect on the target uncertainty objective.

Arena Active always uses an `all_pairs` candidate space, while the influence policy can use `all_pairs`, `all_outcomes`, or `all_outcomes_weighted`. Thus, `all_pairs` is the fairest direct comparison, whereas the outcome-aware spaces evaluate a richer influence action space.

C. Reproducibility Details

For reproducibility, we report the main hyperparameters and compute settings used in our experiments. Unless otherwise stated, all Bradley–Terry models were fit with `hessian_ridge=0.0`. In the grouped Newton refinement, we used a small numerical stabilization ridge of 10^{-8} . For the smooth Kendall’s τ surrogate, we used temperature $T = 0.1$ in the player-level grouped-drop analysis and $T = 0.5$ in the robustness and action-curve analyses. All influence computations used the `1sn` approximation throughout.

All main experiments were run with 4 CPUs and 32 GB memory per task. We used the same hyperparameter settings across datasets unless explicitly noted above.

D. Other Results

D.1. Dataset-level robustness scores

To summarize robustness across datasets, we define normalized audit scores for the three leaderboard conclusions studied in the main text. Let B denote the maximum perturbation budget used for the corresponding audit and let b^* be the minimum number of influence-guided actions required to change the audited criterion after refitting the Bradley–Terry model. If no change occurs within the budget, we set $b^* = B$. We define

$$R_{\text{Top-1}} = \frac{b_{\text{Top-1}}^*}{B_{\text{Top-1}}}, \quad R_{\text{CITrace}} = \frac{U(B_{\text{CITrace}})}{U(0)}, \quad R_\tau = \tau(B_\tau),$$

where $R_{\text{Top-1}}$ measures the normalized budget needed to change the point-estimate Top-1 boundary, R_{CITrace} reports the remaining trace-uncertainty proxy under the fixed CI-trace audit budget, and R_τ reports the Kendall- τ value after the fixed global-ranking audit budget. Thus, lower $R_{\text{Top-1}}$ indicates that fewer actions are needed to change the local leaderboard conclusion, while lower R_{CITrace} and lower R_τ indicate stronger degradation under the corresponding fixed-budget audits.

For a compact overall summary, we also report

$$R_{\text{all}} = \frac{1}{3} (R_{\text{Top-1}} + R_{\text{CITrace}} + R_\tau). \quad (18)$$

This aggregate is not intended to replace the three component scores. Instead, it provides a single coarse audit number while preserving the decomposition into local point-estimate fragility, CI-trace robustness, and global ranking stability.

D.2. Kendall’s τ and trace uncertainty curves

Figures 4–9 report the full perturbation-budget curves for other datasets beyond Arena 55k (shown in the main paper, Figure 2). The same qualitative patterns hold throughout: influence-guided `Flip` causes the steepest τ degradation, `Add` greedy variants dominate trace uncertainty reduction, and random selection is negligible in every case. Datasets with sparser comparison graphs (ATP) exhibit much larger absolute τ degradation—reaching near 0 in the ATP case—confirming that graph density is a key determinant of global robustness.

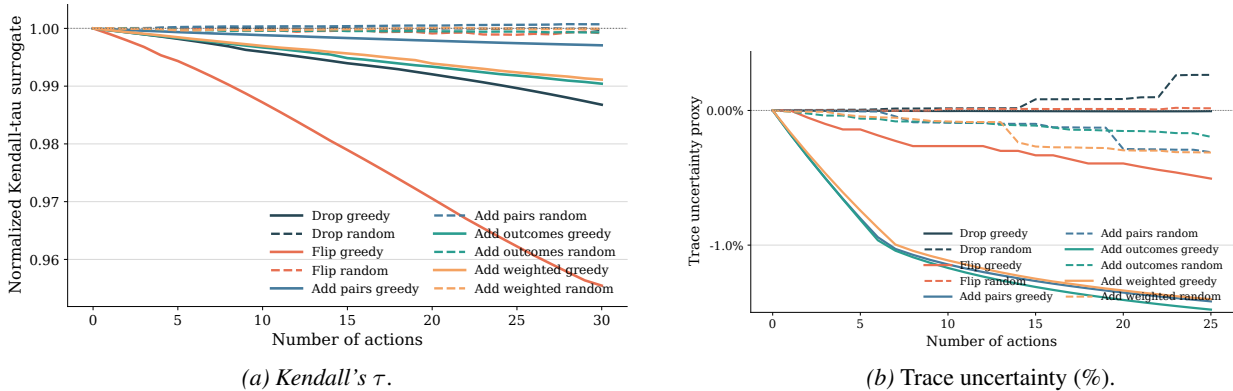


Figure 4. Perturbation-budget curves for **Chatbot Arena LLM Judges** (49,938 matches, 64 models). Influence-guided `Flip` degrades τ most steeply; `Add` greedy variants dominate uncertainty reduction.

A Unified Perturbation Framework for Leaderboard Stability and Manipulation

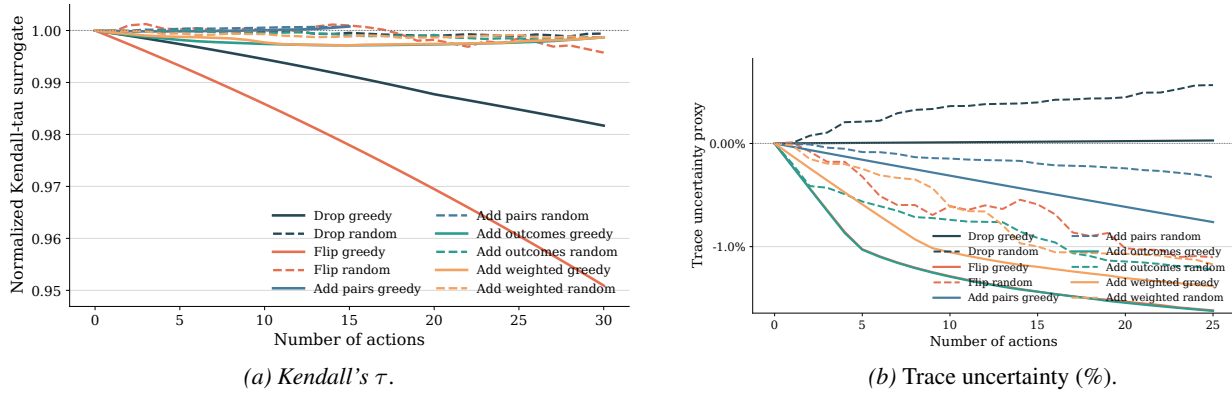


Figure 5. Perturbation-budget curves for **MT-Bench Human Judgments** (3,355 matches). τ degrades more slowly than on crowd-sourced datasets, consistent with Table 1 showing MT-Bench is the most robust dataset.

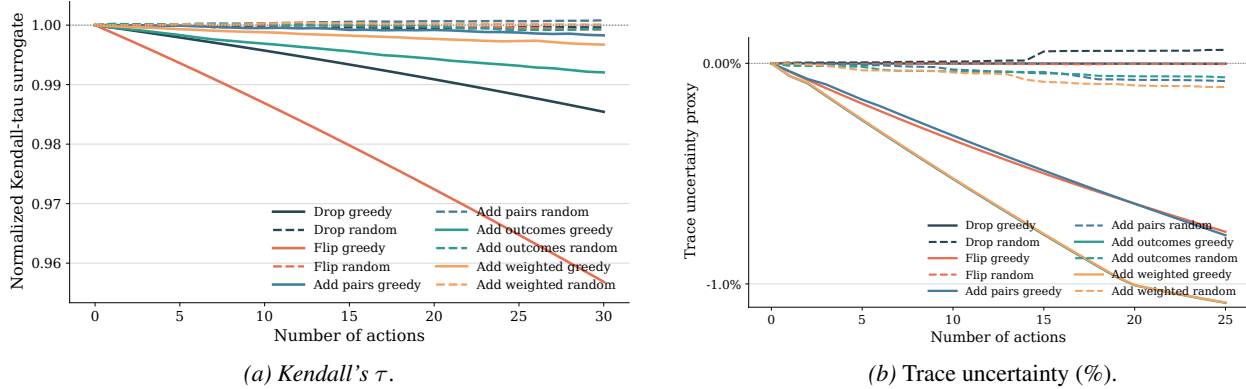


Figure 6. Perturbation-budget curves for **NBA Elo Top-50 Teams** (109,892 matches). The dense graph limits per-action impact on τ , but the Add-dominates-uncertainty reversal persists.

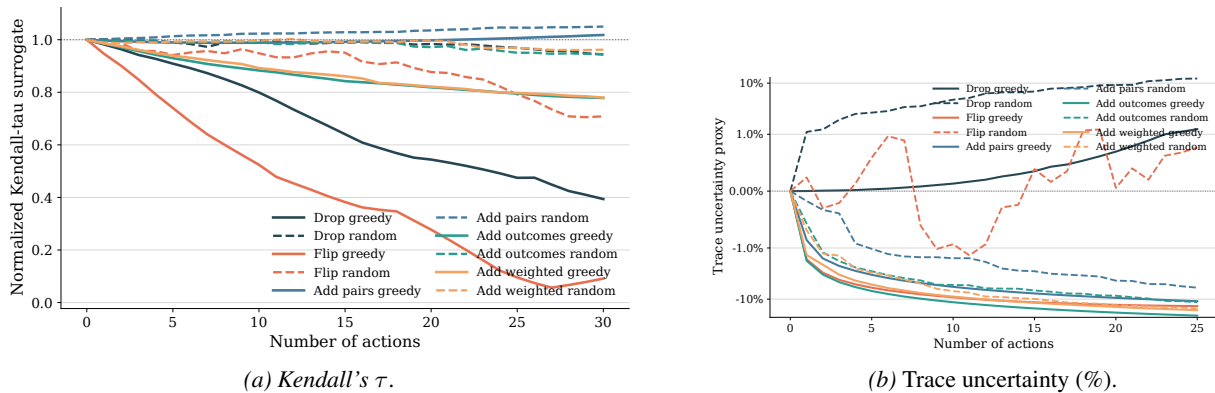


Figure 7. Perturbation-budget curves for **ATP Top-10 Matchups** (278 matches). The extremely sparse graph makes this the least robust dataset: Flip drives τ to near 0 within 30 actions.

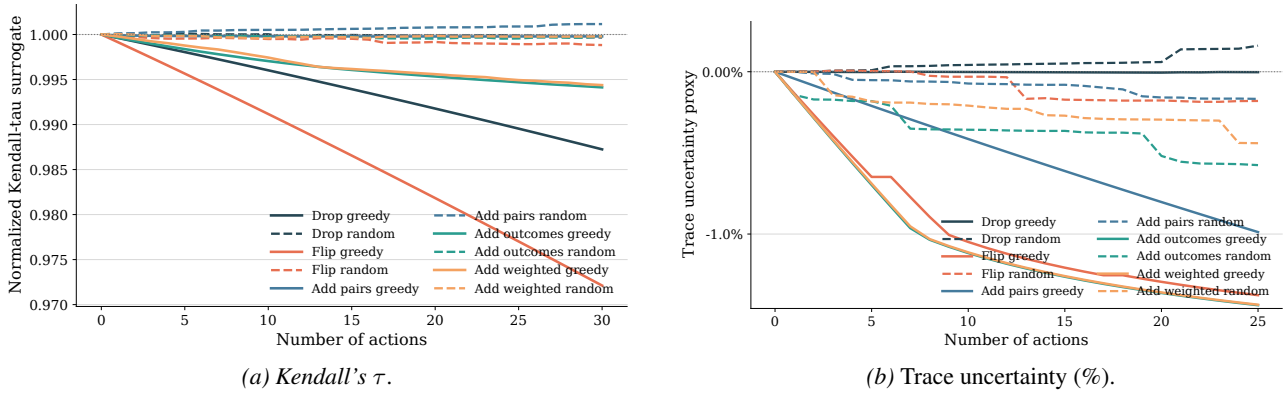


Figure 8. Perturbation-budget curves for **Vision Arena** (29,849 matches). Patterns mirror those of Arena 55k: Flip leads τ degradation and Add variants dominate uncertainty.

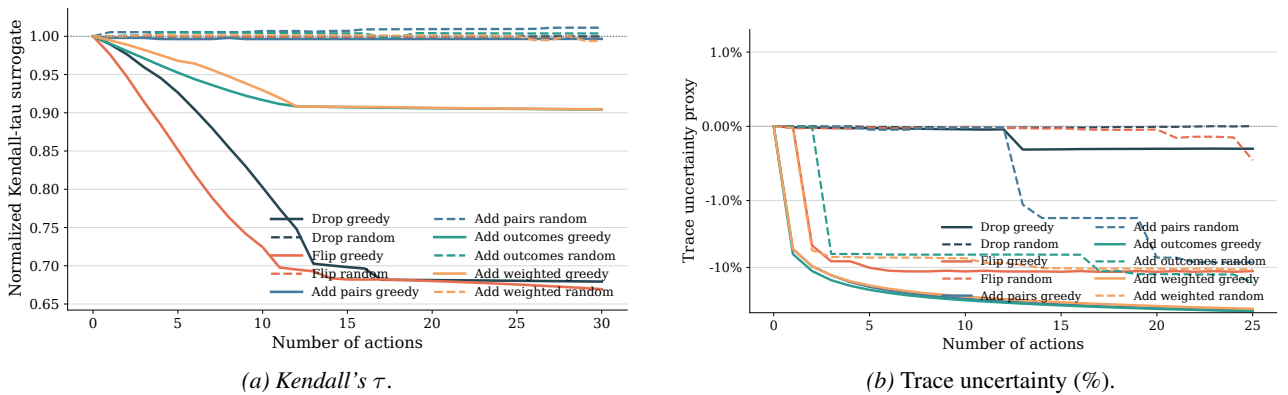


Figure 9. Perturbation-budget curves for **WebDev Arena** (10,501 matches). The sparse graph results in severe τ degradation (≈ 0.65 under Flip) while Add variants still dominate uncertainty.

D.3. Top- k Membership Change: CI-Aware vs. Non-CI-Aware

Figure 10 compares the effect of influence-guided drop perturbations on Top- k membership under the standard point-estimate objective and the stricter CI-aware objective. In the non-CI-aware setting, a membership change is counted as soon as the estimated skill ordering crosses the Top- k boundary, even if the affected models have overlapping confidence intervals; on Arena 55k with $k = 22$, this occurs after only 2 targeted `Drop` actions. By contrast, the CI-aware setting requires a stronger form of manipulation: the promoted model must not only cross the boundary in point estimate, but must do so with sufficient statistical separation from the displaced model, requiring 19 targeted `Drop` actions in the same setting. The comparison therefore separates two notions of robustness: instability of the displayed ranking versus instability that remains significant after accounting for uncertainty. As expected, CI-aware Top- k changes require stronger perturbations, but the fact that targeted drops can still alter membership under this stricter criterion shows that the leaderboard is not only locally sensitive but can also be vulnerable in a statistically meaningful sense.

E. Ablation

E.1. Specification of matches

Influential matches reflect structure, not exposure. Figure 11 reports correlations between one-step match effects and simple match-level covariates across match-specification objectives. We consider four covariates: *match count*, the number of times the same unordered pair appears in the training data; *bridge variance*, a match-level structural diversity score computed from the variance of the opponents faced by the two endpoint players, intended to capture whether the match connects players with broad or heterogeneous comparison neighborhoods; *closeness*, the log-transformed absolute BT skill gap between the two matched players; and *surprise*, the discrepancy between the observed outcome and the probability

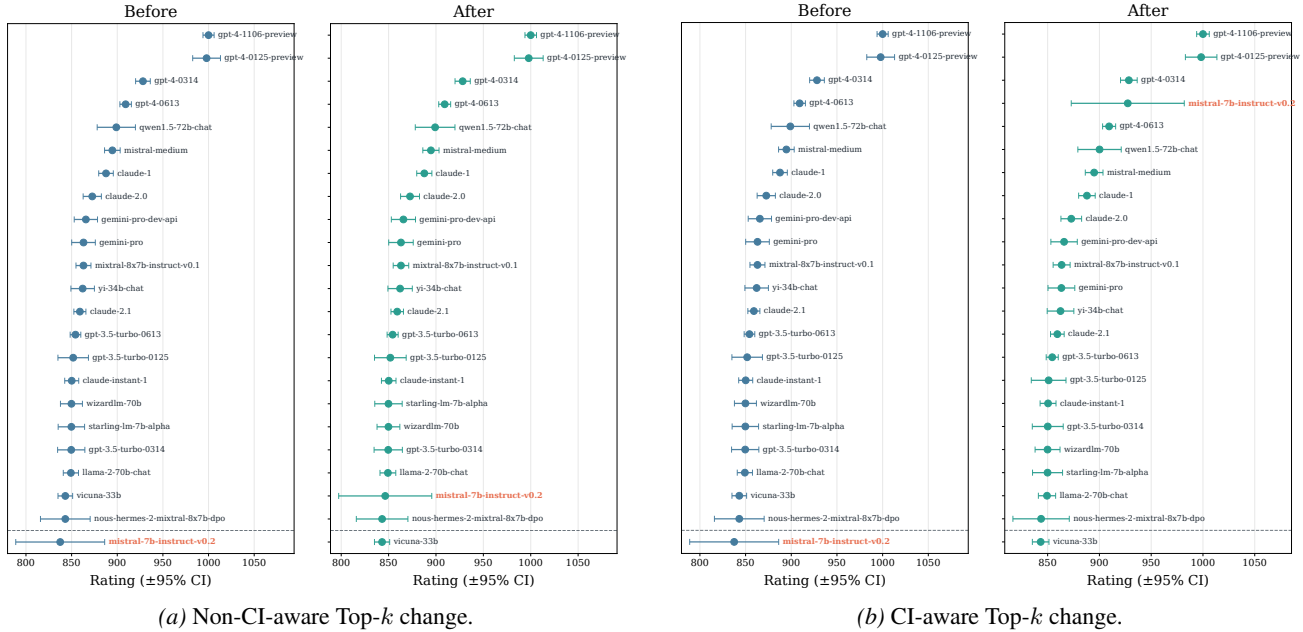


Figure 10. Top- k membership change under point-estimate and CI-aware objectives on Arena 55k. For $k = 22$, influence-guided Drop changes the point-estimate Top- k boundary after 2 actions, while the stricter CI-aware criterion requires 19 actions for the outsider to overtake the insider with non-overlapping confidence intervals. The two panels show that boundary membership is fragile, yet CI-aware success requires a more statistically robust change.

predicted by the fitted BT model. Because the natural direction of degradation differs across objectives, decreasing Kendall’s τ indicates worse global ranking stability, whereas increasing trace uncertainty indicates worse uncertainty, we interpret signed correlations relative to each objective rather than as a universal measure of influence magnitude. Under this objective-specific interpretation, match count is not a stable proxy for influence: frequently observed pairs are not consistently the ones whose perturbation most harms the leaderboard objective. By contrast, bridge variance and closeness often become strongly associated with large effects under flip objectives, suggesting that influence concentrates on structurally important and closely contested comparisons, where reversing an outcome can propagate broadly through the ranking. The surprise covariate is more heterogeneous, with its association changing across add, drop, and flip actions, especially for Kendall’s τ , player uncertainty, and trace uncertainty. Overall, the figure suggests that influential matches are objective-dependent and are shaped more by structural uncertainty and near-boundary comparisons than by raw exposure alone.

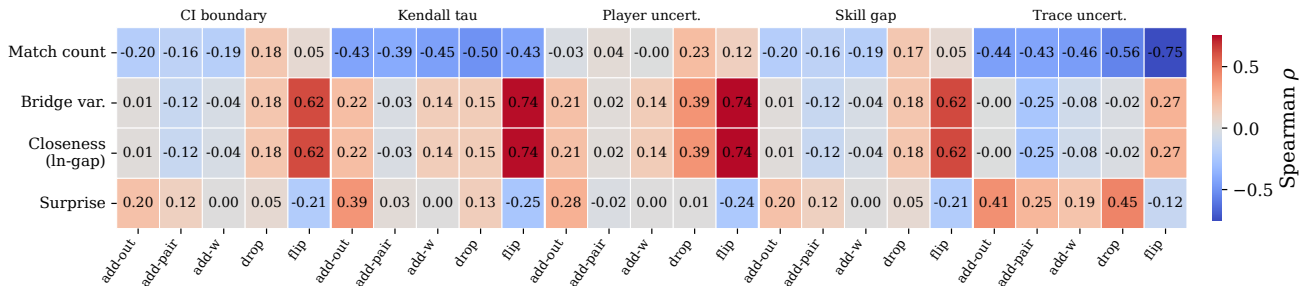


Figure 11. Signed Spearman correlations between one-step match influence scores and four match-level covariates in Chatbot Arena 55k, grouped by objective family and action type. Across most objectives, match count is weakly to strongly negative, while bridge variance and closeness are especially large for flip, highlighting the importance of uncertain, tightly contested comparisons in determining influential matches.

Influence concentrates on specified competitors. Under match-specific objectives, influence rankings concentrate sharply on the intended entities. Every objective that explicitly targets a player or player pair achieves a focus-hit rate of 1.0 at

top-20 selected action across all five action types: `Drop`, `Flip`, `Add-Pairs`, `Add-Outcomes`, and `Add-Weighted`. We test this on Arena55k, for the CI-boundary and skill-gap objectives defined on `gpt-4-1106-preview` versus `gpt-4-0125-preview`, as well as the player-uncertainty objective for `gpt-4-1106-preview`, all top-20 influential matches involve the designated focal player(s). Thus, when the objective specifies which competitors or matchups matter, influence does not spread diffusely across the dataset; it identifies edits tightly aligned with the target player or pair. In contrast, global objectives such as *Kendall's τ* and trace uncertainty have no designated focal player, so this focus-hit metric is not applicable to them.

E.2. Effect of k on top- k membership

To study how the choice of k affects the amount of data needed to change top- k membership, we vary k and measure the required perturbation budget. Figure 12 shows that the effort required to change the top- k ranking in Arena55k is highly non-monotonic in k and depends strongly on the action type. The most striking feature is the sharp peak at $k = 3$, where adding pairs becomes dramatically more expensive than all other interventions, while adding weighted outcomes, adding outcomes, and dropping results remain substantially smaller, and flipping outcomes is consistently the least costly. Beyond $k = 5$, the required number of actions drops quickly for all methods and remains low and fairly stable through $k = 10, 20, 40$, indicating that deeper top- k boundaries are easier to perturb than the top few ranks. Overall, the figure suggests that the top-3 boundary is the most structurally resistant part of the ranking, especially when interventions are constrained to adding new comparisons rather than altering existing ones.

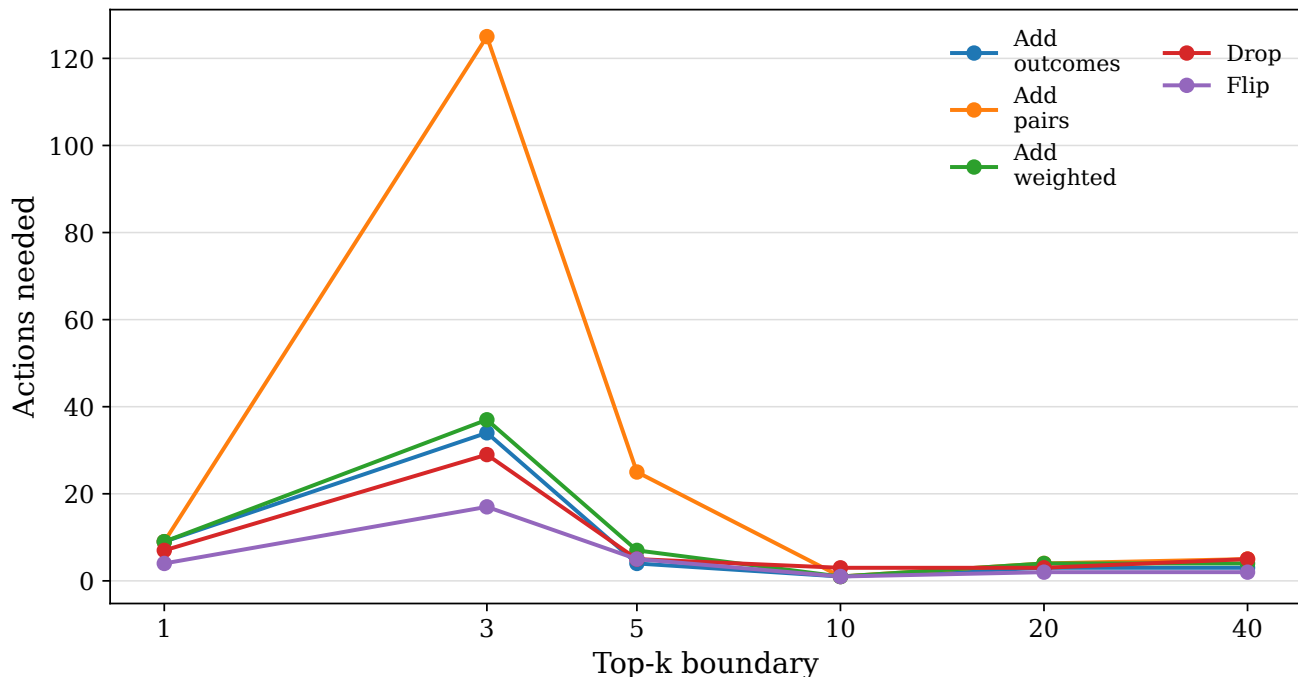


Figure 12. Top- k robustness is boundary-dependent. Minimum number of influence-guided actions required to change the Top- k set on Chatbot Arena 55k as k varies. The Top-3 boundary is the most structurally resistant, requiring roughly 120 actions, while several larger- k boundaries require only a few actions. Thus, robustness is not monotonic in k but depends on the local structure of the leaderboard near each cutoff.

E.3. Structural predictors of player-removal impact

To understand what drives player-removal impact, Figure 13 reports four association statistics between each structural player feature and absolute $|\tau|$ -influence, aggregated across datasets. We consider *degree*, the number of matches a player participates in; *bridge variance*, how unevenly a player connects otherwise weakly connected regions of the graph; *closeness*, how centrally located a player is in terms of shortest-path distances; and *surprise*, how unexpected a player’s outcomes are relative to model predictions.

Degree is the strongest and most consistent predictor: it achieves the highest Pearson ($|r| = 0.41$) and Spearman ($|\rho| = 0.34$)

correlations, and the largest quartile separation (Q4–Q1 mean $z = 0.65$, Cohen’s $d = 0.86$). Bridge variance and closeness show weaker or less consistent associations, while surprise is only weakly associated.

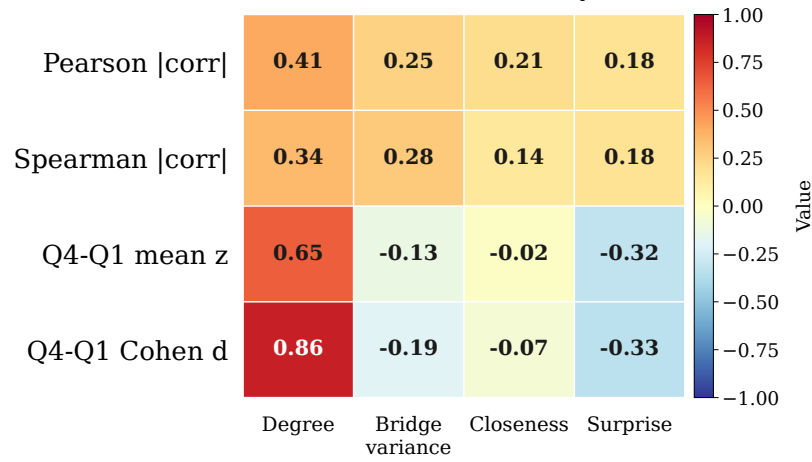


Figure 13. Association between structural player features and absolute $|\tau|$ -influence, aggregated across datasets using Pearson and Spearman correlations, Q4–Q1 mean z -score, and Cohen’s d . Degree is the strongest and most consistent predictor of player-removal impact.