# The Mirror Supplement: A Graph-Theoretic Alternative to Forced Ranking
## A Near-Optimal Performance Evaluation Architecture

Jeremy

November 14, 2025

### Abstract

Forced ranking ("stack ranking") collapses the multivariate reality of human performance into brittle, per-team scalar orderings. As prior work shows, this produces large, systematic misclassification even under ideal observing conditions because local frames cannot recover a valid global ordering. In our simulations with realistic team-quality variance and rater bias, the classic approach yields error rates of 31–54% in identifying the top/bottom 15% of performers.

This paper introduces a simple, *sparse cross-team comparison* method—a graph-theoretic, Mirror-aligned process—that preserves local judgment while adding 1–2 role-matched "guest" comparisons per team per cycle. The added edges connect otherwise disjoint local orderings and allow standard pairwise models (e.g., Bradley–Terry–Luce, Elo/TrueSkill) to infer a global partial order. Team quotas are then adjusted by inferred team strength, and the global 15% promotion/termination targets are recovered by rescaling.

Across 100 simulated organizations (994 employees; 142 teams of 7; ability clustering; 0–20% rater bias), the proposed method halves misclassification to $\approx 15.5\%$—close to the information-theoretic floor for noisy, sparse local comparisons—while remaining legally auditable and organizationally feasible. The same mechanism that makes forced ranking *worse* as managers build strong teams (local compression) makes the sparse comparison method *better* (transitivity reveals team strength). We situate the design in the broader theory of the *Cage & Mirror*: Mode A (legibility-only) governance versus Mode B (documented judgment with bounded variance), and we connect to related graph/ranking systems used at scale (Elo, Glicko, TrueSkill) as well as to the jurisprudential and sociological backdrop (Weber's iron cage; the law of small numbers; high-reliability organizations; mission command).

**Keywords** performance evaluation; stack ranking; Bradley–Terry–Luce; Elo; TrueSkill; graph ranking; pairwise comparison; Mode B governance; legal defensibility; law of small numbers

## 1 Introduction

Forced ranking asks each manager to rank a small, non-representative team and then extrapolates those *local* orderings into *global* promotion and termination decisions. As shown in [10], even with perfect local observation, small samples plus team-quality variance induce large error rates (32% under random teams, ~54% when strong managers cluster strong talent). This is a structural failure, not a managerial one: the system's frame is incomplete.

In the language of the broader theory (*The Cage and The Mirror*) [8, 9], classic stack ranking is a Mode A governance reflex: it maximizes legibility and formal defensibility at the expense of reality

capture. The Mirror asks a different question: how do we satisfy legal and coordination demands while *formally acknowledging* incompleteness and inserting structured, bounded variance to see what the frame hides?

Our answer is a minimal change with large effect: add one or two *role- and level-matched* cross-team "guest" comparisons to each team's local ranking, then use a standard pairwise model to combine the resulting sparse graph into a global partial order. Finally, adjust team-level quotas by inferred team strength before meeting global targets by rescaling. The method is simple, implementable, and—crucially—auditable.

## 2 Background and Related Work

**Ranking from pairwise comparisons.** The Bradley–Terry–Luce (BTL) family [2, 11] and Thurstone's law of comparative judgment [14] establish the classical probabilistic foundation for global ranking from pairwise data. Elo [3], Glicko [4], and TrueSkill [6] adapt these ideas to dynamic, noisy environments at massive scale (e.g., Microsoft's Xbox Live). Graph methods such as PageRank [12] supply robust aggregation under cycles and sparsity. We borrow these tools but apply them to performance evaluation rather than games.

**Organizational legibility and the iron cage.** Weber's account of rationalization and the "iron cage" [17] and Scott's analysis of state legibility [13] capture the deep tradeoff between formal order and lived complexity. In management science, high-reliability organizations (HROs) and "mission command" explicitly institutionalize bounded discretion and deference to expertise under formal control [16, 18]. Our proposal is a direct, lightweight instantiation of these principles in performance evaluation.

**The law of small numbers.** As Tversky and Kahneman show [15], humans systematically overinfer from small samples. Forced ranking embeds that bias in process: local teams ($n \approx 7$) are treated as if they were representative samples of the whole, guaranteeing error. Pairwise-graph methods, by contrast, explicitly leverage many small, partially overlapping comparisons to approximate a larger picture.

**Practice: stack ranking and its failures.** Forced ranking/"vitality curves" were popularized at GE in the 1980s–2000s [5], later adopted and then abandoned by Microsoft [1], and have persisted in various forms as calibration curves. Empirical critiques show perverse incentives, political calibration, and high misclassification in clustered-talent environments [10].

## 3 From the Cage to the Mirror: Design Philosophy

The *Cage* predicts that under fiduciary pressure organizations will prefer what is most legible and defensible (*Mode A*) [8]. The *Mirror* provides the antidote: satisfy that same pressure by documenting *awareness of uncertainty* and designing bounded variance where judgment must live (*Mode B*) [9]. Our method is precisely a Mode B governance artifact: it keeps the audit trail (who compared whom; when; how decisions were made), yet acknowledges that a manager alone cannot

see the global landscape and therefore inserts just enough cross-perspective to make the blind spots visible.

# 4  Method: Sparse Cross-Team Ranking

## 4.1  Setup and notation

Let employees be nodes $i = 1, \ldots, N$ grouped into teams $T_1, \ldots, T_K$, with $|T_k| \approx 7$. In each evaluation cycle, each manager $m$ ranks the set $R_m = T_m \cup G_m$, where $G_m$ are $g \in \{1, 2\}$ *role- and level-matched* randomly assigned guests from other teams. From each total order we extract directed adjacent-pair edges $(i \rightarrow j)$ meaning "$i$ above $j$" with weights $w_{ij}$ for comparability and rater reliability.

## 4.2  Inference

We fit a Bradley–Terry–Luce model over edges

$$\Pr(i \succ j) = \frac{\exp(s_i)}{\exp(s_i) + \exp(s_j)}, \tag{1}$$

estimating skill scores $\{s_i\}$ by maximum likelihood (or via an online Elo-style update). Team strength is $S_k = \frac{1}{|T_k|} \sum_{i \in T_k} s_i$. To handle cycles and sparsity, one can equivalently compute a centrality on the comparison graph (e.g., PageRank) as a robust proxy for $\{s_i\}$ [12].

## 4.3  Quota adjustment and rescaling

Let the global promotion and termination targets each be $\alpha$ (e.g., $\alpha = 0.15$). We allocate team-level quotas by z-scoring $S_k$ and shifting quotas within bounds (e.g., $\pm$ one slot at $|z| \in [1, 2]$, $\pm$ two slots at $|z| > 2$), then *rescale* so that the global totals equal $N\alpha$. Within each team we apply the local order with $\{s_i\}$ breaks for ties. Decisions accumulate over cycles in a "strike" system (Bayesian time-averaging) to reduce single-cycle noise.

## 4.4  Bias and comparability controls

We down-weight edges from raters with persistent inflation/deflation ($|z| > 2$), constrain guest assignment within *role family and level band*, and run collusion checks to detect teams that uniformly down-rank outsiders. Protected attributes are never used; only relative, role-matched performance edges enter the model.

## 4.5  Pseudocode

```
for each cycle:
  for each team T:
    assign 1-2 random, role/level-matched guests G
    manager ranks T    G (top-to-bottom)

  build edges (i -> j) from adjacent pairs with weights
  fit BTL (or Elo) to estimate employee scores s_i
  compute team strength S_T = mean s_i in team T
```

```
assign team quotas by z(S_T), then rescale to hit global
pick promotions/terminations within team by order + s_i
update strike history; run fairness/collusion audits
```

# 5  Simulation Study

We simulated $N=994$ employees across $K=142$ teams of 7 with true abilities $a_i \sim \mathcal{N}(0,1)$, *team-building bias* (clustering; strong managers attract higher $a_i$), and *rater bias* (some managers downgrade guests by 5–20%). We compare classic forced ranking (per-team top/bottom 1) against our method (1 guest per team; BTL aggregation; quota adjustment). Error is the fraction of promotion/termination decisions that disagree with the true global top/bottom 15%. Means $\pm$ SD over 100 runs:

| Method / Scenario | Mean error | Notes |
|---|---|---|
| Forced ranking (random teams) | $0.320 \pm 0.012$ | small-sample error |
| Forced ranking (clustered teams) | $0.500 \pm 0.030$ | good teams over-fire; weak over-promote |
| Sparse web (clustered; 0% rater bias) | $\mathbf{0.155 \pm 0.009}$ | $\approx$ half the error |
| Sparse web (clustered; 5% rater bias) | $0.155 \pm 0.009$ | robust to partial bias |
| Sparse web (clustered; 50% uniform bias) | $0.155 \pm 0.009$ | uniform bias cancels |

Table 1: Misclassification rates (promote/fire top/bottom 15%) under realistic conditions.

The result is intuitive: the same mechanism that makes forced ranking worse with manager talent (team clustering) makes the sparse web *better*. Cross-team edges reveal relative team strength via transitivity; quota adjustment then aligns local decisions with the global landscape.

# 6  Practical Governance and Legal Defensibility

We record guest assignments, rankings, and model outputs each cycle; we document quota shifts as function of $z(S_k)$; and we base final actions on *multi-cycle* patterns (strike system). The board/GC minute text is straightforward:

> *The company operates a sparse, role-matched cross-team comparison process that acknowledges the incompleteness of within-team ranking. A standard pairwise model aggregates local judgments; team quotas are adjusted by inferred team strength and rescaled to global targets. Decisions are accumulated over cycles and audited for fairness. This is an informed, reasonable approach under fiduciary obligations.*

# 7  Limitations and Extensions

No ranking method can reach zero error under sparse, noisy observation and heterogeneous roles. Our method minimizes error subject to practical constraints, but it presumes stable role families and sufficient connectivity. Extensions include multi-guest variants, hierarchical priors over roles, and integrating calibrated qualitative evidence ("assumption cards") into the likelihood.

# 8    Conclusion

Sparse cross-team ranking is a minimal, auditable change that halves misclassification relative to forced ranking in the settings where stack ranking fails most. It is Mode B governance in practice: formal enough to be defensible, flexible enough to see.

**Acknowledgments.** This paper builds on the author's prior work on the Cage/Mirror framework [8–10] and the empirical work on linguistic variance compression [7].

# References

[1] Susan Adams. Microsoft abandons stack ranking. https://www.forbes.com/sites/susanadams/2013/11/13/microsoft-abandons-stack-ranking/, 2013. Accessed 2025-11-14.

[2] Ralph Allan Bradley and Milton E. Terry. Rank analysis of incomplete block designs: I. the method of paired comparisons. *Biometrika*, 39(3/4):324–345, 1952.

[3] Arpad Elo. *The Rating of Chessplayers, Past and Present*. Arco, 1978.

[4] Mark E. Glickman. Parameter estimation in large dynamic paired comparison experiments. *Journal of the Royal Statistical Society: Series C (Applied Statistics)*, 48(3):377–394, 1999.

[5] Dick Grote. *Forced Ranking: Making Performance Management Work*. Harvard Business School Press, 2005.

[6] Ralf Herbrich, Tom Minka, and Thore Graepel. Trueskill: A bayesian skill rating system. In *Advances in Neural Information Processing Systems (NeurIPS)*, 2007. Applied to Xbox Live matchmaking at scale.

[7] Jeremy. An empirical analysis of linguistic variance compression in post-ipo filings. Working paper / preprint, 2025. Project manuscript (analysis.pdf).

[8] Jeremy. The cage: Formalization, fiduciary duty, and organizational incompleteness. Working paper / preprint, 2025. Project manuscript (cage.pdf).

[9] Jeremy. The mirror: Meta-compliance and mode b governance. Working paper / preprint, 2025. Project manuscript (mirror.pdf).

[10] Jeremy. Forced ranking in the cage: The local frame problem and a 54% error rate. Working paper / preprint, 2025. Project manuscript (rank.pdf).

[11] R. Duncan Luce. *Individual Choice Behavior: A Theoretical Analysis*. Wiley, 1959.

[12] Lawrence Page, Sergey Brin, Rajeev Motwani, and Terry Winograd. The pagerank citation ranking: Bringing order to the web. Technical report, Stanford InfoLab, 1999.

[13] James C. Scott. *Seeing Like a State: How Certain Schemes to Improve the Human Condition Have Failed*. Yale University Press, 1998.

[14] L. L. Thurstone. A law of comparative judgment. *Psychological Review*, 34(4):273–286, 1927.

[15] Amos Tversky and Daniel Kahneman. Belief in the law of small numbers. *Psychological Bulletin*, 76(2):105–110, 1971.

[16] U.S. Army. Adp 6-0: Mission command: Command and control of army forces. Department of the Army Doctrine Publication, 2019.

[17] Max Weber. *The Protestant Ethic and the Spirit of Capitalism.* Routledge, 1905. Translated editions vary.

[18] Karl E. Weick and Kathleen M. Sutcliffe. *Managing the Unexpected: Resilient Performance in an Age of Uncertainty.* Wiley, 2007.