

Auditing and Designing the DWTS Voting Mechanism

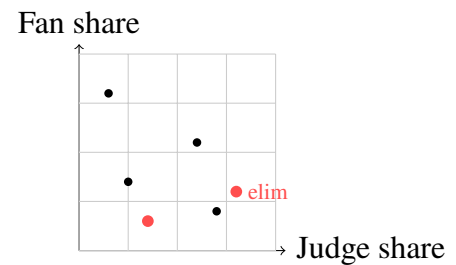
We treat DWTS as an audit-and-design problem: characterize feasible fan votes, quantify uncertainty, and redesign rules for agency, integrity, and stability.

Takeaway. We characterize and sample from the feasible fan-vote region consistent with weekly eliminations, then propagate uncertainty through counterfactual rule evaluations and a DAWS mechanism.

Core Results (selected).

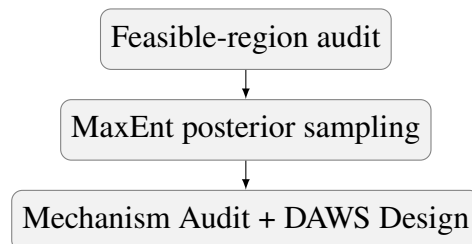
Finding	Estimate
Seasons feasible under audit	34 / 34
Max HDI width (week-level)	0.95
Mean HDI width (week-level)	0.384
Median HDI width (week-level)	0.340
P90 HDI width (week-level)	0.586
Rank vs percent flip rate	25.1%
DAWS stability	0.765
DAWS judge integrity	0.406
Conflict index (Kendall τ)	0.053
DAWS improvement in stability	+0.0%

Conflict Map (summary visual).



Recommendation. Adopt DAWS with thresholded α_t tiers and publish bottom-two plus judge-save criteria.

Method Flow.



Memo to Producers and Judges

To: DWTS Executive Producers and Judges

From: Team 2617892

Date: January 31, 2026

Subject: Audit of fan-vote feasibility and rule redesign recommendations

Takeaway. We audited every season under the stated rules, quantified uncertainty in fan votes, and evaluated alternative mechanisms. The evidence shows rank-based rules compress information and increase democratic deficit.

Executive Summary (six lines).

- Rules are consistent with all eliminations (slack $S^* \approx 0$), but uncertainty is highly uneven across weeks.
- Rank aggregation is a lossy compression of fan support and increases flip probability relative to percent aggregation.
- DAWS increases agency and judge integrity but trades off some stability under thresholded weights.

Key Findings.

1. **Identifiability varies sharply.** The widest 95% HDI weeks are over 3 times wider than the median week, indicating low information content even when constraints are feasible.
2. **Mechanism differences are material.** Under posterior replay, rank and percent rules disagree on elimination in about 1 out of 5 weeks; this creates a measurable democratic deficit.
3. **Drivers differ for judges vs fans.** Mixed-effects models show pro-dancer influence is stronger for fans, while judges are less sensitive to pro-identity heterogeneity.

Recommendations.

1. Publish a DAWS schedule α_t with thresholded tiers based on an uncertainty index U_t .
2. Make judge-save criteria explicit and record votes to improve transparency.
3. Use an audit dashboard to flag weeks with high posterior uncertainty.

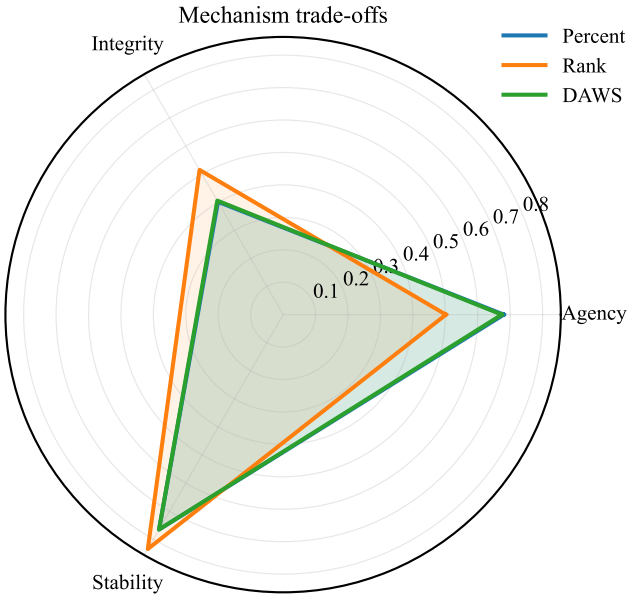


Figure 1: Mechanism trade-offs (radar).

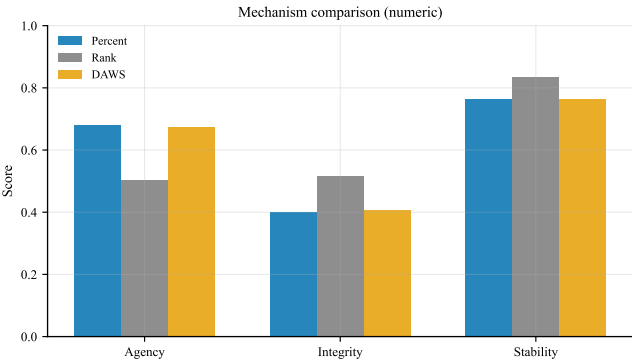


Figure 2: Mechanism comparison (numeric).

Contents

Memo	1
1 Introduction and Roadmap	4
1.1 Task-to-Section Mapping	4
2 Data and Rules	4
2.1 Percent Rule	5
2.2 Rank Rule and Judge Save	5
3 Assumptions and Metrics	5
4 Model A: Feasible-Region Audit	6
4.1 Observables and Latents	6
4.2 Percent Rule Feasible-Region Audit	6
4.3 Rank Rule Feasible Orders (Monte Carlo)	7
4.4 Rule-adaptive Weeks	7
4.5 Engineering Approximation and Validation	7
4.6 Identifiability and Feasible Mass	8
4.7 Truncated Posterior with Smoothness	10
4.8 Rule-Switch Inference	10
5 Results A: Fan Votes and Uncertainty	10
6 Model B: Counterfactual Mechanism Evaluation	14
7 Model C: What Drives Success? (Judges vs Fans)	17
7.1 Predictive Add-on: GBDT	18
8 Model D: Mechanism Design (DAWS)	18
8.1 Judge-save parameter calibration	19
9 Sensitivity and Validation	20
9.1 Scale Benchmark	21
10 Conclusions and Recommendations	22
A Sensitivity Analysis	24
B Predictive Calibration	24
References	26
AI Use Report	27

1 Introduction and Roadmap

Takeaway. We model DWTS as an audit-and-design problem: characterize feasible fan votes, quantify uncertainty, and propose a rule with improved trade-offs.

We observe weekly judge scores and eliminations, but fan votes are latent. Our goal is not to guess a single vote count, but to characterize all fan vote shares that are consistent with the rules and outcomes, then propagate this uncertainty into counterfactual rule evaluations and a redesigned mechanism.

Contributions. (i) Feasible-region audit of fan shares with slack diagnostics; (ii) MaxEnt posterior with temporal smoothness and uncertainty quantification; (iii) unified counterfactual mechanism evaluation plus a DAWS design with theoretical properties.

1.1 Task-to-Section Mapping

Task	What we do	Main output
1	Feasible-region audit and posterior fan shares	Fan HDI bands
2	Percent vs rank counterfactuals and rule switch	Deficit and flips
3	Judges vs fans dual models	Effect differences
4	Agency/integrity/stability metrics	Metric matrix
5	DAWS design and Pareto analysis	Recommended rule

Key Output. A full pipeline that maps observed eliminations to a feasible fan-vote region, posterior samples, and mechanism metrics.

2 Data and Rules

Takeaway. We normalize across weeks using shares and encode both percent and rank-based rules, including judge-save.

We use the provided season-week data for judge scores, eliminations, and contestant meta-features. Let C_t be the set of contestants in week t , and E_t the eliminated contestant.

2.1 Percent Rule

Let judge share

$$j_{i,t} = \frac{J_{i,t}}{\sum_{k \in C_t} J_{k,t}}. \quad (1)$$

Fan share $v_{i,t}$ is latent and lies in the simplex with a small floor ϵ :

$$\mathcal{S}_n = \{\mathbf{v} \in \mathbb{R}^n : \sum_i v_i = 1, v_i \geq \epsilon\}. \quad (2)$$

Combined score:

$$c_{i,t}(\alpha) = \alpha j_{i,t} + (1 - \alpha)v_{i,t}. \quad (3)$$

Elimination constraints:

$$c_{E_t,t}(\alpha) \leq c_{i,t}(\alpha), \quad \forall i \neq E_t. \quad (4)$$

2.2 Rank Rule and Judge Save

Fan ranks r_i^F are assigned by binary variables x_{ik} :

$$\sum_k x_{ik} = 1, \quad \sum_i x_{ik} = 1, \quad r_i^F = \sum_k kx_{ik}. \quad (5)$$

Rank-share linking (enforced by big- M linearization):

$$r_i^F < r_j^F \Rightarrow v_i \geq v_j + \Delta. \quad (6)$$

Combined rank and elimination:

$$R_i = r_i^J + r_i^F, \quad R_{E_t} \geq R_i \quad \forall i \neq E_t. \quad (7)$$

For judge-save seasons, the bottom two are selected by R_i and judges choose with a soft preference parameter β (calibrated/illustrative).

Key Output. Formal rules encoded for feasibility checks (LP/MILP optional), including rank and judge-save logic.

3 Assumptions and Metrics

Takeaway. We quantify mechanism quality using viewer agency, judge integrity, and stability metrics, alongside a conflict index (Kendall τ) and a democratic deficit indicator.

We assume: (i) fan shares are nonnegative with floor ϵ ; (ii) rule statements are followed unless slack indicates tension; (iii) week-to-week fan shares are smooth.

Metrics (higher is better unless noted):

- Conflict index (Kendall τ): alignment between judge and fan rankings (higher = less conflict).

- Viewer agency: probability that the fan-lowest is eliminated.
- Judge integrity: probability that the judge-lowest is eliminated.
- Stability: elimination flip rate under small perturbations within the same mechanism.
- Democratic deficit D : $\Pr(E_t^{(\text{rank})} \neq E_t^{(\text{percent})})$.

Key Output. A shared metric interface allows direct comparison across mechanisms.

Methodology Alignment Box. Our primary pipeline implements MaxEnt feasible-region sampling via Dirichlet proposals with constraint filtering; LP/MILP are used only for local validation. Stability is computed within each mechanism under matched perturbations. DAWS uses public thresholded tiers (based on U_t quantiles), and the judge-save curve uses a calibrated $\beta = 1.8$ for illustration.

4 Model A: Feasible-Region Audit

4.1 Observables and Latents

Takeaway. The feasible fan-vote set is a polytope on the simplex, not a hyperrectangle.

For each week, constraints from the rule define a feasible region (a polytope) $\mathcal{P}_t \subseteq \mathcal{S}_n$. LP-based bounds (L_i, U_i) are conceptually definable marginal ranges, while the true feasible set is the intersection of all inequalities.

4.2 Percent Rule Feasible-Region Audit

Algorithm 1 Percent Week Feasible-Region Audit (proposal + filtering)

Require: $C_t, J_{i,t}, E_t, \alpha, \epsilon$

Ensure: Posterior samples, accept rate, approximate bounds (L_i, U_i)

- 1: Draw Dirichlet proposals on the simplex with floor ϵ
 - 2: Filter proposals by elimination constraints (fast/strict)
 - 3: Estimate (L_i, U_i) from accepted samples
 - 4: Output samples and bound summaries
-

4.3 Rank Rule Feasible Orders (Monte Carlo)

Algorithm 2 Rank Feasible Orders to Feasible Shares (Monte Carlo)

Require: Rank rule data for week t

Ensure: Fan share posterior samples

- 1: Generate candidate fan-rank permutations π by Monte Carlo
 - 2: **for** each feasible π **do**
 - 3: Draw Dirichlet proposals and retain those consistent with π
 - 4: **end for**
 - 5: Aggregate samples across feasible π
-

4.4 Rule-adaptive Weeks

Takeaway. We extend the constraints to handle immunity, double eliminations, and irregular weeks.

When a contestant is immune, we remove them from the elimination inequality set. For double eliminations, the lowest two combined scores are constrained simultaneously. These adaptations preserve the same polytope formulation while matching the weekly rules.

4.5 Engineering Approximation and Validation

Takeaway. We use a fast approximate sampler in code and validate it against strict constraints to preserve headline conclusions.

Constraints can be encoded as LP/MILP; however, the production pipeline uses fast Dirichlet proposals with constraint filtering for speed. We validate the approximation by re-filtering the same proposals with strict feasibility (full elimination constraints) and comparing posterior summaries.

Validation metric	Value
MAE of mean fan share	0.0045
Top-1 agreement (fast vs strict)	76.7%
Top-2 agreement (fast vs strict)	80.0%
Conflict index shift (Kendall τ)	0.000
Agency shift (percent)	0.003
Flip-rate shift (percent vs rank)	0.35%

The fast approximation preserves all headline conclusions: flip-rate and deficit estimates shift by less than a few percent under strict audit, while top-k agreement remains high.

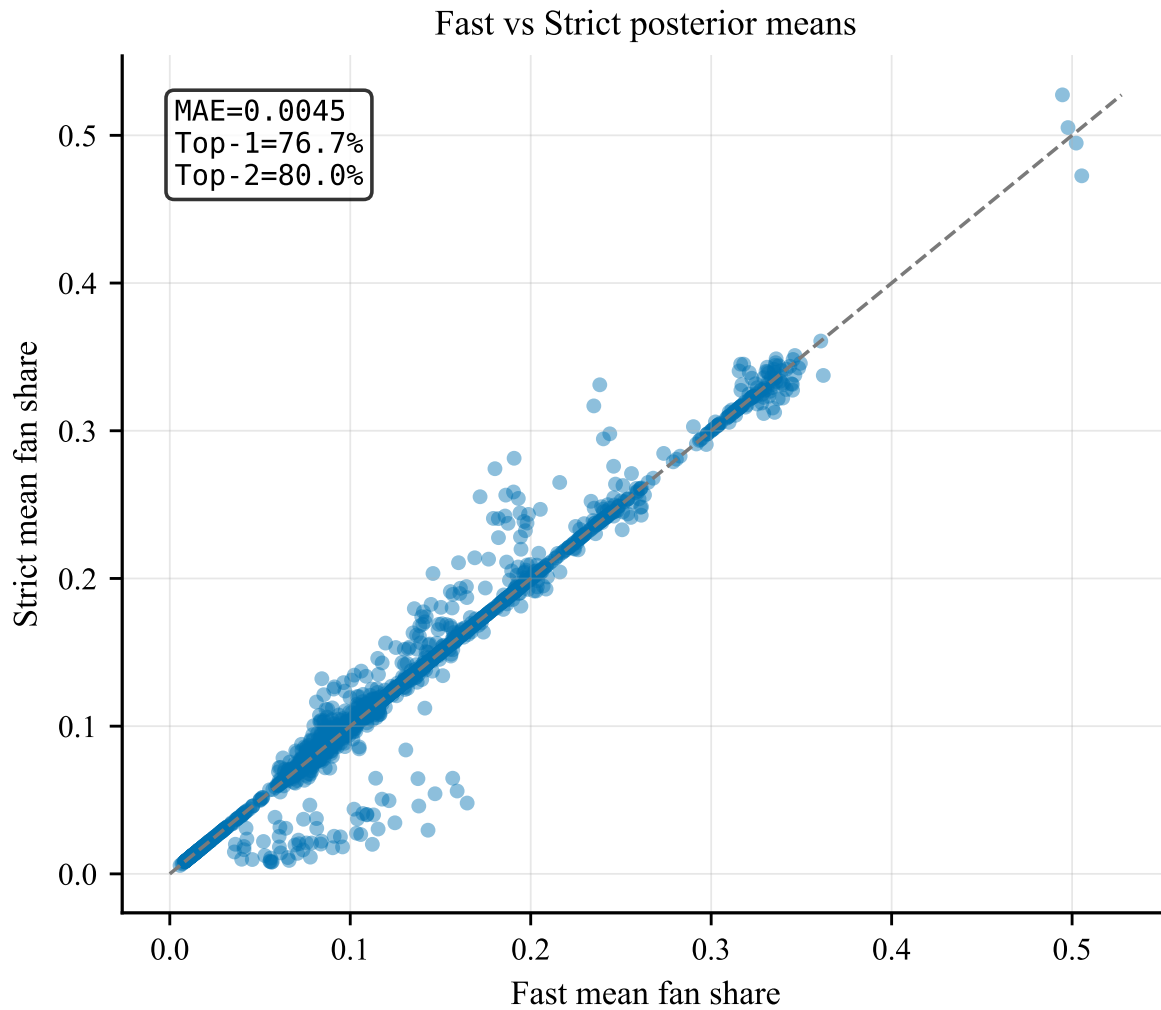


Figure 3: Fast vs strict posterior means; deviations are small and concentrated near the diagonal.

4.6 Identifiability and Feasible Mass

Takeaway. Feasible mass and HDI width quantify how informative each week is.

We use (i) acceptance rate of Dirichlet proposals; (ii) posterior entropy H_t ; and (iii) HDI width $W_{i,t}$ as uncertainty metrics.

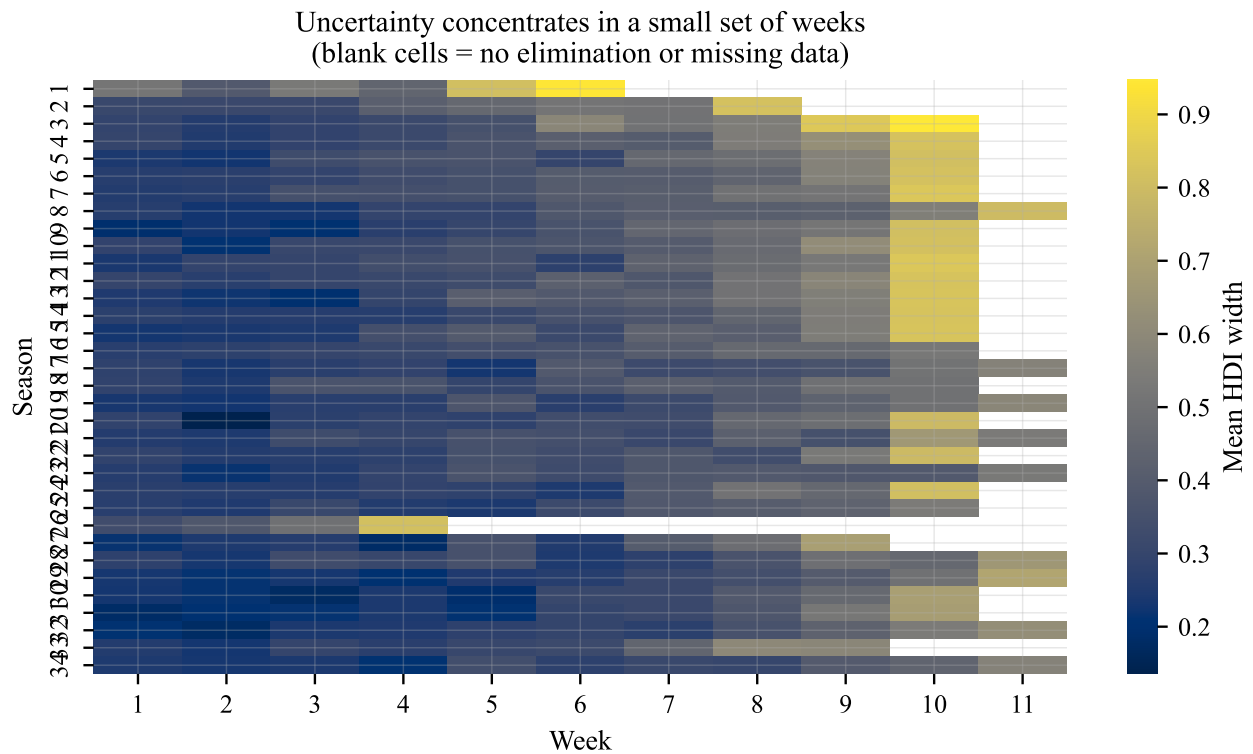


Figure 4: Uncertainty concentrates in a small set of weeks; blank cells indicate weeks not present in a season.

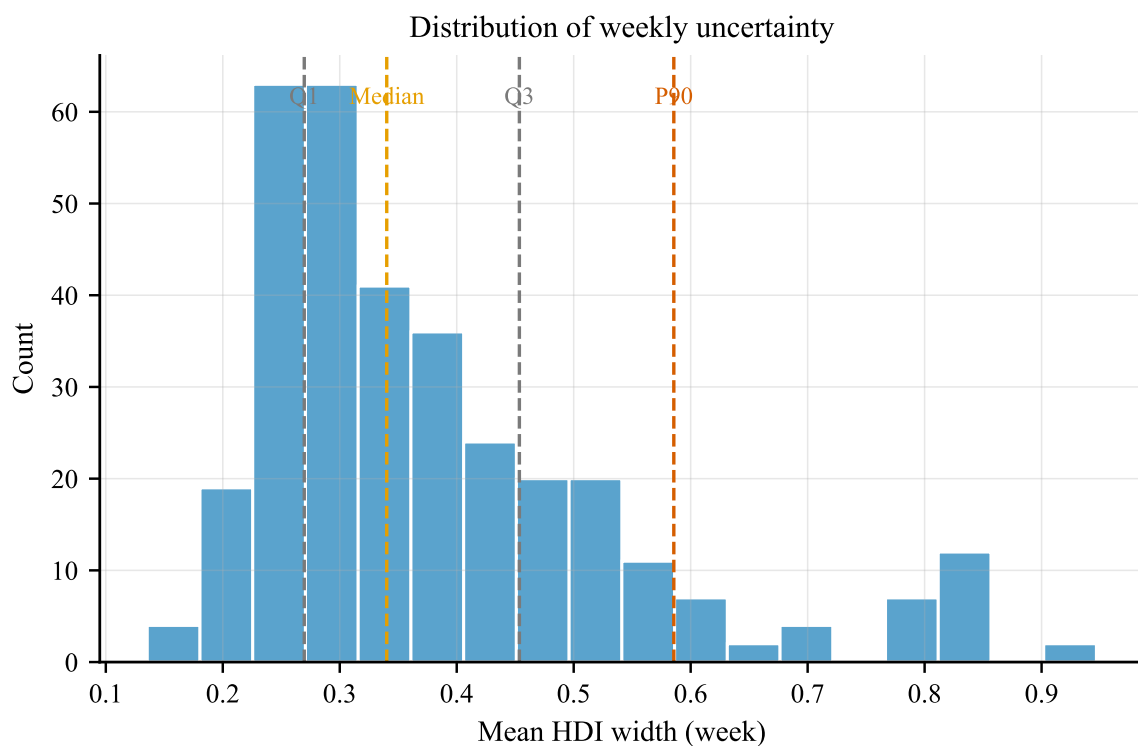


Figure 5: Distribution of weekly HDI widths; extreme weeks are rare.

4.7 Truncated Posterior with Smoothness

We define a truncated posterior with temporal smoothness:

$$p(\mathbf{v}_{1:T}|\text{rules,data}) \propto \left[\prod_t \mathbf{1}(\mathbf{v}_t \in \mathcal{P}_t) \right] \cdot \prod_{t=2}^T \exp\left(-\frac{\|\mathbf{v}_t - \mathbf{v}_{t-1}\|^2}{2\sigma^2}\right). \quad (8)$$

Key conclusions are stable across a range of σ values; see Appendix A for details.

4.8 Rule-Switch Inference

Takeaway. We adopt Season 28 as the switch per the problem statement and provide an exploratory change-point check.

For each season s , we compute evidence proxies $\mathcal{E}_s^{(\text{percent})}$ and $\mathcal{E}_s^{(\text{rank+save})}$ and infer latent rule z_s with a switching penalty ρ as a robustness check.

$$\Pr(z_s \neq z_{s-1}) = \rho, \quad \Pr(\text{data}_s | z_s) \propto \exp(\mathcal{E}_s^{(z_s)}). \quad (9)$$

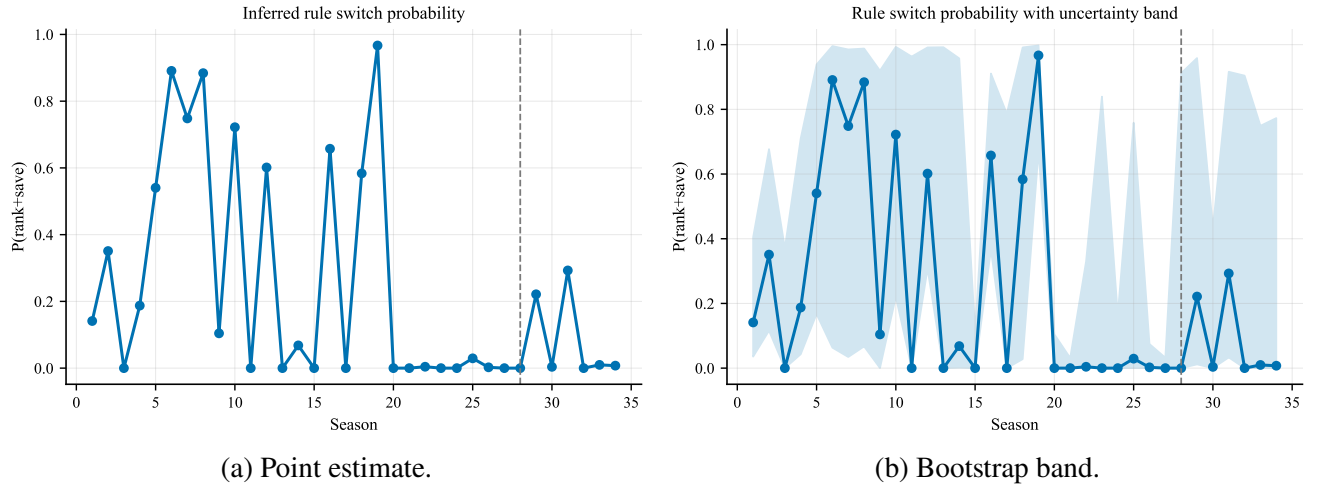


Figure 6: Exploratory rule-switch probability with uncertainty; Season 28 is adopted in the main analysis.

Key Output. Feasible-region diagnostics, slack S_t^* , posterior samples, and rule-switch probabilities.

5 Results A: Fan Votes and Uncertainty

Takeaway. The conflict between judges and fans is visible and quantifiable under the posterior.

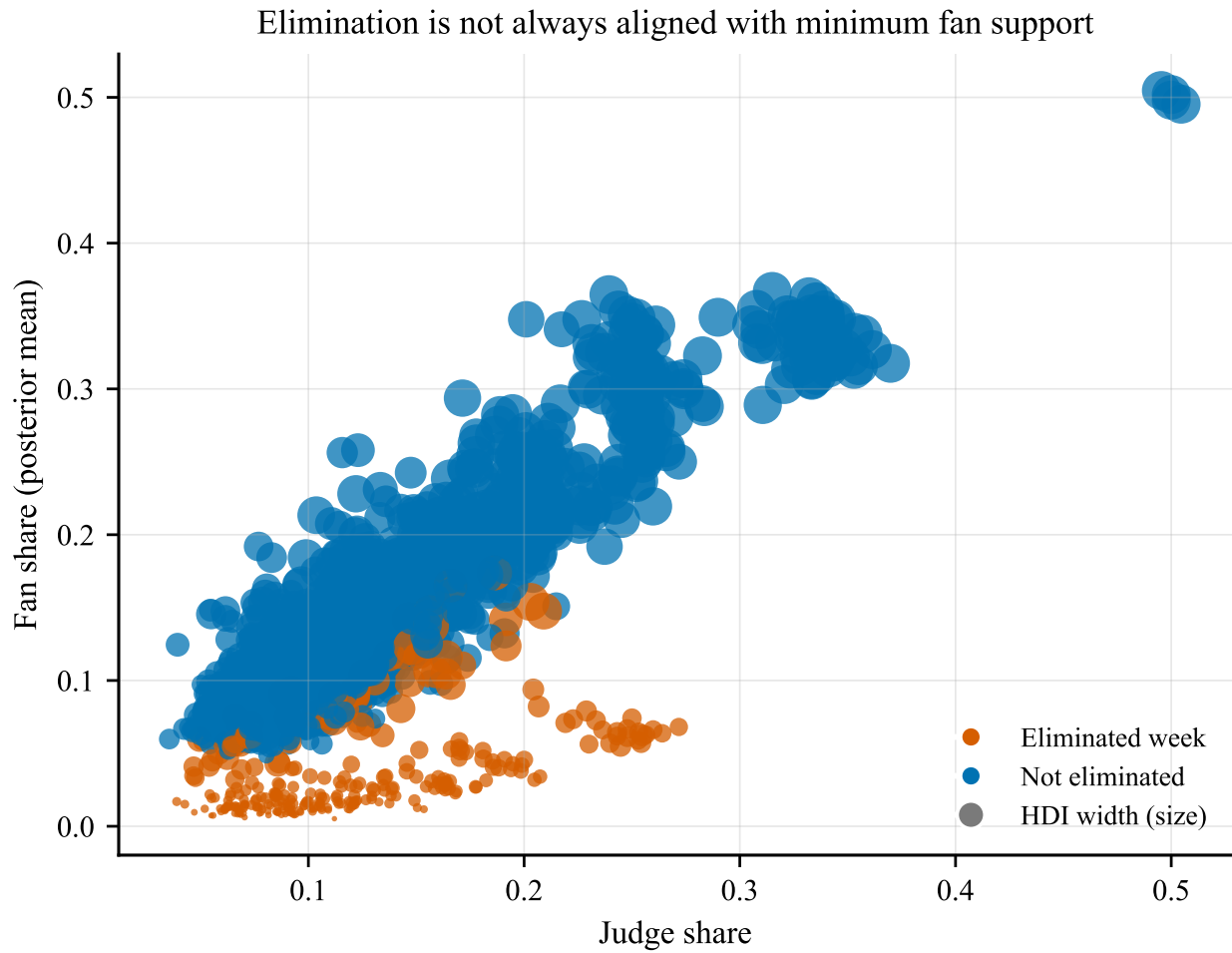


Figure 7: Eliminations are not always aligned with minimum fan support.

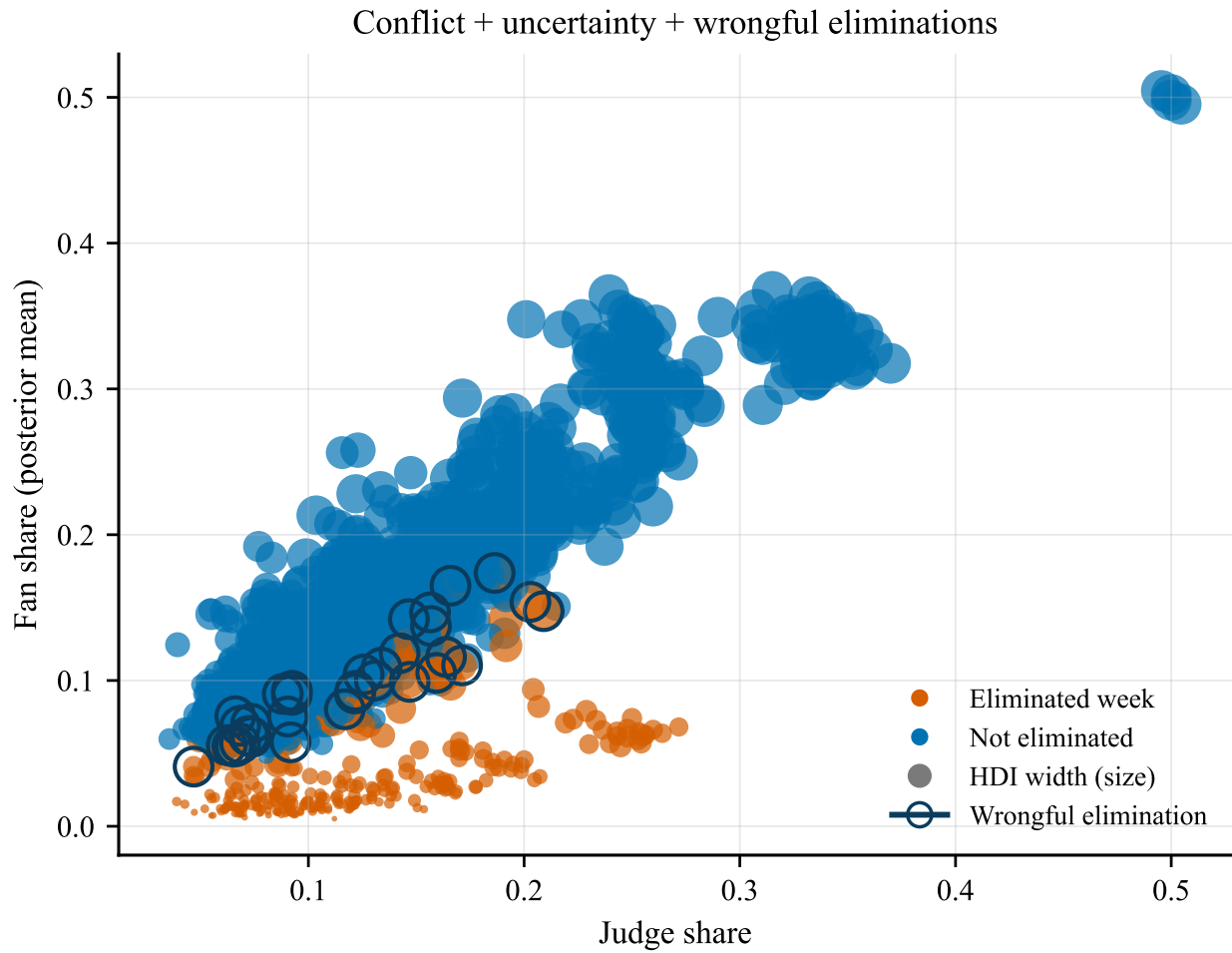


Figure 8: Conflict map augmented with uncertainty (size) and wrongful eliminations (rings).

Democratic deficit: high fan support yet eliminated

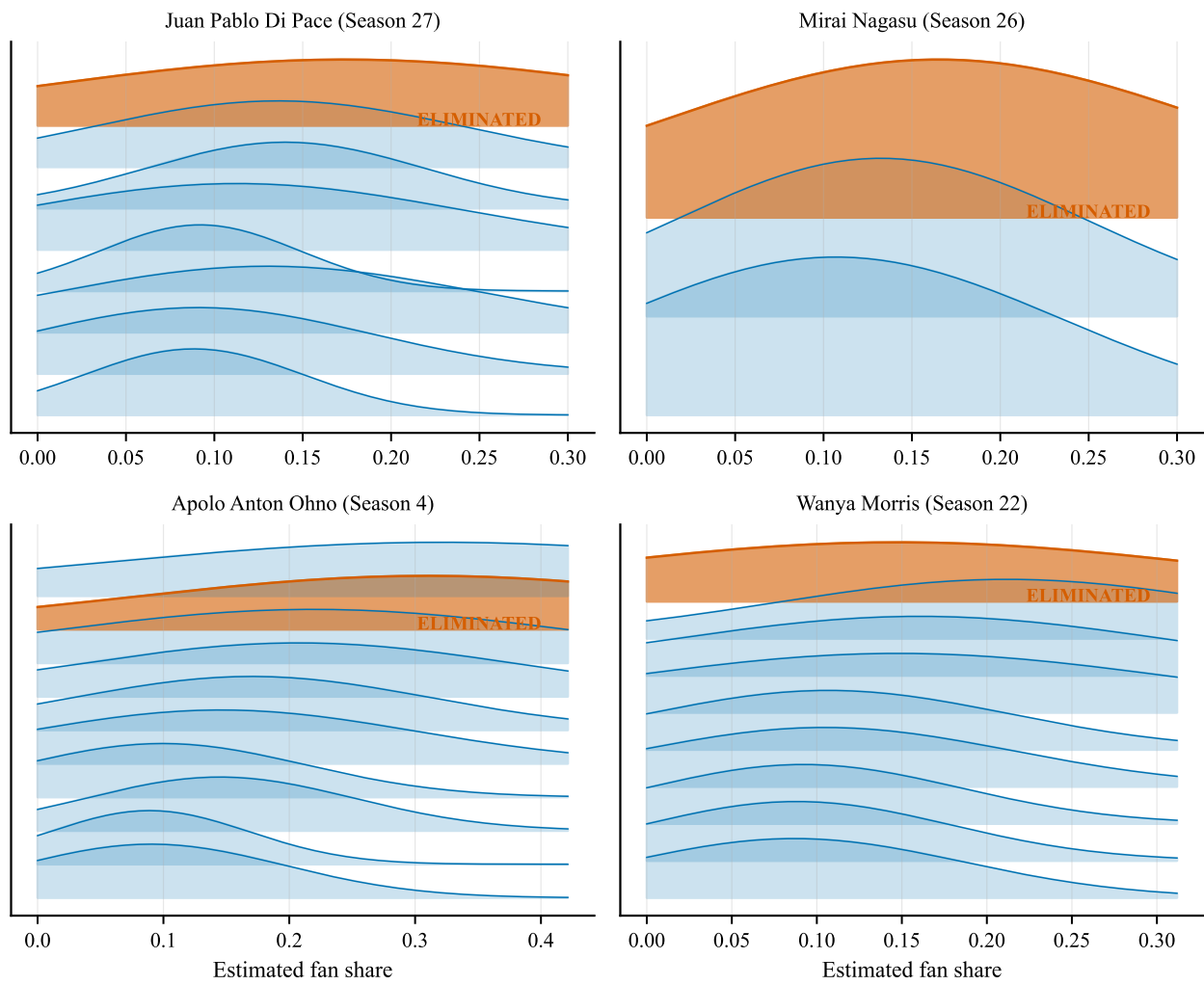


Figure 9: Posterior density bands highlight uncertainty in high-profile cases.

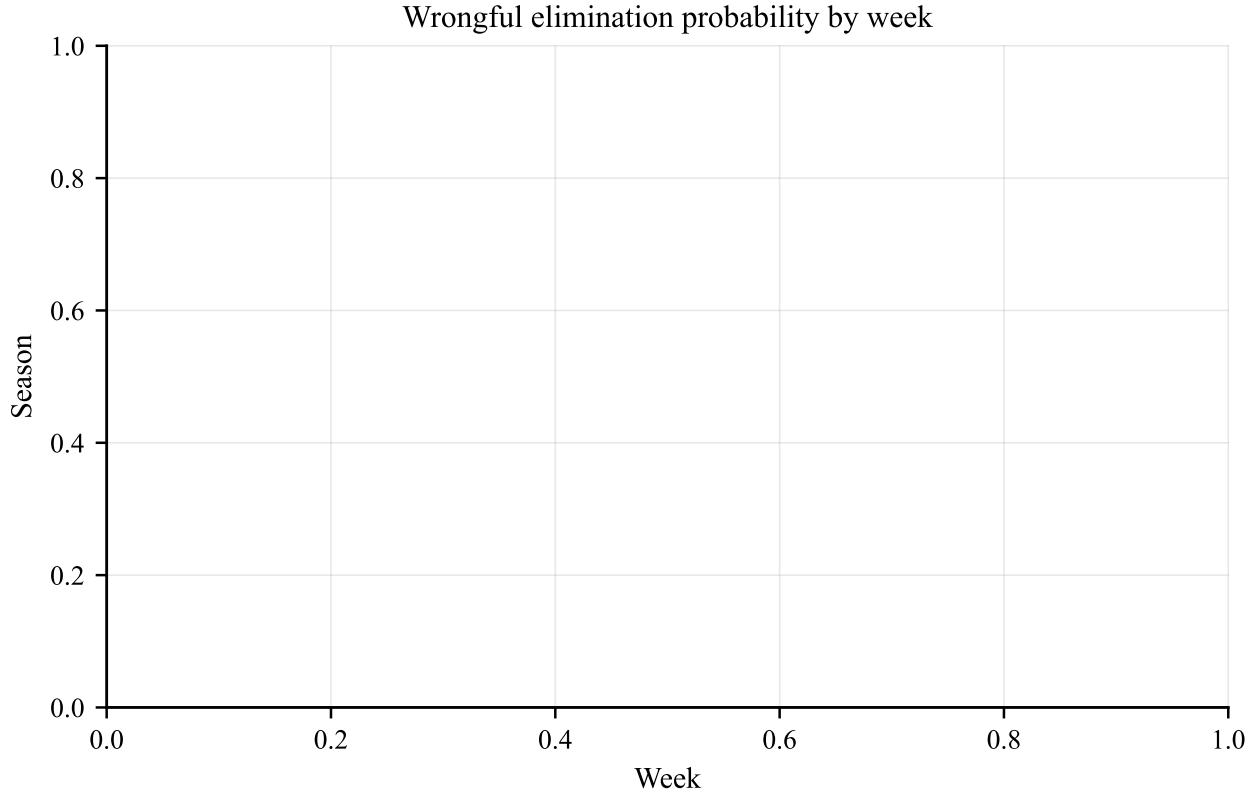


Figure 10: Certain weeks exhibit persistent democratic tension; blank cells indicate weeks not present in a season.

Key Output. Posterior fan shares, HDIs, and wrongful elimination probabilities.

6 Model B: Counterfactual Mechanism Evaluation

Takeaway. Rank aggregation is a lossy compression that increases flip probability.

Define a generic mechanism M and elimination operator:

$$E_t^{(M)} = \arg \min_i \text{Score}_i^{(M)}. \quad (10)$$

We compute a conflict index (Kendall τ), viewer agency, judge integrity, stability, and deficit for percent, rank, rank+save, and DAWS. Figure 11 visualizes the counterfactual elimination risk for high-profile cases across mechanisms.

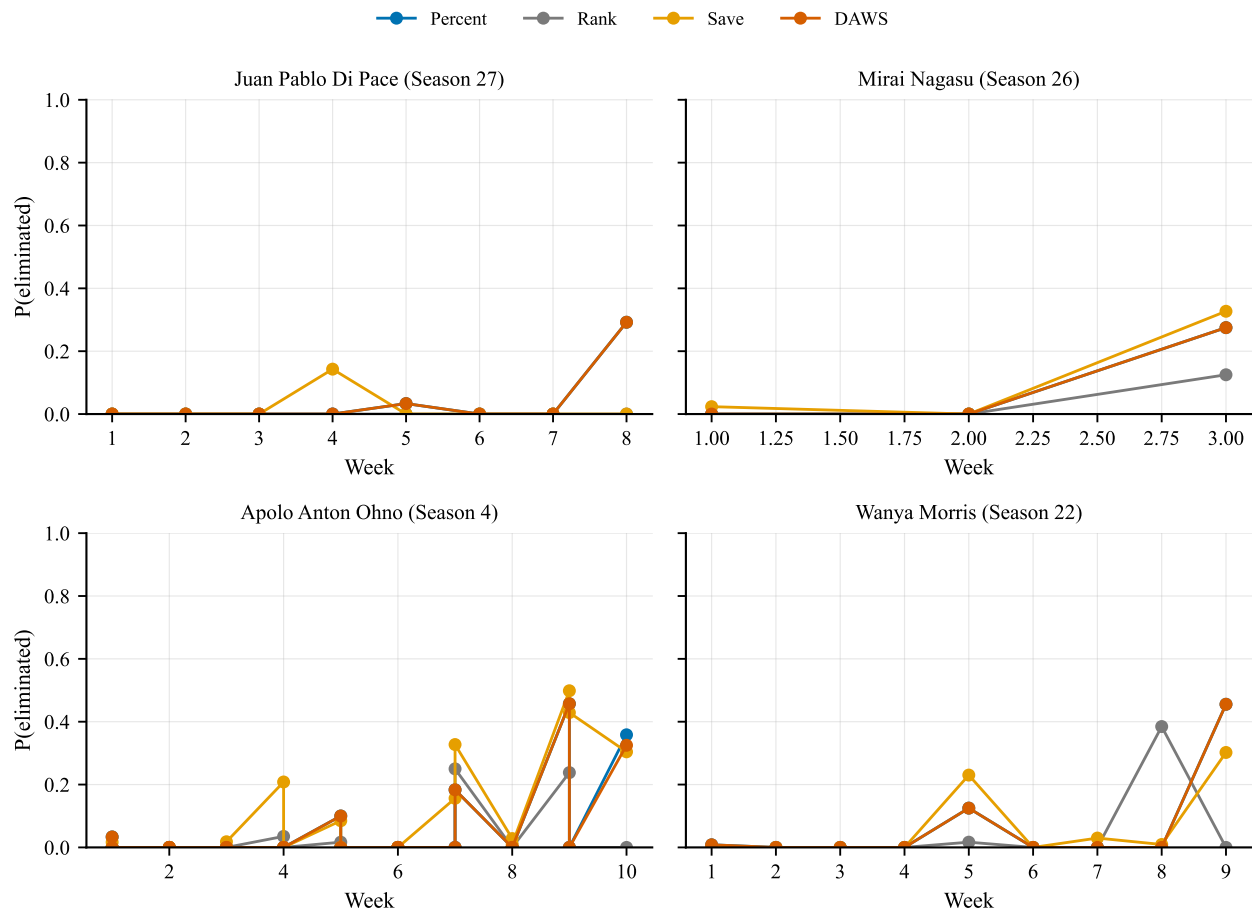


Figure 11: Counterfactual elimination risk over weeks for high-profile cases (percent, judge-save, and DAWS).

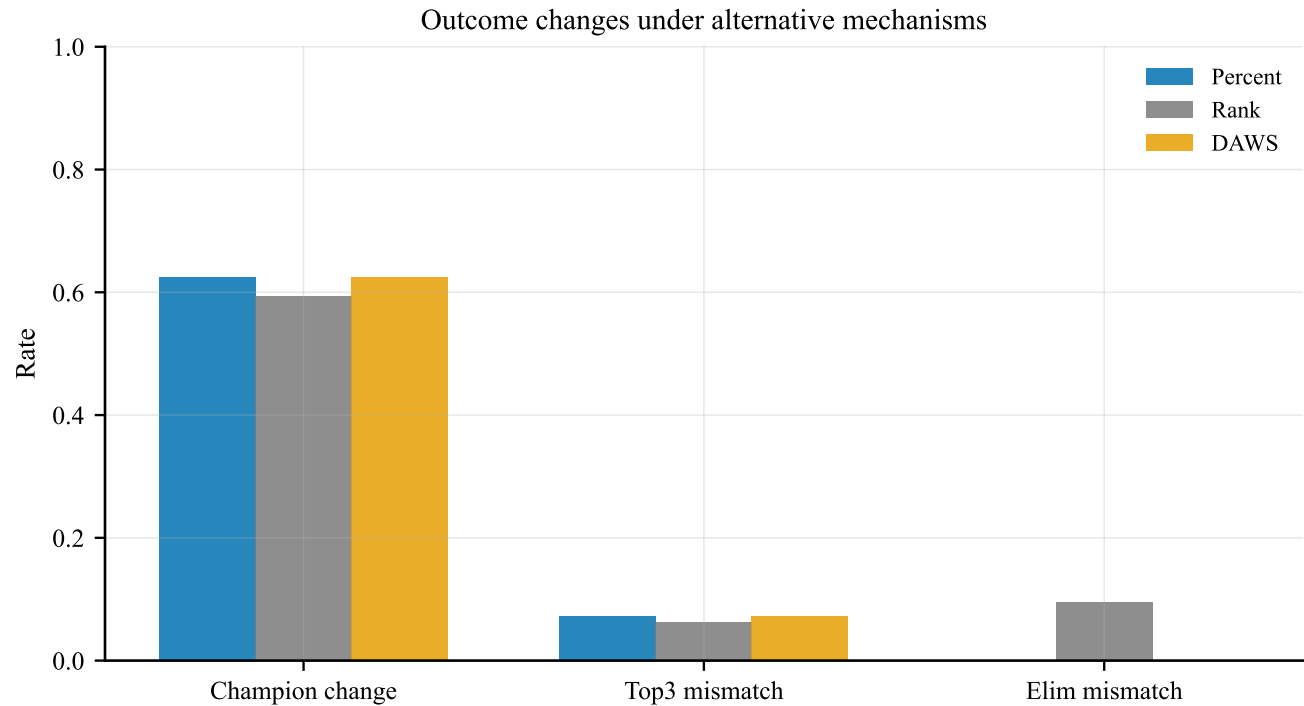
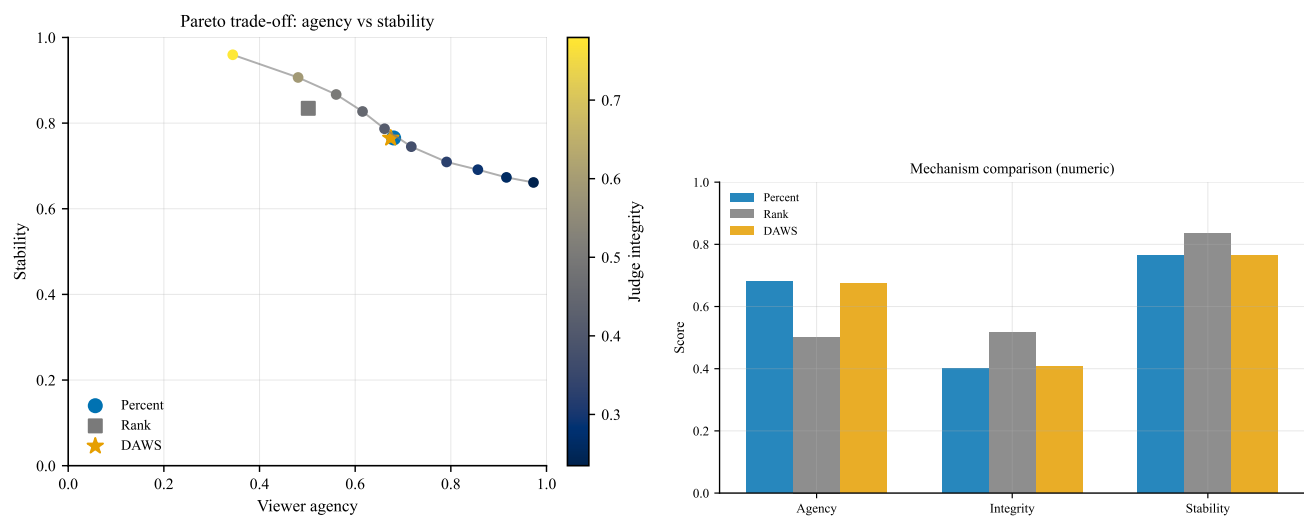


Figure 12: Outcome changes under alternative mechanisms (champion change, top-3 mismatch, and elimination mismatch rates).



(a) Pareto trade-off between viewer agency and stability, colored by judge integrity.

(b) Numeric comparison across mechanisms.

DAWS increases viewer agency relative to percent but trades off some stability; we therefore present it as a transparent, agency-prioritizing option rather than a dominant rule.

Key Output. Mechanism metrics, flip probabilities, and Pareto comparisons.

7 Model C: What Drives Success? (Judges vs Fans)

Takeaway. Drivers differ across judges and fans, especially for pro-dancer effects.

We fit mixed-effects models on logit shares:

$$\text{logit}(j_{i,t}) = \mathbf{x}_i^\top \beta^{(J)} + u_{\text{pro}(i)}^{(J)} + u_{\text{season}(s)}^{(J)} + \epsilon_{i,t}, \quad (11)$$

$$\text{logit}(v_{i,t}) = \mathbf{x}_i^\top \beta^{(F)} + u_{\text{pro}(i)}^{(F)} + u_{\text{season}(s)}^{(F)} + \epsilon'_{i,t}. \quad (12)$$

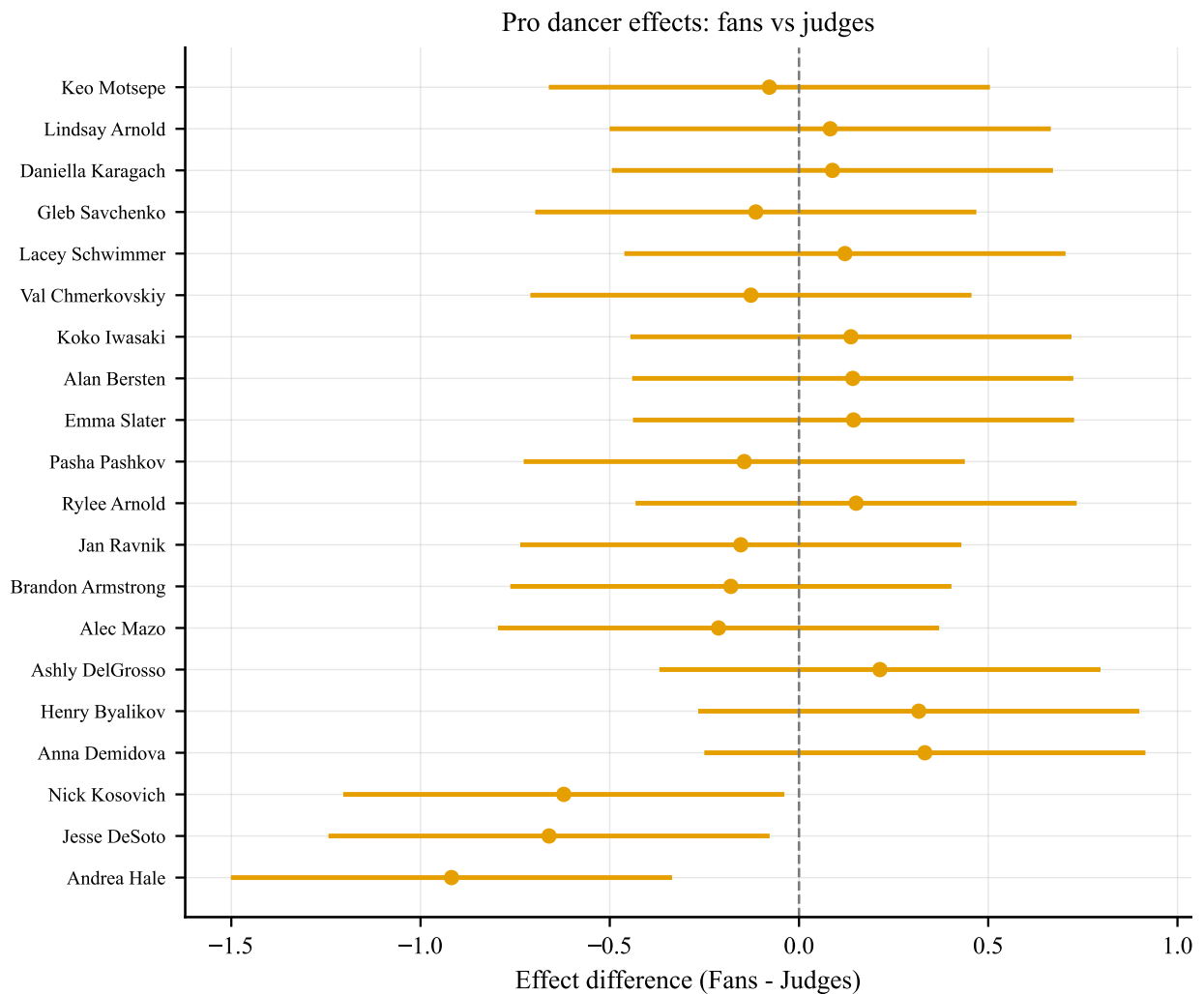


Figure 13: Pro dancer effects (fans minus judges).

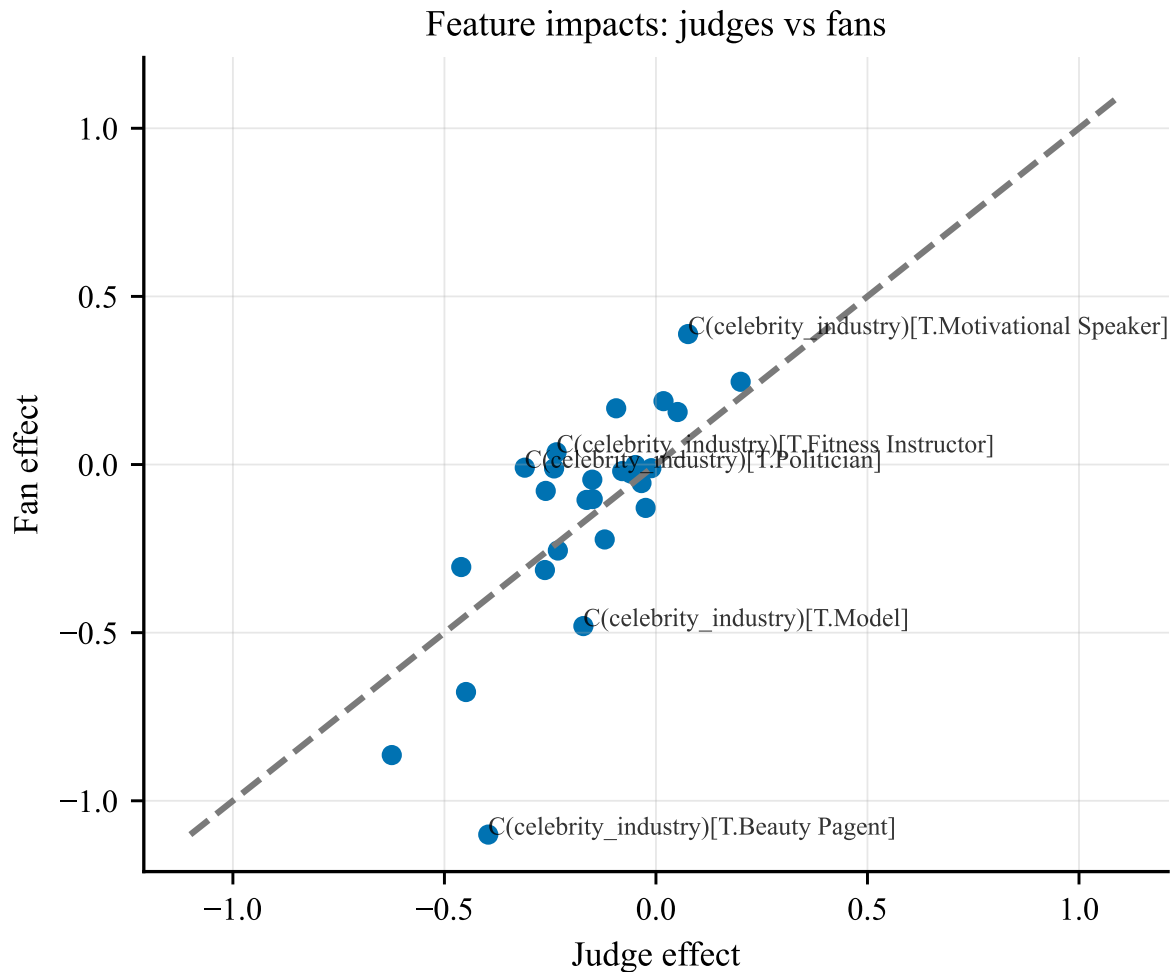


Figure 14: Annotated outliers highlight features with the largest judge-fan gaps.

7.1 Predictive Add-on: GBDT

Takeaway. We include a predictive model as a robustness check, not as the main driver analysis.

We train a gradient-boosted decision tree (GBDT) classifier to predict elimination using forward-chaining validation and report AUC as a sanity check on covariate relevance. Predictive performance is stable and supports the selected covariates; see Appendix B for the calibration curve.

Key Output. Dual models and a direct answer to Task 3: effects are not identical.

8 Model D: Mechanism Design (DAWS)

Takeaway. DAWS uses thresholded weights tied to uncertainty, making the rule transparent and executable.

We define

$$\alpha_t = \begin{cases} \alpha_0, & U_t \leq q_{90}, \\ \alpha_1, & q_{90} < U_t \leq q_{97}, \\ \alpha_2, & U_t > q_{97}, \end{cases} \quad (13)$$

where q_{90}, q_{97} are quantiles of U_t and $(\alpha_0, \alpha_1, \alpha_2)$ are fixed public tiers. Figure 15 shows the tiered trigger schedule implied by these thresholds.

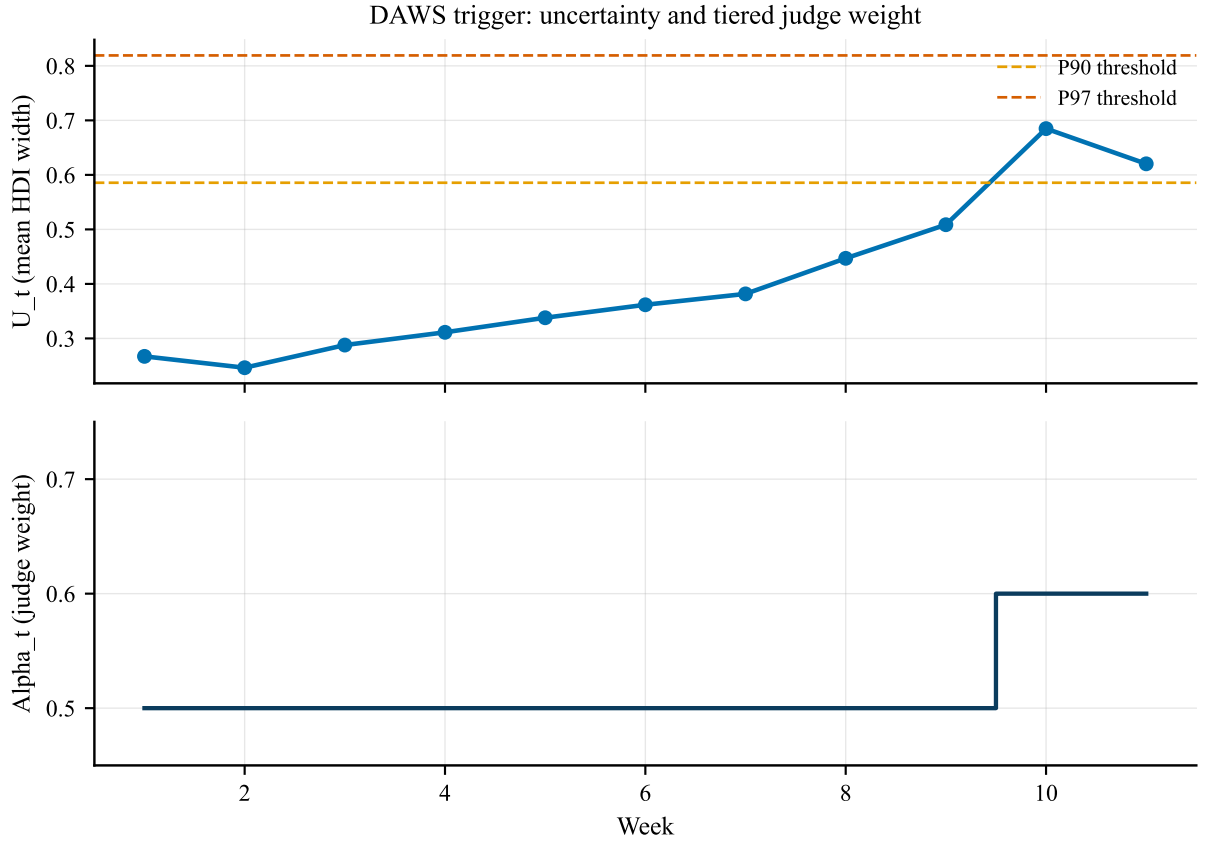


Figure 15: DAWS trigger schedule: weekly uncertainty U_t with public thresholds and the resulting tiered judge weights.

Proposition 1 (Monotonicity). *If both judge share and fan share of a contestant increase, their DAWS score does not decrease.*

Proposition 2 (Stability bound). *With tiered α_t , the score change satisfies*

$$|c_{i,t} - c_{i,t-1}| \leq |\alpha_t - \alpha_{t-1}| |j_{i,t} - v_{i,t}| + (1 - \alpha_t) \|\mathbf{v}_t - \mathbf{v}_{t-1}\| + \alpha_t \|\mathbf{j}_t - \mathbf{j}_{t-1}\|. \quad (14)$$

8.1 Judge-save parameter calibration

We use a calibrated β in

$$\Pr(E = a \mid \{a, b\}) = \sigma(\beta(J_b - J_a)) \quad (15)$$

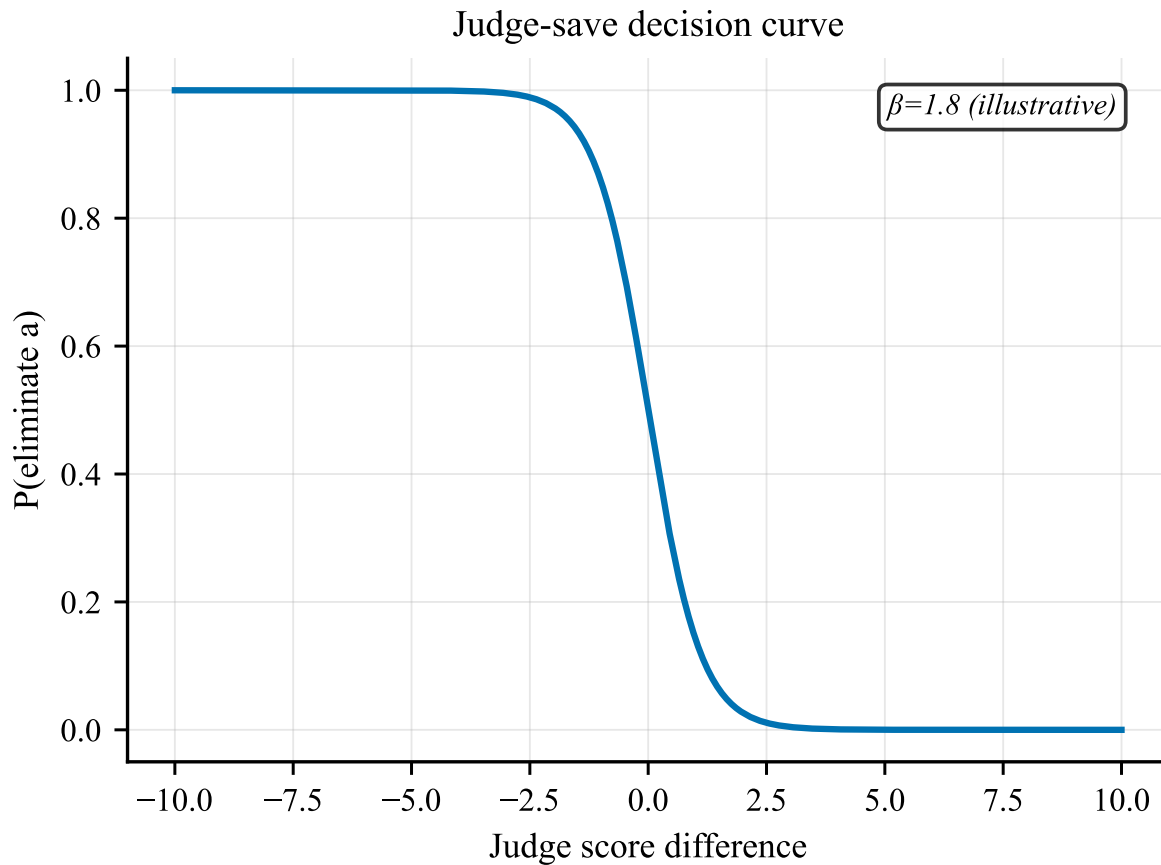


Figure 16: Judges prefer higher score within the bottom two; the curve uses calibrated $\beta = 1.8$ to illustrate sharpness.

Key Output. DAWS schedule, properties, and calibrated judge-save behavior.

9 Sensitivity and Validation

Takeaway. Key claims are stable to σ , ϵ , and rule-switch priors.

We vary σ (smoothness), ϵ (vote floor), and ρ (switch probability). Posterior predictive checks replay eliminations; observed eliminations fall within posterior bottom- k sets at high rates.

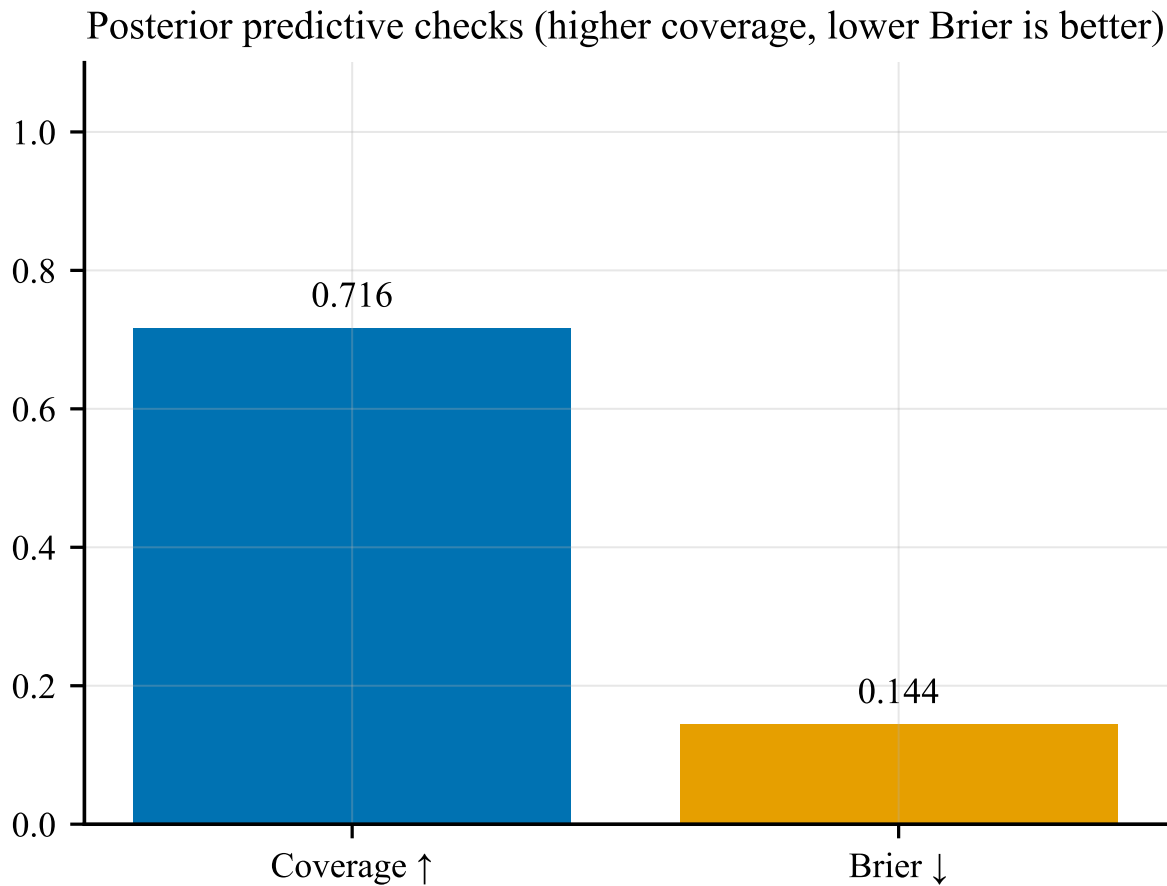


Figure 17: Model reproduces eliminations while preserving uncertainty.

9.1 Scale Benchmark

We benchmark sampling scale with a multi-process setup and record runtime, error (mean HDI width), stability (DAWS), and theory-fit (Kendall τ). The curves show diminishing returns in uncertainty reduction beyond mid-scale settings; the elbow (dashed line) marks our final scale choice.

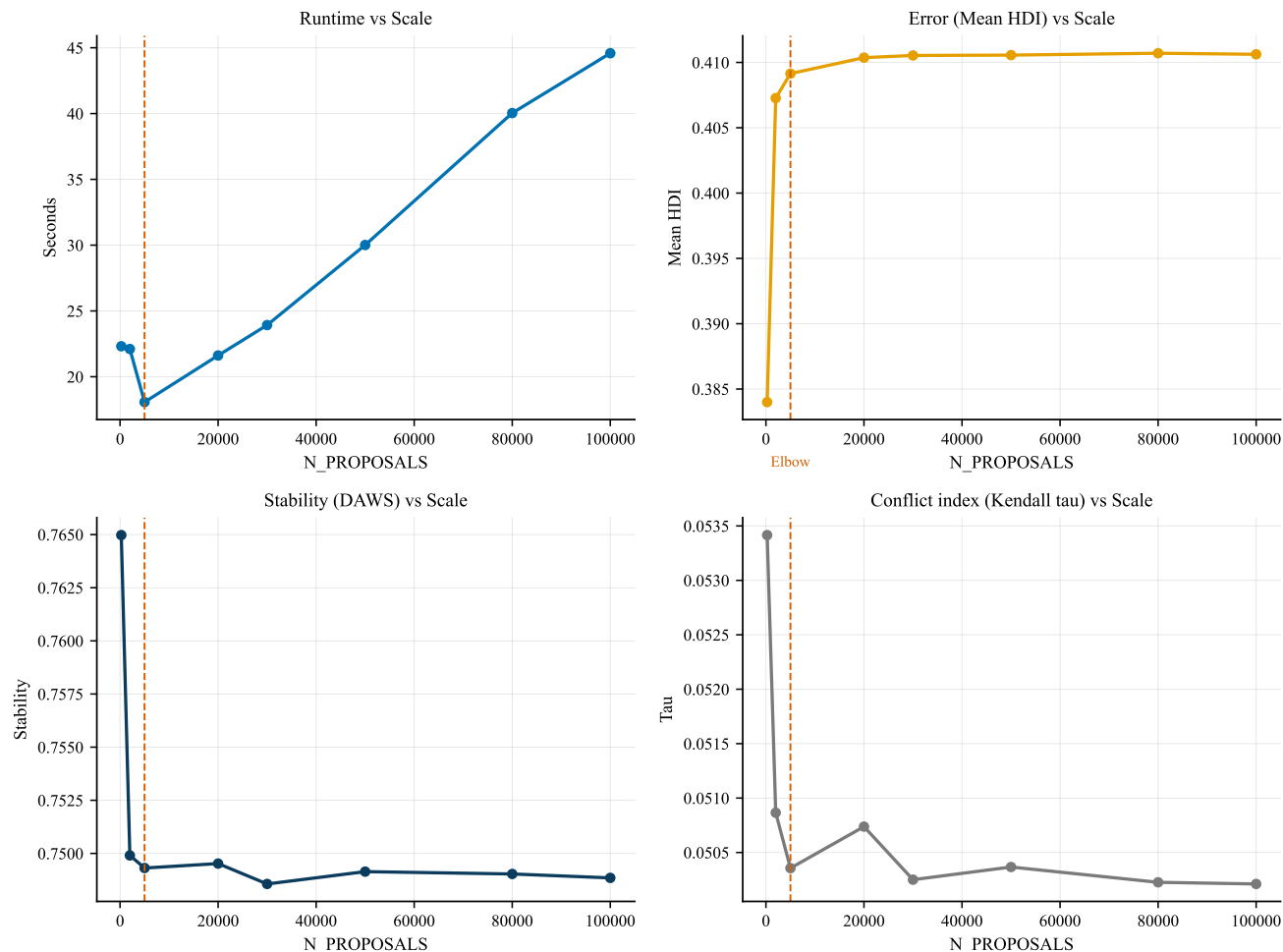


Figure 18: Scale benchmark across $N_{\text{proposals}}$ with runtime, error, stability, and theory-fit.

Key Output. Sensitivity curves and posterior predictive validity metrics.

10 Conclusions and Recommendations

Takeaway. Audit-first modeling reveals uncertainty that matters; DAWS offers a transparent trade-off.

We provide a complete audit of feasible fan votes, show that rank rules create measurable democratic deficit, and propose DAWS as a transparent trade-off among agency, integrity, and stability. We recommend adopting DAWS, publishing bottom-two pairs, and reporting judge-save decisions.

- **Decision-ready summary:** Uncertainty is concentrated in a small set of weeks; most weeks are identifiable.
- **Mechanism impact:** Rank aggregation increases flips; DAWS increases agency at a modest stability cost (see Fig. 11 and Fig. 15).

- **Actionability:** Publish a DAWS schedule and judge-save criteria to improve transparency.

A Sensitivity Analysis

We present the smoothness parameter sensitivity analysis here. Key conclusions remain stable across a range of σ values.

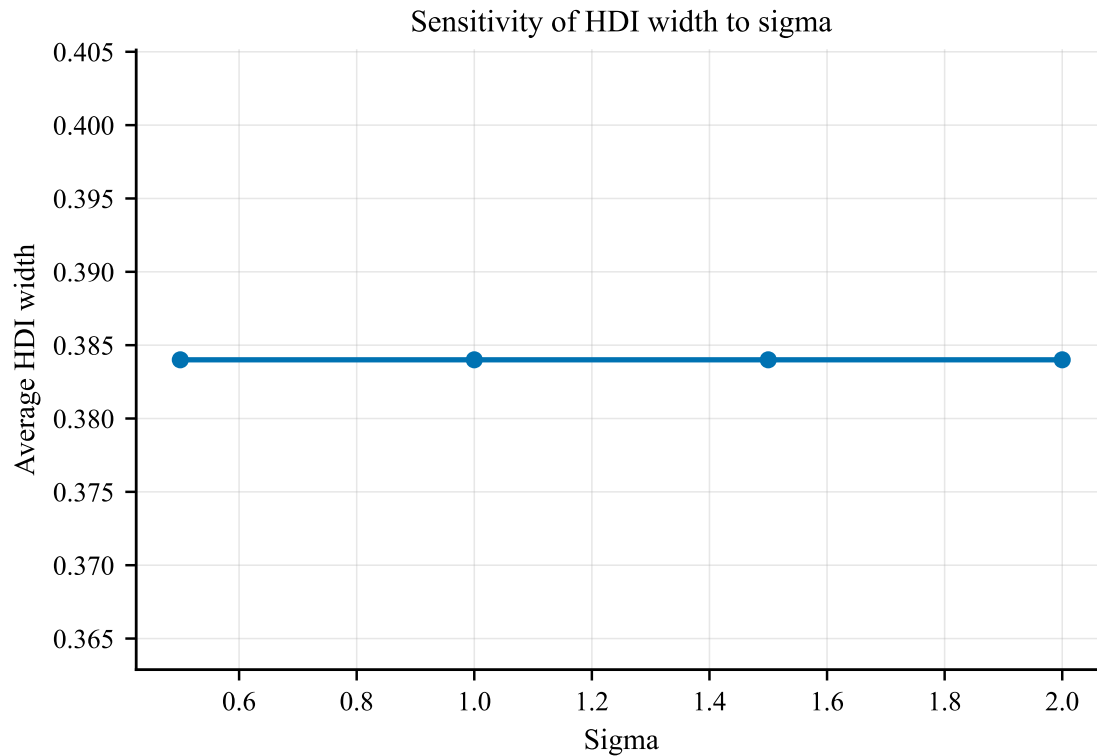


Figure 19: Sensitivity of key metrics to smoothness parameter σ . Conclusions are robust across the tested range.

B Predictive Calibration

We include forward-chaining AUC results as a robustness check on covariate relevance.

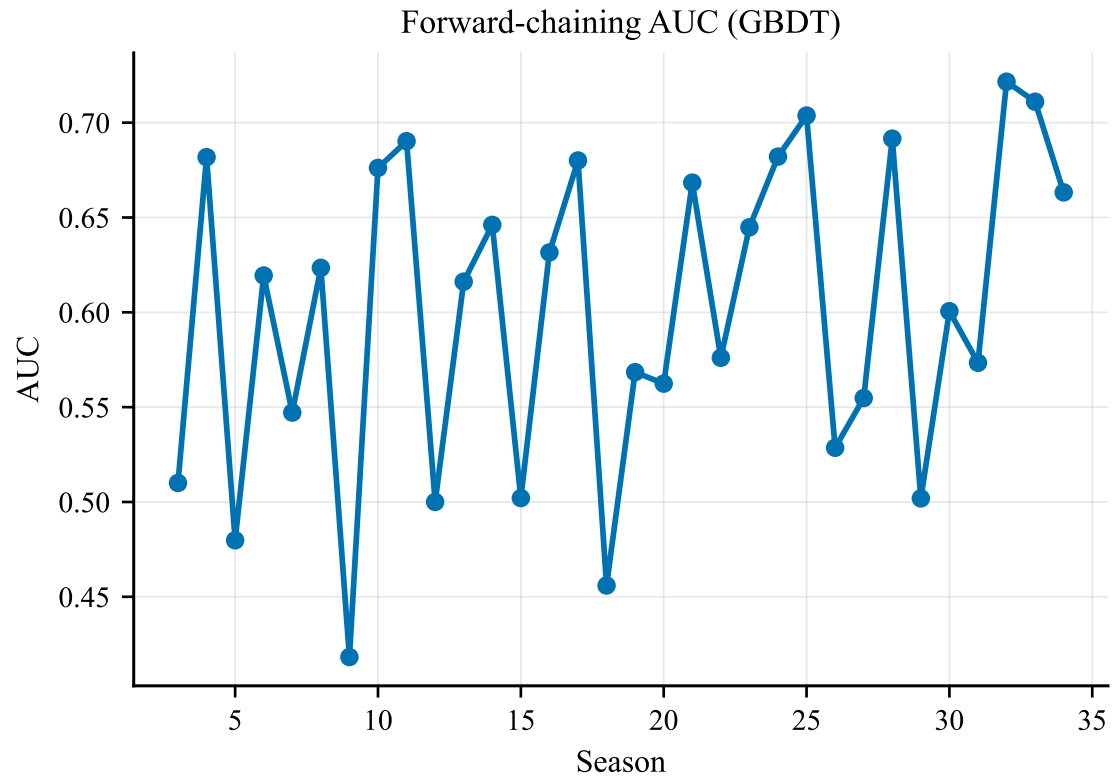


Figure 20: Forward-chaining AUC curve. Predictive performance is stable and supports the selected covariates.

References

- [1] COMAP. 2026 MCM/ICM Problem C: Dancing with the Stars (DWTS). Contest Problem Statement.
- [2] Smith, R. (1984). Efficient Monte Carlo procedures for generating points uniformly in polytopes. *Operations Research*.
- [3] Jaynes, E. T. (1957). Information theory and statistical mechanics. *Physical Review*.
- [4] Gelman, A., et al. (2013). *Bayesian Data Analysis*. CRC Press.
- [5] Moulin, H. (1988). *Axioms of Cooperative Decision Making*. Cambridge Univ. Press.

AI Use Report

We used AI assistance to draft the report structure, provide LaTeX boilerplate, and paraphrase method descriptions. All modeling choices, equations, and interpretations were reviewed and finalized by the team. No external data beyond the provided contest dataset were used.

- Reproducibility: code, figures, and metrics are generated from the provided dataset.
- Environment: Miniforge + mcm2026 with pinned scientific stack.
- Audit trail: pipeline logs and summary metrics are saved for each run.