

Reward Bias Substitution: Single-Axis Mitigations Shift Optimization Pressure

MAX LAMPARTH, Stanford University, USA

DANIEL FEIN, Stanford University, USA

ANDREAS HAUPT, Stanford University, USA

MARCEL HUSSING, University of Pennsylvania, USA

MYKEL KOCHENDERFER, Stanford University, USA

Single-axis mitigations of reward-model biases (e.g., reducing reliance of the proxy reward on length, sycophancy, or style) can rotate optimization pressure onto correlated proxies rather than eliminate it, a failure mode we call reward bias substitution. We formalize the underlying measurement-vs-optimization gap between the audit distributions where mitigations are validated and the policy distributions where optimization realizes their effects. We introduce a taxonomy, instantiated in closed form, classifying single-axis mitigation outcomes into successful mitigation, bias substitution, overcorrection, silent non-op, and audit-distribution sensitivity. We prove that single-axis mitigation methods cannot be validated by audit-distribution-only evaluation: successful mitigation, bias substitution, and overcorrection produce structurally identical observables under ranking accuracy and win-rate scoring, regardless of benchmarks richness. Augmenting evaluation with policy-induced distributions while tracking multiple biases provably closes the gap. We give actionable prescriptions for mitigation methods and benchmarks. Across published preference-learning mitigation work, no method we survey reports the evidence needed to certify successful mitigation. We demonstrate bias substitution in language model RLHF, where a length penalty during GRPO training compresses responses as intended yet redirects optimization pressure onto confidence calibration, driving the trained policy into overconfidence. Our experiments also show that a published length-debiasing operator zeros pooled reward-length correlation but flips sign within-prompt on three of four SOTA reward models with true reward degrading on two, and that length-sycophancy coupling reverses under human-LLM judge disagreement across eight model families.

CCS Concepts: • **Computing methodologies** → **Natural language processing; Reinforcement learning.**

Additional Key Words and Phrases: RLHF, reward hacking, preference learning, evaluations, reward modeling, theory

ACM Reference Format:

Max Lamparth, Daniel Fein, Andreas Haupt, Marcel Hussing, and Mykel Kochenderfer. 2026. Reward Bias Substitution: Single-Axis Mitigations Shift Optimization Pressure. 1, 1 (May 2026), 52 pages. <https://doi.org/10.1145/nnnnnnnn.nnnnnnnn>

1 Introduction

Reinforcement Learning from Human Feedback (RLHF) and related preference-learning methods rely on human annotations to shape model behavior on non-verifiable objectives, but the resulting reward models and policies are confounded by spurious correlations in either the preferences or their actions. Both cases are examples of reward hacking as caused by the reward identifiability problem [5, 37, 72, 73, 78], leading to misgeneralization in terms of

Authors' Contact Information: Max Lamparth, lamparth@stanford.edu, Stanford University, Stanford, California, USA; Daniel Fein, Stanford University, Stanford, California, USA; Andreas Haupt, Stanford University, Stanford, California, USA; Marcel Hussing, University of Pennsylvania, Philadelphia, Pennsylvania, USA; Mykel Kochenderfer, Stanford University, Stanford, California, USA.

Permission to make digital or hard copies of all or part of this work for personal or classroom use is granted without fee provided that copies are not made or distributed for profit or commercial advantage and that copies bear this notice and the full citation on the first page. Copyrights for components of this work owned by others than the author(s) must be honored. Abstracting with credit is permitted. To copy otherwise, or republish, to post on servers or to redistribute to lists, requires prior specific permission and/or a fee. Request permissions from permissions@acm.org.

© 2026 Copyright held by the owner/author(s). Publication rights licensed to ACM.

Manuscript submitted to ACM

Manuscript submitted to ACM

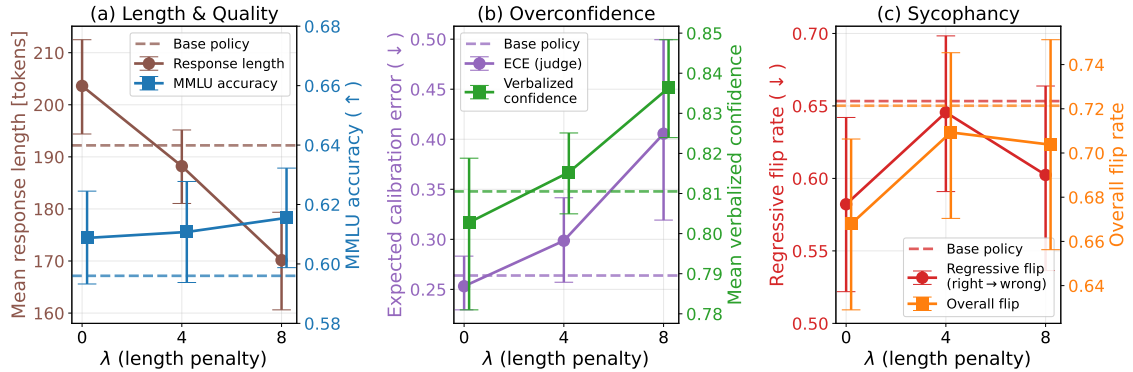


Fig. 1. RLHF of *Llama-3.2-3B-Instruct* with *Skywork-Reward-V2-Llama-3.1-8B* using GRPO ($\beta = 2e-2$, $LR = 3e-5$, 600 steps) on *UltraFeedback* with four random seeds for each $\lambda \in \{0, 4, 8\}$, 95% bootstrapped confidence intervals. To mitigate length, we modify the reward \hat{R} in the RLHF objective as $\hat{R}(x, y) = R_{RM}(x, y) - \lambda n_{\text{tok}}/100$. Training with $\lambda = 0$ does not break confidence calibration. As length is decreased the optimization pressure rotates onto the calibration axis, while MMLU accuracy and sycophantic behavior are preserved. See Section F.1 for details.

the original behavioral-shaping intention. A recurring example in language RLHF is reward-length correlation [e.g. 59, 71] of language reward models and verbosity of language model policies in RHLF and DPO methods [14, 33, 42, 64]. The so-called *length bias* issue has been observed with different proposed mitigation strategies, however, most prior reward bias mitigation works do not fully resolve or *even overcorrect* [21], consider at most two confounding factors simultaneously [21], or evaluate within fixed audit distributions that leave substitution onto uncovered axes undetectable [22, 86]. Naively, one would want the proxy reward to track the true reward and be invariant to irrelevant surface features (length, formatting, style), varying only with the causally relevant features annotators care about (correctness, helpfulness, informativeness), motivating *reward feature invariance*. However, surface features labeled spurious are typically partially informative about the true reward and coupled with each other, so enforcing strict invariance removes signal alongside bias. As we demonstrate in this paper, ignoring these couplings produces **reward bias substitution**, where single-axis mitigations (e.g., debiasing length) tend to rotate optimization pressure into other correlated proxies rather than eliminating it. Any single-axis mitigation asserts that the targeted feature is spurious, so bias substitution can arise whether or not that assertion is identifiable from preference data. Figure 1 shows this issue concretely. Applying a length penalty during GRPO training compresses response length as intended, yet the freed optimization pressure drives the policy into the untargeted overconfidence axis, while multiple-choice accuracy holds and free-form accuracy decreases (see also Section F.1). As a result, we show that almost all existing reward bias mitigation papers cannot determine whether they have successfully mitigated, overcorrected, or merely substituted reliance onto an uncovered proxy on the evidence they report. Current reward/preference benchmarks cannot adjudicate between these outcomes either, risking unintended post-deployment harms of, e.g., silently substituting length bias for sycophancy bias in user interactions. We address these gaps by formalizing the failure mode and characterizing when standard evaluation can detect it. **This paper contributes the following:**

- **We formalize bias substitution as failure mode of single-axis mitigation**, in which removing spurious reliance at the evaluation distribution where mitigations are validated leaves optimization pressure free to

rotate onto correlated proxies at the trained policy-induced distribution. We identify the *measurement-versus-optimization gap* as the blindspot in prior reward-hacking definitions (see Section C).

- **We introduce a regime taxonomy** classifying all mechanistically distinct single-axis mitigation outcomes: successful mitigation (R0, with a contaminated sub-case R0_{cont}), bias substitution (R1), overcorrection (R2), silent non-op (R3), and audit-distribution sensitivity (R4). Each regime is instantiated in closed form and validated empirically, with additional deployment-empirical instances for R1, R2, and R4.
- **We prove a matched impossibility–sufficiency pair for audit-distribution evaluation.** No benchmark functional over audit-distribution observables can reliably distinguish successful mitigation from bias substitution or overcorrection, but, augmenting the benchmark input with policy-induced distributions provably does. We derive necessary prescriptions for *certifiable reward bias mitigation claims* and benchmark improvements.
- **We provide empirical evidence for our framework.** We demonstrate bias substitution directly in language model RLHF, where a GRPO length penalty trades response length for policy overconfidence (R1). We further instantiate the measurement-versus-optimization gap on a length-debiasing operator across five reward models. We provide the first systematic characterization of length-sycophancy statistical dependence across four labeling regimes and eight model families, finding that the coupling reverses sign under human–LLM judge disagreement (R4). Beyond our experiments, we identify previously isolated findings as instances of R1 and R2 across preference-learning mitigation work.

2 RLHF Invariances and Formalizing Bias Substitution

We work in the single-turn contextual-bandit reduction of RLHF [14, 42, 76] with prompts $x \in \mathcal{X}$ drawn i.i.d. from a fixed context distribution \mathcal{D} , responses $y \in \mathcal{Y}$ sampled $y \sim \pi(\cdot | x)$, true reward $R(x, y)$, learned proxy $\tilde{R}(x, y)$, and the *plain return* $J(\pi, R) = \mathbb{E}_{x \sim \mathcal{D}, y \sim \pi(\cdot | x)} [R(x, y)]$ to maximize. The KL-regularized RLHF training objective is

$$J_{\text{RLHF}}(\pi; \tilde{R}) = \mathbb{E}_{x \sim \mathcal{D}} [\mathbb{E}_{y \sim \pi(\cdot | x)} [\tilde{R}(x, y)] - \beta D_{\text{KL}}(\pi(\cdot | x) \| \pi_{\text{ref}}(\cdot | x))]. \quad (1)$$

For any reward \tilde{R} the KL-regularized optimum is the softmax policy [63, 64]

$$\pi_{\beta}^{\star}(\tilde{R})(y | x) \propto \pi_{\text{ref}}(y | x) \exp(\tilde{R}(x, y) / \beta). \quad (2)$$

This optimal KL regularized policy is not invariant to reward scale. Let $\mu_{\pi}(x, y) = \mathcal{D}(x)\pi(y | x)$ denote the policy-induced distribution on $\mathcal{X} \times \mathcal{Y}$ given a policy π (the one-step occupancy measure under \mathcal{D}). In addition to the KL anchor π_{ref} -defined measure $\mu_{\pi_{\text{ref}}}$, we fix a *diagnostic measure* μ_{diag} on $\mathcal{X} \times \mathcal{Y}$, used for reliance estimation and correlation measurement. Natural choices include $\mu_{\pi_{\text{ref}}}$ itself (coupling measurement and optimization references) and annotator-conditioned audit distributions $\mu_{\text{diag}}^{\text{human}}, \mu_{\text{diag}}^{\text{LLM}}$, reflecting that reliance estimates depend on who labels and formalizing the human-versus-LLM-judge gap (see also Section E.3). A standing regularity Assumption E.1 is invoked throughout for well-defined softmax policies and policy-level expectations. Reward modeling in RLHF involves invariances at two levels. First, what preference data identifies about the reward (data-level), and second, how the KL-regularized policy responds to reward transformations (policy-level). At the data level, the *reward identifiability problem* leaves the learned reward underdetermined. Many reward functions explain the same preference data and preference comparisons identify the true reward only up to a prompt-only additive shift under the Bradley-Terry likelihood, leaving the cardinal scale loosely pinned [72]. This identifiability gap allows the reward model to absorb spurious feature correlations consistent with the observed rankings. At the policy level, KL-regularized optimization is sensitive to cardinal reward values, since $\tilde{R} \rightarrow c\tilde{R}$ at fixed β is equivalent to $\beta \rightarrow \beta/c$, so two reward models agreeing on all pairwise preferences can still induce different

157 policies. Combined with data-level identifiability, this means confounded features can be consistent with the observed
 158 rankings while still distorting the cardinal values that determine the optimized policy. For example, applied to length
 159 bias, this implies that preference data alone cannot distinguish causally relevant length from spurious length without
 160 additional assumptions or interventions [72]. Partial identifiability of reward under Bradley-Terry has been formally
 161 characterized [72] and reward misidentification empirically demonstrated on continuous-control benchmarks [78]. Both
 162 concern what preference data fails to identify. The downstream question of how single-axis mitigation operators that
 163 target identified spurious features behave under optimization has been observed empirically in supervised vision [48],
 164 but not formally characterized in any setting, and its consequences for reward-model evaluation, mitigation methods,
 165 and downstream policy optimization remain unaddressed for RLHF. We work within this identifiability gap, classifying
 166 mitigation outcomes by quantities that remain meaningful when the spurious/structural partition and cardinal scale are
 167 underdetermined.

171 Thus, we decompose the proxy reward \tilde{R} along a fixed set of interpretable surface features Φ (length, sycophancy,
 172 formatting), partitioned at an audit distribution μ_{diag} into spurious features Φ_{sp} and structurally relevant features Φ_{struct}
 173 by whether the true reward R depends on them (see Section A for the full formal derivation). A single-axis mitigation M_i
 174 is a projection-type operator, abstracting deployed linear-probe, calibration, and disentanglement methods, that zeros
 175 the proxy’s linear reliance g_i on a targeted spurious feature at μ_{diag} . This creates a measurement-versus-optimization
 176 gap, as the mitigation is validated at μ_{diag} , where $g_i = 0$ by construction, but the optimized KL-regularized policy
 177 realizes its effects at the policy-induced distribution μ_{π^*} , where reliance generically reopens and optimization pressure
 178 rotates onto correlated features. Classifying outcomes by whether exploitation rotates onto another spurious feature
 179 ($\Delta_j \neq 0$) and by the sign of the true-reward change (ΔJ) yields five regimes, successful mitigation (R0), contaminated
 180 success (R0_{cont}), bias substitution (R1), overcorrection (R2), and silent non-op (R3), plus a transversal regime R4 whose
 181 label flips with the choice of μ_{diag} . We then prove a matched impossibility-sufficiency pair (Theorems A.13 and A.14).
 182 No benchmark functional over audit-distribution observables can separate R0 from R0_{cont}, R1, or R2, even granted
 183 oracle access to R , whereas augmenting the input with policy-induced distributions provably recovers the regime. This
 184 motivates our central prescription, that methods and benchmarks evaluate mitigations at policy-induced distributions
 185 and report off-target Δ_j alongside ΔJ . Full definitions, proofs, and closed-form instantiations are in Section A.

190 3 Bias Substitution in the Wild

192 Bias substitution is already pervasive in the published RLHF/DPO mitigation literature, miscategorized as isolated
 193 anomalies, capability trade-offs, or judge noise. Deployed reward models satisfy its precondition of correlated, partially
 194 informative spurious features, and we add length-sycophancy coupling measurements across eight model families,
 195 where the effect reverses from +154.3 characters under human-LLM judge agreement to −43.1 under disagreement, the
 196 empirical signature of R4 (Section F.3). To show the failure arises end to end, we run GRPO with a length penalty on
 197 *Llama-3.2-3B-Instruct* (Figure 1, Section F.1). As the penalty rises, mean response length falls from 204 to 170 tokens and
 198 MMLU accuracy holds, yet expected calibration error climbs from 0.25 to 0.41 and TriviaQA accuracy drops from 0.56
 199 to 0.42, while a $\lambda = 0$ control leaves calibration intact. Optimization pressure rotates off length onto overconfidence,
 200 instantiating harmful R1. A published length-debiasing operator likewise zeros pooled reward-length correlation
 201 (0.316 \rightarrow 0.037) yet flips within-prompt sign on three of four SOTA reward models, with true reward degrading on two
 202 (Section F.2). Surveying the literature under our prescriptions, no published method provides evidence sufficient to
 203 certify successful mitigation (R0), available evidence points to R0_{cont}, R1, or R2, and every benchmark we examined falls
 204 inside the impossibility class of Theorem A.13. See Section B for details and Section D for a broader impact analysis.

References

- [1] Anthropic. 2026. Models overview. Blog post. Accessed: 2026-05-06.
- [2] Jan Betley, Niels Warncke, Anna Sztyber-Betley, Daniel Tan, Xuchan Bao, Martín Soto, Megha Srivastava, Nathan Labenz, and Owain Evans. 2026. Training large language models on narrow tasks can lead to broad misalignment. *Nature* 649 (2026), 584–589. doi:10.1038/s41586-025-09937-5
- [3] Yuyan Bu, Liangyu Huo, Yi Jing, and Qing Yang. 2025. Beyond Excess and Deficiency: Adaptive Length Bias Mitigation in Reward Models for RLHF. In *Findings of the Association for Computational Linguistics: NAACL 2025*. Association for Computational Linguistics, Albuquerque, New Mexico, 3091–3098. doi:10.18653/v1/2025.findings-naacl.169
- [4] Jianfeng Cai, Jinhua Zhu, Ruopei Sun, Yue Wang, Li Li, Wengang Zhou, and Houqiang Li. 2026. Disentangling Length Bias in Preference Learning via Response-Conditioned Modeling. In *The Fourteenth International Conference on Learning Representations*. <https://openreview.net/forum?id=hKxYESOzen>
- [5] Stephen Casper, Xander Davies, Claudia Shi, Thomas Krendl Gilbert, Jérémy Scheurer, Javier Rando, Rachel Freedman, Tomek Korbak, David Lindner, Pedro Freire, Tony Tong Wang, Samuel Marks, Charbel-Raphael Segerie, Michah Carroll, Andi Peng, Phillip J.K. Christoffersen, Mehul Damani, Stewart Slocum, Usman Anwar, Anand Siththaranjan, Max Nadeau, Eric J Michaud, Jacob Pfau, Dmitrii Krashennikov, Xin Chen, Lauro Langosco, Peter Hase, Erdem Biyik, Anca Dragan, David Krueger, Dorsa Sadigh, and Dylan Hadfield-Menell. 2023. Open Problems and Fundamental Limitations of Reinforcement Learning from Human Feedback. *Transactions on Machine Learning Research* (2023). <https://openreview.net/forum?id=bx24KpJ4Eb> Survey Certification, Featured Certification.
- [6] Guiming Hardy Chen, Shunian Chen, Ziche Liu, Feng Jiang, and Benyou Wang. 2024. Humans or LLMs as the Judge? A Study on Judgement Bias. In *Proceedings of the 2024 Conference on Empirical Methods in Natural Language Processing*. Association for Computational Linguistics, Miami, Florida, USA, 8301–8327. doi:10.18653/v1/2024.emnlp-main.474
- [7] Longze Chen, Lu Wang, Renke Shan, Ze Gong, Run Luo, Jiaming Li, Jing Luo, Qiyao Wang, and Min Yang. 2026. Learning Ordinal Probabilistic Reward from Preferences. In *The Fourteenth International Conference on Learning Representations*. <https://openreview.net/forum?id=0Vf5trUAVF>
- [8] Lichang Chen, Chen Zhu, Jihui Chen, Davit Soselia, Tianyi Zhou, Tom Goldstein, Heng Huang, Mohammad Shoeybi, and Bryan Catanzaro. 2024. ODIN: disentangled reward mitigates hacking in RLHF. In *Proceedings of the 41st International Conference on Machine Learning (ICML'24)*. JMLR.org, Article 312, 18 pages.
- [9] Myra Cheng, Cinoo Lee, Pranav Khadpe, Sunny Yu, Dyllan Han, and Dan Jurafsky. 2026. Sycophantic AI decreases prosocial intentions and promotes dependence. *Science* 391 (2026), eaec8352. doi:10.1126/science.aec8352
- [10] Myra Cheng, Sunny Yu, Cinoo Lee, Pranav Khadpe, Lujain Ibrahim, and Dan Jurafsky. 2026. ELEPHANT: Measuring and understanding social sycophancy in LLMs. In *The Fourteenth International Conference on Learning Representations*. <https://openreview.net/forum?id=igbRHKiEiAs>
- [11] Yeshwanth Cherapanamjeri, Constantinos Costis Daskalakis, Gabriele Farina, and Sobhan Mohammadpour. 2026. Learning Correlated Reward Models: Statistical Barriers and Opportunities. In *The Fourteenth International Conference on Learning Representations*. <https://openreview.net/forum?id=TbEyl6krsY>
- [12] Wei-Lin Chiang, Lianmin Zheng, Ying Sheng, Anastasios N. Angelopoulos, Tianle Li, Dacheng Li, Banghua Zhu, Hao Zhang, Michael I. Jordan, Joseph E. Gonzalez, and Ion Stoica. 2024. Chatbot arena: an open platform for evaluating LLMs by human preference. In *Proceedings of the 41st International Conference on Machine Learning (ICML'24)*. JMLR.org, Article 331, 30 pages.
- [13] Brian Christian, Jessica A F Thompson, Elle, Vincent Adam, Hannah Rose Kirk, Christopher Summerfield, and Tsvetomira Dumbalska. 2026. Reward Models Inherit Value Biases from Pretraining. In *The Fourteenth International Conference on Learning Representations*. <https://openreview.net/forum?id=dT399j1Azv>
- [14] Paul F. Christiano, Jan Leike, Tom B. Brown, Miljan Martic, Shane Legg, and Dario Amodei. 2017. Deep reinforcement learning from human preferences. In *Proceedings of the 31st International Conference on Neural Information Processing Systems (NIPS'17)*. Curran Associates Inc., Red Hook, NY, USA, 4302–4310.
- [15] Karl Cobbe, Vineet Kosaraju, Mohammad Bavarian, Mark Chen, Heewoo Jun, Lukasz Kaiser, Matthias Plappert, Jerry Tworek, Jacob Hilton, Reiichiro Nakano, Christopher Hesse, and John Schulman. 2021. Training Verifiers to Solve Math Word Problems. arXiv:2110.14168 [cs.LG] <https://arxiv.org/abs/2110.14168> arXiv:2110.14168.
- [16] Ganqu Cui, Lifan Yuan, Ning Ding, Guanming Yao, Wei Zhu, Yuan Ni, Guotong Xie, Zhiyuan Liu, and Maosong Sun. 2023. UltraFeedback: Boosting Language Models with High-quality Feedback. arXiv:2310.01377 [cs.CL]
- [17] Magda Dubois, Cozmin Ududec, Christopher Summerfield, and Lennart Luettgau. 2026. Ask don't tell: Reducing sycophancy in large language models. arXiv:2602.23971 [cs.HC] <https://arxiv.org/abs/2602.23971>
- [18] Yann Dubois, Percy Liang, and Tatsunori Hashimoto. 2024. Length-Controlled AlpacaEval: A Simple Debiasing of Automatic Evaluators. In *First Conference on Language Modeling*. <https://openreview.net/forum?id=CybBmzWBX0>
- [19] Jacob Eisenstein, Chirag Nagpal, Alekh Agarwal, Ahmad Beirami, Alexander Nicholas D'Amour, Krishnamurthy Dj Dvijotham, Adam Fisch, Katherine A Heller, Stephen Robert Pfohl, Deepak Ramachandran, Peter Shaw, and Jonathan Berant. 2024. Helping or Herding? Reward Model Ensembles Mitigate but do not Eliminate Reward Hacking. In *First Conference on Language Modeling*. <https://openreview.net/forum?id=5u1GpUkKtG>
- [20] Aaron Fanous, Jacob Goldberg, Ank Agarwal, Joanna Lin, Anson Zhou, Sonnet Xu, Vasiliki Bikia, Roxana Daneshjou, and Sanmi Koyejo. 2025. SycEval: Evaluating LLM Sycophancy. In *Proceedings of the Eighth AAAI/ACM Conference on AI, Ethics, and Society*, Vol. 8. 893–900. doi:10.1609/aies.v8i1.36598

- 261 [21] Daniel Fein, Max Lamparth, Violet Xiang, Mykel J. Kochenderfer, and Nick Haber. 2026. One Bias After Another: Mechanistic Reward Shaping and
262 Persistent Biases in Language Reward Models. arXiv:2603.03291 [cs.CL] <https://arxiv.org/abs/2603.03291>
- 263 [22] Xuan Feng, Bo An, Tianlong Gu, Liang Chang, Fengrui Hao, Peipeng Yu, and Shuai Zhao. 2025. C2PO: Diagnosing and Disentangling Bias Shortcuts
264 in LLMs. arXiv:2512.23430 [cs.CL] <https://arxiv.org/abs/2512.23430>
- 265 [23] Lukas Fluri, Leon Lang, Alessandro Abate, Patrick Forré, David Krueger, and Joar Max Viktor Skalse. 2025. The Perils of Optimizing Learned Reward
266 Functions: Low Training Error Does Not Guarantee Low Regret. In *Proceedings of the 42nd International Conference on Machine Learning (Proceedings
267 of Machine Learning Research, Vol. 267)*, Aarti Singh, Maryam Fazel, Daniel Hsu, Simon Lacoste-Julien, Felix Berkenkamp, Tegan Maharaj, Kiri
268 Wagstaff, and Jerry Zhu (Eds.). PMLR, 17306–17377. <https://proceedings.mlr.press/v267/fluri25a.html>
- 269 [24] Jiayi Fu, Xuandong Zhao, Chengyuan Yao, Heng Wang, Qi Han, and Yanghua Xiao. 2025. Reward Shaping to Mitigate Reward Hacking in RLHF. In
270 *ICML 2025 Workshop on Reliable and Responsible Foundation Models*. <https://openreview.net/forum?id=62A4d5Mokc>
- 271 [25] Leo Gao, John Schulman, and Jacob Hilton. 2023. Scaling laws for reward model overoptimization. In *Proceedings of the 40th International Conference
272 on Machine Learning (ICML '23)*. JMLR.org, Article 437, 32 pages.
- 273 [26] Robert Geirhos, Jörn-Henrik Jacobsen, Claudio Michaelis, Richard Zemel, Wieland Brendel, Matthias Bethge, and Felix A. Wichmann. 2020. Shortcut
274 Learning in Deep Neural Networks. *Nature* 2 (2020), 665–673. doi:10.1038/s42256-020-00257-z
- 275 [27] Aaron Grattafiori, Abhimanyu Dubey, Abhinav Jauhri, Abhinav Pandey, Abhishek Kadian, Ahmad Al-Dahle, Aiesha Letman, Akhil Mathur, Alan
276 Schelten, Alex Vaughan, Amy Yang, Angela Fan, Anirudh Goyal, Anthony Hartshorn, Aobo Yang, Archi Mitra, Archie Sravankumar, Artem Korenev,
277 Arthur Hinsvark, Arun Rao, Aston Zhang, Aurelien Rodriguez, Austen Gregerson, Ava Spataru, Baptiste Roziere, et al. 2024. The Llama 3 Herd of
278 Models. arXiv:2407.21783 [cs.CL] <https://arxiv.org/abs/2407.21783> arXiv:2407.21783.
- 279 [28] Jiseung Hong, Grace Byun, Seungone Kim, and Kai Shu. 2025. Measuring Sycophancy of Language Models in Multi-turn Dialogues. In *Findings of
280 the Association for Computational Linguistics: EMNLP 2025*. Association for Computational Linguistics, Suzhou, China, 2239–2259. doi:10.18653/v1/
281 2025.findings-emnlp.121
- 282 [29] Zhengyu Hu, Linxin Song, Jieyu Zhang, Zheyuan Xiao, Tianfu Wang, Zhengyu Chen, Nicholas Jing Yuan, Jianxun Lian, Kaize Ding, and Hui Xiong.
283 2025. Explaining Length Bias in LLM-Based Preference Evaluations. In *Findings of the Association for Computational Linguistics: EMNLP 2025*.
284 Association for Computational Linguistics, Suzhou, China, 6763–6794. doi:10.18653/v1/2025.findings-emnlp.358
- 285 [30] Zeyu Huang, Zihan Qiu, Zili Wang, Edoardo Ponti, and Ivan Titov. 2025. Post-hoc Reward Calibration: A Case Study on Length Bias. In *The
286 Thirteenth International Conference on Learning Representations*. <https://openreview.net/forum?id=lu8RytBaji>
- 287 [31] Lujain Ibrahim, Katherine M. Collins, Sunnie S. Y. Kim, Anka Reuel, Max Lamparth, Kevin Feng, Lama Ahmad, Prajna Soni, Alia El Kattan, Merlin
288 Stein, Siddharth Swaroop, Iliia Sucholutsky, Andrew Strait, Q. Vera Liao, and Umang Bhatt. 2025. Measuring and mitigating overreliance is necessary
289 for building human-compatible AI. arXiv:2509.08010 [cs.CY] <https://arxiv.org/abs/2509.08010>
- 290 [32] Lujain Ibrahim, Franziska Sofia Hafner, and Luc Rocher. 2026. Training language models to be warm can reduce accuracy and increase sycophancy.
291 *Nature* 652 (2026), 1159–1165. doi:10.1038/s41586-026-10410-0
- 292 [33] Timo Kaufmann, Paul Weng, Viktor Bengs, and Eyke Hüllermeier. 2025. A Survey of Reinforcement Learning from Human Feedback. *Transactions
293 on Machine Learning Research* (2025). <https://openreview.net/forum?id=f70klurx4b> Survey Certification.
- 294 [34] Michael Kearns, Seth Neel, Aaron Roth, and Zhiwei Steven Wu. 2018. Preventing Fairness Gerrymandering: Auditing and Learning for Subgroup
295 Fairness. In *Proceedings of the 35th International Conference on Machine Learning (Proceedings of Machine Learning Research, Vol. 80)*. PMLR,
296 2564–2572.
- 297 [35] Emaan Bilal Khan, Amy Winecoff, Miranda Bogen, and Dylan Hadfield-Menell. 2026. Safety Drift After Fine-Tuning: Evidence from High-Stakes
298 Domains. arXiv:2604.24902 [cs.CY] <https://arxiv.org/abs/2604.24902>
- 299 [36] Hyeonji Kim, Sujeong Oh, and Sanghack Lee. 2025. Mitigating Length Bias in RLHF through a Causal Lens. arXiv:2511.12573 [cs.CL] <https://arxiv.org/abs/2511.12573>
- 300 [37] Kuno Kim, Shivam Garg, Kirankumar Shiragur, and Stefano Ermon. 2021. Reward Identification in Inverse Reinforcement Learning. In *Proceedings
301 of the 38th International Conference on Machine Learning (Proceedings of Machine Learning Research, Vol. 139)*. PMLR, 5496–5505.
- 302 [38] Ryan Koo, Minhwa Lee, Vipul Raheja, Jong Inn Park, Zae Myung Kim, and Dongyeop Kang. 2024. Benchmarking Cognitive Biases in Large
303 Language Models as Evaluators. In *Findings of the Association for Computational Linguistics: ACL 2024*. Association for Computational Linguistics,
304 Bangkok, Thailand, 517–545. doi:10.18653/v1/2024.findings-acl.29
- 305 [39] Ashwin Kumar, Yuzi He, Aram H Markosyan, Bobbie Chern, and Imanol Arrieta-Ibarra. 2025. Detecting Prefix Bias in LLM-based Reward Models.
306 In *Proceedings of the 2025 ACM Conference on Fairness, Accountability, and Transparency (FAccT '25)*. Association for Computing Machinery, New
307 York, NY, USA, 3196–3206. doi:10.1145/3715275.3732204
- 308 [40] Thomas Kwa, Drake Thomas, and Adrià Garriga-Alonso. 2024. Catastrophic Goodhart: regularizing RLHF with KL divergence does not mitigate
309 heavy-tailed reward misspecification. In *The Thirty-eighth Annual Conference on Neural Information Processing Systems*. [https://openreview.net/
310 forum?id=UXuBzWoZGK](https://openreview.net/forum?id=UXuBzWoZGK)
- 311 [41] Cassidy Laidlaw, Shivam Singhal, and Anca Dragan. 2025. Correlated Proxies: A New Definition and Improved Mitigation for Reward Hacking. In
312 *The Thirteenth International Conference on Learning Representations*. <https://openreview.net/forum?id=msEr27EejF>
- [42] Nathan Lambert. 2026. Reinforcement Learning from Human Feedback. arXiv:2504.12501 [cs.LG] <https://arxiv.org/abs/2504.12501>
- [43] Nathan Lambert, Valentina Pyatkin, Jacob Morrison, LJ Miranda, Bill Yuchen Lin, Khyathi Chandu, Nouha Dziri, Sachin Kumar, Tom Zick,
Yejin Choi, Noah A. Smith, and Hannaneh Hajishirzi. 2025. RewardBench: Evaluating Reward Models for Language Modeling. In *Findings of*

- 313 *the Association for Computational Linguistics: NAACL 2025*. Association for Computational Linguistics, Albuquerque, New Mexico, 1755–1797.
314 doi:10.18653/v1/2025.findings-naacl.96
- 315 [44] Jixuan Leng, Chengsong Huang, Banghua Zhu, and Jiaxin Huang. 2025. Taming Overconfidence in LLMs: Reward Calibration in RLHF. In *The*
316 *Thirteenth International Conference on Learning Representations*. <https://openreview.net/forum?id=10tg0jzsdL>
- 317 [45] Gengxu Li, Tingyu Xia, Yi Chang, and Yuan Wu. 2025. Length-Controlled Margin-Based Preference Optimization without Reference Model.
318 arXiv:2502.14643 [cs.CL] <https://arxiv.org/abs/2502.14643>
- 319 [46] Tianle Li, Anastasios Angelopoulos, and Wei-Lin Chiang. 2024. Does Style Matter? Disentangling Style and Substance in Chatbot Arena. Blog post.
320 Accessed: 2026-05-06.
- 321 [47] Zhuo Li, Pengyu Cheng, Zhechao Yu, Feifei Tong, Anningzhe Gao, Tsung-Hui Chang, Xiang Wan, Erchao Zhao, Xiaoxi Jiang, and Guanjun Jiang.
322 2025. Eliminating Inductive Bias in Reward Models with Information-Theoretic Guidance. arXiv:2512.23461 [cs.LG] <https://arxiv.org/abs/2512.23461>
- 323 [48] Zhiheng Li, Ivan Evtimov, Albert Gordo, Caner Hazirbas, Tal Hassner, Cristian Canton Ferrer, Chenliang Xu, and Mark Ibrahim. 2023. A Whac-a-Mole
324 Dilemma: Shortcuts Come in Multiples Where Mitigating One Amplifies Others. In *Proceedings of the IEEE/CVF Conference on Computer Vision and*
Pattern Recognition (CVPR).
- 325 [49] Chris Yuhao Liu, Liang Zeng, Yuzhen Xiao, Jujie He, Jiakai Liu, Chaojie Wang, Rui Yan, Wei Shen, Fuxiang Zhang, Jiacheng Xu, Yang Liu, and Yahui
326 Zhou. 2025. Skywork-Reward-V2: Scaling Preference Data Curation via Human-AI Synergy. arXiv:2507.01352 [cs.CL] <https://arxiv.org/abs/2507.01352>
327 arXiv:2507.01352.
- 328 [50] Zixuan Liu, Xiaolin Sun, and Zizhan Zheng. 2026. Robust Optimization for Mitigating Reward Hacking with Correlated Proxies. In *The Fourteenth*
329 *International Conference on Learning Representations*. <https://openreview.net/forum?id=O3shkBMW2s>
- 330 [51] Junru Lu, Jiazheng Li, Siyu An, Meng Zhao, Yulan He, Di Yin, and Xing Sun. 2024. Eliminating Biased Length Reliance of Direct Preference
331 Optimization via Down-Sampled KL Divergence. In *Proceedings of the 2024 Conference on Empirical Methods in Natural Language Processing*.
Association for Computational Linguistics, 1047–1067. doi:10.18653/v1/2024.emnlp-main.60
- 332 [52] Saumya Malik, Valentina Pyatkin, Sander Land, Jacob Morrison, Noah A. Smith, Hannaneh Hajishirzi, and Nathan Lambert. 2026. RewardBench 2:
333 Advancing Reward Model Evaluation. In *The Fourteenth International Conference on Learning Representations*. [https://openreview.net/forum?id=](https://openreview.net/forum?id=fb0G86Dewb)
334 [fb0G86Dewb](https://openreview.net/forum?id=fb0G86Dewb)
- 335 [53] Yu Meng, Mengzhou Xia, and Danqi Chen. 2024. SimPO: simple preference optimization with a reference-free reward. In *Proceedings of the 38th*
336 *International Conference on Neural Information Processing Systems (NIPS '24)*. Curran Associates Inc., Article 3946, 38 pages.
- 337 [54] Rajiv Movva, Smitha Milli, Sewon Min, and Emma Pierson. 2026. What's In My Human Feedback? Learning Interpretable Descriptions of Preference
338 Data. In *The Fourteenth International Conference on Learning Representations*. <https://openreview.net/forum?id=sC6A1bFDUf>
- 339 [55] Ignavier Ng, Patrick Blöbaum, Siddharth Bhandari, Kun Zhang, and Shiva Kasiviswanathan. 2025. Debiasing Reward Models by Representation
340 Learning with Guarantees. arXiv:2510.23751 [cs.LG] <https://arxiv.org/abs/2510.23751>
- 341 [56] Nvidia, ., Bo Adler, Niket Agarwal, Ashwath Aithal, Dong H. Anh, Pallab Bhattacharya, Annika Brundyn, Jared Casper, Bryan Catanzaro, Sharon
342 Clay, Jonathan Cohen, Sirshak Das, Ayush Dattagupta, Olivier Delalleau, Leon Derczynski, Yi Dong, Daniel Egert, Ellie Evans, Aleksander Fieck,
343 Denys Fridman, Shaona Ghosh, Boris Ginsburg, Igor Gitman, Tomasz Grzegorzec, Robert Hero, Jining Huang, Vibhu Jawa, Joseph Jennings,
344 Aastha Jhunjhunwala, John Kamalu, Sadaf Khan, Oleksii Kuchaiev, Patrick LeGresley, Hui Li, Jiwei Liu, Zihan Liu, Eileen Long, Ameya Sunil
345 Mahabaleshwarkar, Somshubra Majumdar, James Maki, Miguel Martinez, Maer Rodrigues de Melo, Ivan Moshkov, Deepak Narayanan, Sean
346 Narenthiran, Jesus Navarro, Phong Nguyen, Osvald Nitski, Vahid Noroozi, Guruprasad Nutheti, Christopher Parisien, Jupinder Parmar, Mostofa
347 Patwary, Krzysztof Pawelec, Wei Ping, Shrimai Prabhunoye, Rajarshi Roy, Trisha Saar, Vasanth Rao Naik Sabavat, Sanjeev Satheesh, Jane Polak
348 Scowcroft, Jason Sewall, Pavel Shamis, Gerald Shen, Mohammad Shoeybi, Dave Sizer, Misha Smelyanskiy, Felipe Soares, Makesh Narsimhan Sreedhar,
349 Dan Su, Sandeep Subramanian, Shengyang Sun, Shubham Toshniwal, Hao Wang, Zhilin Wang, Jiaxuan You, Jiaqi Zeng, Jimmy Zhang, Jing Zhang,
350 Vivienne Zhang, Yian Zhang, and Chen Zhu. 2024. NemoTron-4 340B Technical Report. arXiv:2406.11704 [cs.CL] <https://arxiv.org/abs/2406.11704>
- 351 [57] OpenAssistant. 2023. OpenAssistant/reward-model-deberta-v3-large-v2. <https://huggingface.co/OpenAssistant/reward-model-deberta-v3-large-v2>.
352 Hugging Face model card. Accessed: 2026-01-06.
- 353 [58] Henry Papadatos and Rachel Freedman. 2024. Linear Probe Penalties Reduce LLM Sycophancy. In *Workshop on Socially Responsible Language*
354 *Modelling Research*. <https://openreview.net/forum?id=6N2yES22rG>
- 355 [59] Ryan Park, Rafael Rafailov, Stefano Ermon, and Chelsea Finn. 2024. Disentangling Length from Quality in Direct Preference Optimization. In
356 *Findings of the Association for Computational Linguistics: ACL 2024*. Association for Computational Linguistics, Bangkok, Thailand, 4998–5017.
357 doi:10.18653/v1/2024.findings-acl.297
- 358 [60] Judea Pearl and Elias Bareinboim. 2022. *External Validity: From Do-Calculus to Transportability Across Populations* (1 ed.). Association for Computing
359 Machinery, New York, NY, USA, 451–482. <https://doi.org/10.1145/3501714.3501741>
- 360 [61] Juan Perdomo, Tijana Zrnic, Celestine Mender-Dünner, and Moritz Hardt. 2020. Performative Prediction. In *Proceedings of the 37th International*
361 *Conference on Machine Learning (Proceedings of Machine Learning Research, Vol. 119)*. PMLR, 7599–7609.
- 362 [62] Ethan Perez, Sam Ringer, Kamile Lukosiute, Karina Nguyen, Edwin Chen, Scott Heiner, Craig Pettit, Catherine Olsson, Sandipan Kundu, Saurav
363 Kadavath, Andy Jones, Anna Chen, Benjamin Mann, Brian Israel, Bryan Seethor, Cameron McKinnon, Christopher Olah, Da Yan, Daniela Amodei,
364 Dario Amodei, Dawn Drain, Dustin Li, Eli Tran-Johnson, Guro Khundadze, Jackson Kernion, James Landis, Jamie Kerr, Jared Mueller, Jeeyoon
Hyun, Joshua Landau, Kamal Ndousse, Landon Goldberg, Liane Lovitt, Martin Lucas, Michael Sellitto, Miranda Zhang, Neerav Kingsland, Nelson
Elhage, Nicholas Joseph, Noemi Mercado, Nova DasSarma, Oliver Rausch, Robin Larson, Sam McCandlish, Scott Johnston, Shauna Kravec, Sheer

- 365 El Showk, Tamera Lanham, Timothy Telleen-Lawton, Tom Brown, Tom Henighan, Tristan Hume, Yuntao Bai, Zac Hatfield-Dodds, Jack Clark,
366 Samuel R. Bowman, Amanda Askell, Roger Grosse, Danny Hernandez, Deep Ganguli, Evan Hubinger, Nicholas Schiefer, and Jared Kaplan. 2023.
367 Discovering Language Model Behaviors with Model-Written Evaluations. In *Findings of the Association for Computational Linguistics: ACL 2023*.
368 Association for Computational Linguistics, 13387–13434. doi:10.18653/v1/2023.findings-acl.847
- 369 [63] Jan Peters and Stefan Schaal. 2007. Reinforcement learning by reward-weighted regression for operational space control. In *Proceedings of the 24th*
370 *International Conference on Machine Learning (ICML '07)*. Association for Computing Machinery, 745–750. doi:10.1145/1273496.1273590
- 371 [64] Rafael Rafailov, Archit Sharma, Eric Mitchell, Christopher D Manning, Stefano Ermon, and Chelsea Finn. 2023. Direct Preference Optimization:
372 Your Language Model is Secretly a Reward Model. In *Thirty-seventh Conference on Neural Information Processing Systems*. [https://openreview.net/
forum?id=HPuSLXJaa9](https://openreview.net/forum?id=HPuSLXJaa9)
- 373 [65] Keita Saito, Akifumi Wachi, Koki Wataoka, and Youhei Akimoto. 2023. Verbosity Bias in Preference Labeling by Large Language Models. In *NeurIPS*
374 *2023 Workshop on Instruction Tuning and Instruction Following*. <https://openreview.net/forum?id=magEgFpK1y>
- 375 [66] Harshay Shah, Kaustav Tamuly, Aditi Raghunathan, Prateek Jain, and Praneeth Netrapalli. 2020. The Pitfalls of Simplicity Bias in Neural Networks.
376 In *Advances in Neural Information Processing Systems (NeurIPS)*.
- 377 [67] Zhihong Shao, Peiyi Wang, Qihao Zhu, Runxin Xu, Junxiao Song, Xiao Bi, Haowei Zhang, Mingchuan Zhang, Y. K. Li, Y. Wu, and Daya Guo. 2024.
378 DeepSeekMath: Pushing the Limits of Mathematical Reasoning in Open Language Models. arXiv:2402.03300 [cs.CL] <https://arxiv.org/abs/2402.03300>
- 379 [68] Itai Shapira, Gerdus Benade, and Ariel D. Procaccia. 2026. How RLHF Amplifies Sycophancy. arXiv:2602.01002 [cs.AI] <https://arxiv.org/abs/2602.01002>
- 380 [69] Mrinank Sharma, Meg Tong, Tomasz Korbak, David Duvenaud, Amanda Askell, Samuel R. Bowman, Esin DURMUS, Zac Hatfield-Dodds, Scott R
381 Johnston, Shauna M Kravec, Timothy Maxwell, Sam McCandlish, Kamal Ndousse, Oliver Rausch, Nicholas Schiefer, Da Yan, Miranda Zhang, and
382 Ethan Perez. 2024. Towards Understanding Sycophancy in Language Models. In *The Twelfth International Conference on Learning Representations*.
383 <https://openreview.net/forum?id=tvhaxkMKAn>
- 384 [70] Wei Shen, Rui Zheng, Wenyu Zhan, Jun Zhao, Shihan Dou, Tao Gui, Qi Zhang, and Xuanjing Huang. 2023. Loose lips sink ships: Mitigating
385 Length Bias in Reinforcement Learning from Human Feedback. In *The 2023 Conference on Empirical Methods in Natural Language Processing*.
<https://openreview.net/forum?id=q6qctdUwCX>
- 386 [71] Prasann Singhal, Tanya Goyal, Jiacheng Xu, and Greg Durrett. 2024. A Long Way to Go: Investigating Length Correlations in RLHF. In *First*
387 *Conference on Language Modeling*. <https://openreview.net/forum?id=G8LaO1P0xv>
- 388 [72] Joar Skalse, Matthew Farrugia-Roberts, Stuart Russell, Alessandro Abate, and Adam Gleave. 2023. Invariance in policy optimisation and partial
389 identifiability in reward learning. In *Proceedings of the 40th International Conference on Machine Learning (ICML '23)*. JMLR.org, Article 1328, 26 pages.
- 390 [73] Joar Skalse, Nikolaus H. R. Howe, Dmitrii Krashenninikov, and David Krueger. 2022. Defining and characterizing reward hacking. In *Proceedings of*
391 *the 36th International Conference on Neural Information Processing Systems (NIPS '22)*. Curran Associates Inc., Article 687, 12 pages.
- 392 [74] Ruike Song, Zeen Song, Huijie Guo, and Wenwen Qiang. 2025. Causal Reward Adjustment: Mitigating Reward Hacking in External Reasoning via
393 Backdoor Correction. arXiv:2508.04216 [cs.LG] <https://arxiv.org/abs/2508.04216>
- 394 [75] Pragya Srivastava, Harman Singh, Rahul Madhavan, Gandharv Patil, Sravanti Addepalli, Arun Suggala, Rengarajan Aravamudhan, Soumya Sharma,
395 Anirban Laha, Aravindan Raghuvver, Karthikeyan Shanmugam, and Doina Precup. 2026. Robust Reward Modeling via Causal Rubrics. In *The*
Fourteenth International Conference on Learning Representations. <https://openreview.net/forum?id=oP99jQiDYp>
- 396 [76] Nisan Stiennon, Long Ouyang, Jeff Wu, Daniel M. Ziegler, Ryan Lowe, Chelsea Voss, Alec Radford, Dario Amodei, and Paul Christiano. 2020.
397 Learning to summarize from human feedback. In *Proceedings of the 34th International Conference on Neural Information Processing Systems (NIPS '20)*.
398 Curran Associates Inc., Article 253, 14 pages.
- 399 [77] Katherine Tian, Eric Mitchell, Allan Zhou, Archit Sharma, Rafael Rafailov, Huaxiu Yao, Chelsea Finn, and Christopher Manning. 2023. Just Ask
400 for Calibration: Strategies for Eliciting Calibrated Confidence Scores from Language Models Fine-Tuned with Human Feedback. In *Proceedings*
401 *of the 2023 Conference on Empirical Methods in Natural Language Processing*. Association for Computational Linguistics, Singapore, 5433–5442.
402 doi:10.18653/v1/2023.emnlp-main.330
- 403 [78] Jeremy Tien, Jerry Zhi-Yang He, Zackory Erickson, Anca Dragan, and Daniel S. Brown. 2023. Causal Confusion and Reward Misidentification
404 in Preference-Based Reward Learning. In *The Eleventh International Conference on Learning Representations*. [https://openreview.net/forum?id=
R0Xxvr_X3ZA](https://openreview.net/forum?id=R0Xxvr_X3ZA)
- 405 [79] Leandro von Werra, Younes Belkada, Lewis Tunstall, Edward Beeching, Tristan Thrush, Nathan Lambert, Shengyi Huang, Kashif Rasul, and Quentin
406 Galloudec. 2020. *TRL: Transformers Reinforcement Learning*. <https://github.com/huggingface/trl>
- 407 [80] Chaoqi Wang, Zhuokai Zhao, Yibo Jiang, Zhaorun Chen, Chen Zhu, Yuxin Chen, Jiayi Liu, Lizhu Zhang, Xiangjun Fan, Hao Ma, and Sinong Wang.
408 2025. Beyond Reward Hacking: Causal Rewards for Large Language Model Alignment. arXiv:2501.09620 [cs.LG] <https://arxiv.org/abs/2501.09620>
- 409 [81] Zhilin Wang, Yi Dong, Jiaqi Zeng, Virginia Adams, Makesh Narsimhan Sreedhar, Daniel Egert, Olivier Delalleau, Jane Scowcroft, Neel Kant,
410 Aidan Swope, and Oleksii Kuchaiev. 2024. HelpSteer: Multi-attribute Helpfulness Dataset for SteerLM. In *Proceedings of the 2024 Conference of the*
411 *North American Chapter of the Association for Computational Linguistics: Human Language Technologies (Volume 1: Long Papers)*. Association for
412 Computational Linguistics, 3371–3384. doi:10.18653/v1/2024.naacl-long.185
- 413 [82] Koki Wataoka, Tsubasa Takahashi, and Ryokan Ri. 2024. Self-Preference Bias in LLM-as-a-Judge. In *Neurips Safe Generative AI Workshop 2024*.
414 <https://openreview.net/forum?id=tLZZZigPjX>
- 415 [83] Jiaxin Wen, Ruiqi Zhong, Akbir Khan, Ethan Perez, Jacob Steinhardt, Minlie Huang, Samuel R. Bowman, He He, and Shi Feng. 2025. Language
416 Models Learn to Mislead Humans via RLHF. In *The Thirteenth International Conference on Learning Representations*. <https://openreview.net/forum/>

- 417 [id=xJljiPE6dg](#)
- 418 [84] Yuchen Yan, Jin Jiang, Zhenbang Ren, Yijun Li, Xudong Cai, Yang Liu, Xin Xu, Mengdi Zhang, Jian Shao, Yongliang Shen, Jun Xiao, and Yueting
- 419 Zhuang. 2026. VerifyBench: Benchmarking Reference-based Reward Systems for Large Language Models. In *The Fourteenth International Conference*
- 420 *on Learning Representations*. <https://openreview.net/forum?id=JfsjGmuFxz>
- 421 [85] Jiayi Ye, Yanbo Wang, Yue Huang, Dongping Chen, Qihui Zhang, Nuno Moniz, Tian Gao, Werner Geyer, Chao Huang, Pin-Yu Chen, Nitesh V
- 422 Chawla, and Xiangliang Zhang. 2025. Justice or Prejudice? Quantifying Biases in LLM-as-a-Judge. In *The Thirteenth International Conference on*
- 423 *Learning Representations*. <https://openreview.net/forum?id=3GTtZFiajM>
- 424 [86] Wenqian Ye, Guangtao Zheng, and Aidong Zhang. 2026. Rectifying Shortcut Behaviors in Preference-based Reward Learning. In *The Thirty-ninth*
- 425 *Annual Conference on Neural Information Processing Systems*. <https://openreview.net/forum?id=m51t6RKFGH>
- 426 [87] Yaowen Ye, Cassidy Laidlaw, and Jacob Steinhardt. 2025. Iterative Label Refinement Matters More than Preference Optimization under Weak
- 427 Supervision. In *The Thirteenth International Conference on Learning Representations*. <https://openreview.net/forum?id=q5EZ7gKcnW>
- 428 [88] Dongkeun Yoon, Seungone Kim, Sohee Yang, Sunkyoung Kim, Soyeon Kim, Yongil Kim, Eunbi Choi, Yireun Kim, and Minjoon Seo. 2025. Reasoning
- 429 Models Better Express Their Confidence. In *The Thirty-ninth Annual Conference on Neural Information Processing Systems*. [https://openreview.net/](https://openreview.net/forum?id=rbBtoVnduo)
- 430 [forum?id=rbBtoVnduo](https://openreview.net/forum?id=rbBtoVnduo)
- 431 [89] Danlong Yuan, Tian Xie, Shaohan Huang, Zhuocheng Gong, Huishuai Zhang, Chong Luo, Furu Wei, and Dongyan Zhao. 2026. Shorten After You're
- 432 Right: Lazy Length Penalties for Reasoning RL. arXiv:2505.12284 [cs.AI] <https://arxiv.org/abs/2505.12284>
- 433 [90] Yusen Zhang, Sarkar Snigdha Sarathi Das, and Rui Zhang. 2025. Demystify Verbosity Compensation Behavior of Large Language Models. In
- 434 *Proceedings of the 2nd Workshop on Uncertainty-Aware NLP (UncertainNLP 2025)*. Association for Computational Linguistics, Suzhou, China, 160–178.
- 435 [doi:10.18653/v1/2025.uncertainlp-main.14](https://doi.org/10.18653/v1/2025.uncertainlp-main.14)
- 436 [91] Kangwen Zhao, Jianfeng Cai, Jinhua Zhu, Ruopei Sun, Dongyun Xue, Wengang Zhou, Li Li, and Houqiang Li. 2025. Bias Fitting to Mitigate Length
- 437 Bias of Reward Model in RLHF. arXiv:2505.12843 [cs.LG] <https://arxiv.org/abs/2505.12843>
- 438 [92] Lianmin Zheng, Wei-Lin Chiang, Ying Sheng, Siyuan Zhuang, Zhanghao Wu, Yonghao Zhuang, Zi Lin, Zhuohan Li, Dacheng Li, Eric P. Xing, Hao
- 439 Zhang, Joseph E. Gonzalez, and Ion Stoica. 2023. Judging LLM-as-a-judge with MT-bench and Chatbot Arena. In *Proceedings of the 37th International*
- 440 *Conference on Neural Information Processing Systems (NIPS '23)*. Curran Associates Inc., Article 2020, 29 pages.
- 441 [93] Kaitlyn Zhou, Jena D. Hwang, Xiang Ren, and Maarten Sap. 2024. Relying on the Unreliable: The Impact of Language Models' Reluctance to Express
- 442 Uncertainty. In *Proceedings of the 62nd Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*. Association for
- 443 Computational Linguistics, 3623–3643. [doi:10.18653/v1/2024.acl-long.198](https://doi.org/10.18653/v1/2024.acl-long.198)

443 A Formalizing Bias Substitution

444 To formalize bias substitution, we turn the data- and policy-level invariances of Section 2 into a concrete failure-mode

445 taxonomy, prove what audit-distribution evaluation can and cannot certify, and derive prescriptions for mitigation

446 method claims and reward model benchmarks

447

448

449 A.1 Feature Map and Spurious vs. Structurally Relevant Features

450 In this work, we are interested in categorizing and analyzing the effects of bias mitigation strategies given a multiplicity

451 of axes that we may care to optimize for. We will refer to these axes as (surface) features. Before we can disentangle

452 their interactions with rewards, we need to define features.

453

454

455 **Definition A.1** (Feature map). A feature map is a measurable function $\phi : \mathcal{X} \times \mathcal{Y} \rightarrow \mathbb{R}$ capturing an interpretable surface

456 attribute of a response. To form a set of features that we care about, we fix a finite ordered tuple of feature maps ϕ_1, \dots, ϕ_K

457 and collect them into the vector-valued map $\Phi = (\phi_1, \dots, \phi_K)^\top : \mathcal{X} \times \mathcal{Y} \rightarrow \mathbb{R}^K$.

458

459 We use Φ both as the column-vector-valued map $(\phi_1, \dots, \phi_K)^\top$ and, by slight abuse, as the set $\{\phi_1, \dots, \phi_K\}$ when

460 membership ($\phi_i \in \Phi$) or partitions ($\Phi = \Phi_{\text{sp}} \sqcup \Phi_{\text{struct}}$) are more natural.

461 To instantiate our feature sets, we will make the natural assumption that the chosen feature sets capture distinct

462 response attributes under the diagnostic distribution. For example, features such as length, sycophancy, politeness, or

463 hedging may be correlated, but they should not collapse into the same measurement across diagnostic responses. This

464 rules out duplicate or exactly redundant features and ensures that each feature represents a separately identifiable axis

465 of response behavior.

466

467

468

469 **Assumption A.2** (Non-degeneracy). *The Gram matrix $G := \mathbb{E}_{\mu_{\text{diag}}} [\Phi\Phi^\top] \in \mathbb{R}^{K \times K}$, with entries $G_{ij} = \mathbb{E}_{\mu_{\text{diag}}} [\phi_i\phi_j]$, is*
 470 *positive definite: $G \succ 0$.*
 471

472 With these distinct axes in place, we can ask which attributes the reward actually depends on.

473 **Definition A.3** (Spurious vs. structurally relevant at μ_{diag}). *Treat R as a function on $\mathcal{X} \times \mathcal{Y}$ and assume feature realizability*
 474 *(Assumption E.2). Feature ϕ_i is spurious with respect to R at μ_{diag} if*
 475

$$476 \mathbb{E}_{\mu_{\text{diag}}} [R(x, y) \mid x, (\phi_j(x, y))_{j \neq i}] \text{ does not depend on } \phi_i(x, y) \text{ for } \mu_{\text{diag}}\text{-a.e.} \quad (3)$$

478 Otherwise ϕ_i is structurally relevant at μ_{diag} . Write $\Phi_{\text{sp}} = \{\phi_i \in \Phi : \phi_i \text{ spurious w.r.t. } R \text{ at } \mu_{\text{diag}}\}$ and $\Phi_{\text{struct}} = \Phi \setminus \Phi_{\text{sp}}$ for
 479 the induced partition of Φ .
 480

481 See Section E.4 for the equivalent $o(\varepsilon)$ formulation under ε -mixture perturbations, capturing the intuition that R
 482 has no first-order dependence on ϕ_i at μ_{diag} . The corresponding causal reading and partial-identifiability relationship
 483 are given in Section E.1. Neither spuriousness nor structural relevance is fully identifiable from preference data alone
 484 [72]. Under partial informativeness, natural features like length act as mediators of multiple mechanisms. The binary
 485 partition of Φ from Definition A.3 may therefore conservatively classify them as Φ_{struct} when they correlate with R
 486 under μ_{diag} (Section E.3). This partition of Φ also matches existing single-axis mitigations acting on whole features rather
 487 than within-feature decompositions [e.g. 6, 21, 30, 58]. Finer partitions require stronger assumptions (interventions,
 488 gold labels, distributional invariance) outside our regime.
 489
 490
 491

492 A.2 Single-Axis Mitigation and Measurement-Versus-Optimization Gap

493 Several prominent reward-model mitigations identify a feature direction, estimate \tilde{R} 's reliance on it, and subtract its
 494 contribution as, e.g., linear probes [21, 58], non-linear calibration [30], or architectural disentanglement [6]. To abstract
 495 away from these implementation details, we first define a population-level notion of reliance, given by the coefficients
 496 of \tilde{R} 's best linear approximation in Φ under the diagnostic distribution.
 497
 498

499 **Definition A.4** (Linear reliance). *The linear reliance of \tilde{R} on Φ at μ_{diag} is*

$$500 g(\tilde{R}; \mu_{\text{diag}}) = (\mathbb{E}_{\mu_{\text{diag}}} [\Phi\Phi^\top])^{-1} \mathbb{E}_{\mu_{\text{diag}}} [\Phi \tilde{R}] \in \mathbb{R}^K, \quad (4)$$

501 yielding the $L^2(\mu_{\text{diag}})$ -orthogonal decomposition $\tilde{R} = \sum_i g_i \phi_i + \tilde{R}_\perp$ with $\mathbb{E}_{\mu_{\text{diag}}} [\phi_i \tilde{R}_\perp] = 0$.
 502
 503

504 Note g is a statistic of the triple $(\tilde{R}, \Phi, \mu_{\text{diag}})$, not a property of \tilde{R} alone and evaluating at $\mu_{\pi^*} \neq \mu_{\text{diag}}$ yields a different
 505 vector, driving later regime distinctions in our taxonomy.
 506
 507

508 **Definition A.5** (Single-axis mitigation). *A single-axis mitigation targeting feature i is an operator $M_i : \tilde{R} \mapsto \tilde{R}'$ with*
 509 *$|g_i(\tilde{R}'; \mu_{\text{diag}})| < |g_i(\tilde{R}; \mu_{\text{diag}})|$. The canonical projection instance is*

$$510 M_i(\tilde{R})(x, y) = \tilde{R}(x, y) - g_i(\tilde{R}; \mu_{\text{diag}}) \phi_i(x, y). \quad (5)$$

511 We single out the projection instance above as canonical, because it zeros g_i at μ_{diag} exactly. The canonical M_i
 512 depends on μ_{diag} through $g_i(\tilde{R}; \mu_{\text{diag}})$. We write $M_i^{\mu_{\text{diag}}}$ when this dependence matters.
 513
 514

515 **Lemma A.6** (Single-axis identity at μ_{diag}). *For the canonical M_i , $g_i(M_i(\tilde{R}); \mu_{\text{diag}}) = 0$ and $g_j(M_i(\tilde{R}); \mu_{\text{diag}}) = g_j(\tilde{R}; \mu_{\text{diag}})$*
 516 *for $j \neq i$.*
 517

518 **PROOF.** Direct computation using $\mathbb{E}_{\mu_{\text{diag}}} [\Phi\phi_i]$ as the i -th column of $\mathbb{E}_{\mu_{\text{diag}}} [\Phi\Phi^\top]$. □
 519

The identity justifies calling M_i *single-axis*, as at μ_{diag} , mitigation moves \tilde{R} along exactly the i -th coordinate of g -space.¹ We emphasize that M_i is an *associational* operator, as it removes ϕ_i -reliance at μ_{diag} in the projection sense, not the contribution of ϕ_i along any causal path to R (see Appendix E.1). Our resulting taxonomy of classifies the resulting outcomes by what the *optimizer does at μ_{π^*} , not by what is observable at μ_{diag} .*

Measurement-versus-optimization gap. Single-axis diagnostics evaluate M_i at μ_{diag} , where $|g_i| = 0$ by construction (Lemma A.6). The optimizing policy realizes M_i 's effects at μ_{π^*} , where in general $g_i(M_i(\tilde{R}); \mu_{\pi^*}) \neq 0$ and $g_j(M_i(\tilde{R}); \mu_{\pi^*}) \neq g_j(\tilde{R}; \mu_{\pi^*})$ for $j \neq i$. Mitigation moves the proxy along an axis defined at the audit distribution, but optimization responds to a vector defined at a different distribution. This gap is the structural mechanism enabling bias substitution, and it is invisible to any diagnostic that operates at μ_{diag} alone. Section A.5 shows it cannot be closed by any audit-distribution-only evaluation and a closed-form instance is given in Sections E.7 and E.9.

Gauge invariance. The linear reliance g is not invariant under prompt-only reward shifts $\tilde{R}(x, y) \mapsto \tilde{R}(x, y) + b(x)$, while the KL-regularized optimum is. This matters for PPO-style reward whitening [42], where reimplementations of the same mitigation can disagree on g_i . Section E.6 gives a gauge-invariant g_{cent} and verifies that our taxonomy transfers.

Scale invariance. Applying M_i also changes $\|\tilde{R}\|_{L^2(\mu_{\text{diag}})}$, and KL-regularized policies are not invariant to reward scale [72], so this scale change interacts with optimization non-trivially. Section E.5 gives the corrected scale identity and a scale-invariant variant M_i^{norm} .

A.3 Exploitation Shift and Bias Substitution

The measurement-versus-optimization gap admits qualitatively distinct failure modes in the optimized policy. With $\pi = \pi_{\beta}^*(\tilde{R})$ and $\pi' = \pi_{\beta}^*(\tilde{R}')$ for the corresponding pre- and post-mitigation KL-regularized optimal policies, we classify outcomes of a single-axis mitigation M_i along two axes: the true reward difference at the KL-regularized optimum,

$$\Delta J := J(\pi', R) - J(\pi, R),$$

and whether optimization pressure rotates onto another feature to create a change.

Definition A.7 (Feature exploitation gap). *The change in policy-induced expectation of feature ϕ_j between policies π, π' is*

$$\Delta_j(\pi, \pi') = \mathbb{E}_{\mu_{\pi'}}[\phi_j] - \mathbb{E}_{\mu_{\pi}}[\phi_j]. \quad (6)$$

Crossing these two axes yields five mechanistically distinct regimes (R0, R0_{cont}, R1, R2, R3), see Section E.3 for cell-by-cell verification. We also generalize to ε -banded versions of Δ_j and ΔJ in Section E.8 for empirical studies (finite samples, uncertainties). The R0–R3 regimes classify mitigations by quantities depending on R and on Φ_{sp} , neither of which is fully identifiable from preference data alone. Rotation conditions reference only Φ_{sp} since rotation onto structurally relevant features is not penalized.

Throughout, fix $\beta > 0$, let \tilde{R} be a proxy reward satisfying the regularity Assumption E.1, let $\phi_i \in \Phi_{\text{sp}}$ be the targeted feature, and let M_i be a single-axis mitigation (Definition A.5) targeting ϕ_i , with $\tilde{R}' = M_i(\tilde{R})$ likewise satisfying Assumption E.1.

¹The canonical M_i can also be seen as the population Frisch-Waugh-Lovell partial-out of ϕ_i from \tilde{R} at μ_{diag} .

Definition A.8 (R0 Successful and contaminated mitigation). M_i is in regime **R0** if

$$(i) \quad \Delta_j(\pi, \pi') = 0 \quad \text{for all } \phi_j \in \Phi_{sp} \setminus \{\phi_i\}, \quad (7)$$

$$(ii) \quad \Delta J > 0. \quad (8)$$

If (ii) holds but (i) is replaced by $\Delta_j(\pi, \pi') \neq 0$, M_i is in regime **R0_{cont}** (contaminated success).

In R0, mitigation strictly improves true reward at the optimized policy without redirecting first-moment exploitation from feature ϕ_i onto another spurious feature. In R0_{cont}, true reward improves despite first-moment redirection onto another spurious axis. The substitution mechanism is active but its true-reward cost is dominated by gains from reducing ϕ_i . R0_{cont} is the case most easily missed by audit-distribution evaluation, which registers improvement but not rotation.

Definition A.9 (R1 Bias substitution). M_i is in regime **R1** if

$$(i) \quad \Delta_j(\pi, \pi') \neq 0 \quad \text{for some } \phi_j \in \Phi_{sp} \setminus \{\phi_i\}, \quad (9)$$

$$(ii) \quad \Delta J \leq 0. \quad (10)$$

We distinguish the sub-cases neutral substitution ($\Delta J = 0$) and harmful substitution ($\Delta J < 0$).

R1 is the regime enabled specifically by the measurement-versus-optimization gap, as $|g_i(M_i(\tilde{R}); \mu_{\text{diag}})| = 0$ at the audit distribution by construction (see Lemma A.6), while $\Delta J \leq 0$ at the optimized policy. Audit-distribution diagnostics that target ϕ_i therefore register apparent success regardless of the sign of ΔJ .

Definition A.10 (R2 Overcorrection). M_i is in regime **R2** if

$$(i) \quad \Delta_j(\pi, \pi') = 0 \quad \text{for all } \phi_j \in \Phi_{sp} \setminus \{\phi_i\}, \quad (11)$$

$$(ii) \quad \Delta J < 0. \quad (12)$$

In R2, true reward strictly degrades without substitution onto another spurious feature. This effect can be caused by either pushing the targeted feature past zero at the optimization distribution (scale overshoot, fixable by rescaling) or removing a genuinely informative component of a partially informative feature (target misspecification, generically not fixable by rescaling). First-moment ϕ_i observables do not separate the two origins, so they share a regime label. Section E.12 distinguishes them via $\Delta J(c)$ across partial mitigations Mc_i .

Definition A.11 (R3 Silent non-op). M_i is in regime **R3** if

$$(i) \quad \Delta_j(\pi, \pi') = 0 \quad \text{for all } \phi_j \in \Phi_{sp} \setminus \{\phi_i\}, \quad (13)$$

$$(ii) \quad \Delta J = 0. \quad (14)$$

In R3, the mitigation altered the proxy in ways the audit distribution may register but the optimized policy does not express. R3 is a generic outcome of projection-type mitigations whenever the targeted feature's optimization footprint was already small, in which case $D_{\text{KL}}(\pi' \parallel \pi)$ is correspondingly small at the relevant β .

The R0–R3 classification fixes a single audit distribution μ_{diag} , but cannot cover the failure mode where the regime label changes when μ_{diag} does, e.g., across human and LLM-judge audits. Thus:

Definition A.12 (R4 Audit-distribution sensitivity). *A mitigation construction M_i exhibits audit-distribution sensitivity if, holding $(\tilde{R}, R, \Phi_{sp}, \beta)$ fixed, there exist $\mu_{diag}^{(1)} \neq \mu_{diag}^{(2)}$ such that*

$$\pi_{\beta}^*(M_i^{(1)}(\tilde{R})) \quad \text{and} \quad \pi_{\beta}^*(M_i^{(2)}(\tilde{R})) \quad \text{fall into different R0–R3 regimes, where } M_i^{(\ell)} := M_i^{\mu_{diag}^{(\ell)}}.$$

R4 is regime-labeled (rather than treated as a property of the deployment setting, like π_{ref} -sensitivity) because μ_{diag} is a designer-controlled input to the mitigation pipeline² and formalizes the observation that failure-mode classification is a property of $(R, \tilde{R}, M_i, \mu_{diag})$. See Figure 2 for example transition.

Definitions A.8–A.12 cover the mechanistic outcome space, with closed-form instantiations of all regimes given in the linear-Gaussian and quadratic-nonlinear settings of Sections E.7 and E.9. The R0–R3 conditions reference only Δ_j and ΔJ , both invariant under the prompt-only shifts. The canonical M_i is not gauge-invariant, but M_i^{cent} (Section E.6) restores classification invariance.

Example Instances. R1 is directly instantiated by the controlled multi-shortcut demonstration in supervised vision [48], where mitigating one labeled shortcut amplifies the unlabeled at fixed optimization distribution. Our GRPO experiment in Section F.1 provides an in-domain harmful-R1 instance under language RLHF, where a length penalty drives the policy into overconfidence while free-form factual accuracy falls. R2 instantiations include length-penalized reward modeling with accuracy loss [3] and capability degradation in DPO bias-mitigation variants [22]. R4 appears in the systematic human-versus-LLM-judge verbosity gap [6, 65]. The length-bias sign flip in SOTA reward models [21] places in R2 on the targeted axis, or R1 under the conservative reading with off-target axes unmeasured. Section E.9 instantiates all regimes and subregimes in closed form under quadratic non-linearity.

A.4 Impossibility and Sufficiency for Audit-Distribution Evaluation

The measurement-versus-optimization gap raises the question whether any benchmark functional defined at the audit distribution can distinguish successful mitigation (R0) from contaminated success ($R0_{cont}$), substitution (R1), or overcorrection (R2). We prove that they cannot and show that augmenting evaluation with policy-induced distributions suffices (see also Sections E.10 and E.11).

Benchmark class. Let \mathcal{B} denote the class of benchmark functionals depending only on the joint distribution of $(\tilde{R}, M_i(\tilde{R}), R, \Phi)$ under μ_{diag} , with μ_{diag} taken as the empirical evaluation measure. This class subsumes ranking accuracy, pairwise win-rate, preference-prediction calibration, reward–target correlation, the linear-reliance statistic g of Definition A.4, and any composite of these. RewardBench [43], RewardBench2’s headline accuracy [52], AlpacaEval [18], and Chatbot Arena [12] all lie in \mathcal{B} . Including R as an input strengthens our impossibility result, since R appears identically across the four instances and standard benchmarks lack oracle access to it.

THEOREM A.13 (AUDIT-DISTRIBUTION INSUFFICIENCY). *Fix $\beta > 0$. There exist four choices of π_{ref} with $\mu_{diag}, \tilde{R}, M_i, R, \Phi, \Phi_{sp}$, and β held fixed s.t.*

- (i) *the joint distribution of $(\tilde{R}, M_i(\tilde{R}), R, \Phi)$ under μ_{diag} is identical across the four instances, so B takes the same value on all four for every $B \in \mathcal{B}$;*
- (ii) *the optimized policies $\pi_{\beta}^*(M_i(\tilde{R}))$ fall into regimes R0, $R0_{cont}$, R1, and R2 respectively, in the sense of Definitions A.8–A.10.*

² μ_{diag} is non-performative in the sense of Perdomo et al. [61].

The proof (Section E.10) constructs four reference policies in a linear-Gaussian setting with fixed Gram, so audit-side observables coincide while policy-side first moments place the optimized policies in distinct regimes. Since only π_{ref} varies, the same reward model can land in any of R0, R0_{cont}, R1, R2 depending on the reference policy it is deployed against, matching the finding of Malik et al. [52] that high audit scores can fail to transfer to PPO under RM-policy lineage mismatch.

Theorem A.13 establishes the *distributional blindspot*. A second, independent blindspot, the *functional blindspot*, affects ordinal benchmarks: $\mathcal{B}_{\text{ord}} \subseteq \mathcal{B}$ (functionals invariant under monotone transformations of \tilde{R}) is blind to cardinal-scale shifts that KL-regularized optimization tracks (Corollary E.4). Our experiments in Section F.2 realizes this on a non-linear operator using the post-hoc length calibration of Huang et al. [30]. We find that it drives pooled reward-length correlation near zero ($0.316 \rightarrow 0.037$ across five RMs), while within-prompt correlation, which BoN top-1 selection responds to, lands at 0.116 with sign flips on three of four SOTA RMs and $\Delta J < -\epsilon_j$ on two cells.

We show that augmenting the benchmark input with the policy-induced distributions provably suffices. Let \mathcal{B}^+ denote the class of benchmark functionals depending on the joint distribution of $(\tilde{R}, M_i(\tilde{R}), R, \Phi)$ under each of $\mu_{\text{diag}}, \mu_{\pi_\beta^*}(\tilde{R}), \mu_{\pi_\beta^*}(M_i(\tilde{R}))$.

THEOREM A.14 (AUDIT-DISTRIBUTION SUFFICIENCY). Fix $\beta > 0$ and let $\pi = \pi_\beta^*(\tilde{R}), \pi' = \pi_\beta^*(M_i(\tilde{R}))$. With Δ_j and ΔJ as in Definition A.7 and Section A.3, define $B^* \in \mathcal{B}^+$ by

$$B^* = \begin{cases} \text{R0} & \text{if } \Delta J > 0 \text{ and } \Delta_j = 0 \text{ for all } \phi_j \in \Phi_{\text{sp}} \setminus \{\phi_i\}, \\ \text{R0}_{\text{cont}} & \text{if } \Delta J > 0 \text{ and } \Delta_j \neq 0 \text{ for some } \phi_j \in \Phi_{\text{sp}} \setminus \{\phi_i\}, \\ \text{R1} & \text{if } \Delta J \leq 0 \text{ and } \Delta_j \neq 0 \text{ for some } \phi_j \in \Phi_{\text{sp}} \setminus \{\phi_i\}, \\ \text{R2} & \text{if } \Delta J < 0 \text{ and } \Delta_j = 0 \text{ for all } \phi_j \in \Phi_{\text{sp}} \setminus \{\phi_i\}. \end{cases} \quad (15)$$

Under Assumptions E.1, A.2, and E.2, with $M_i(\tilde{R})$ likewise satisfying Assumption E.1, for every tuple $(\pi_{\text{ref}}, \tilde{R}, M_i, R, \Phi, \Phi_{\text{sp}})$ with $(\pi, \pi') \in \text{R0} \cup \text{R0}_{\text{cont}} \cup \text{R1} \cup \text{R2}$ in the sense of Definitions A.8–A.10, B^* returns the correct regime label.

Proof (Section E.11): Δ_j and ΔJ are first moments under μ_π and $\mu_{\pi'}$, both in \mathcal{B}^+ . The four branches are mutually exclusive by Definitions A.8–A.10, where the sign of ΔJ separates R0 and R0_{cont} from R1 and R2, and the rotation condition separates within each pair. A fifth branch ($\Delta J = 0 \wedge \Delta_j = 0$) extends the classifier to R3, with $D_{\text{KL}}(\pi' \parallel \pi)$ separating policy-relevant from policy-irrelevant R3. A finite-sample version under noise-floor ϵ -bands follows from the same construction (Section E.8).

What this pair establishes. Theorems A.13 and A.14 tell us that audit-only inputs cannot separate R0, R0_{cont}, R1, and R2, **even when R itself is granted as an oracle input**. Augmenting with policy-induced distributions and the spurious-feature partition closes the gap.

A.5 Implications for Mitigations Methods and Evaluations

We translate the findings of Theorems A.13 and A.14 into prescriptions for new reward bias mitigation works. We state the corresponding prescriptions for benchmark developers, the multi-axis recommendations, and the operator-variant caveats ($M_i^{\text{norm}}, M_i^{\text{cent}}$) in Section E.12.

Method paper prescriptions. A paper claiming successful reward bias mitigation should:

- (1) *Evaluate at policy-induced distributions and report* (Δ_j, Δ_j) . By Theorem A.14 this input separates R_0 , $R_{0\text{cont}}$, R_1 , and R_2 , as no audit-only input does (Theorem A.13). Concretely: run the mitigated reward model in BoN or short PPO against a fixed reference policy.
- (2) *Instrument off-target Δ_j on every $\phi_j \in \Phi_{\text{sp}} \setminus \{\phi_i\}$ at μ_{π^*}* . Reporting only on-target $|g_i| \approx 0$ leaves R_0 structurally indistinguishable from $R_{0\text{cont}}$ and R_1 , regardless of audit-side strength.
- (3) *Report cardinal scale $\|\tilde{R}\|_{L^2(\mu_{\text{diag}})}$ pre/post mitigation*. By Corollary E.4, ordinal scores are invariant under $\tilde{R} \mapsto c\tilde{R}$ ($c > 0$) while $J(\pi_{\beta^*}(c\tilde{R}), R)$ varies with c . Since mitigation induces rescaling (Section E.5), ordinal-only evaluation misses scale-driven regime shifts.
- (4) *Document π_{ref} -sensitivity*. The construction of Theorem A.13 produces four distinct regimes by varying only π_{ref} for fixed $(\tilde{R}, M_i, \mu_{\text{diag}})$. Validate across different reference policies.

Concrete detection procedures for the regime-specific gaps not closed by B^* are in Section E.12.

Joint adoption. R_0 certification is a joint property of the mitigation operator and the evaluation protocol. Method papers cannot unilaterally provide policy-distribution evaluation without supporting benchmark infrastructure and benchmarks cannot unilaterally provide off-target Δ_j measurement without method papers specifying Φ_{sp} . Joint adoption moves the field from R_0 claims audit-side scores cannot certify (Theorem A.13) to R_0 claims this framework certifies sufficient (Theorem A.14). Extended prescriptions, including benchmark developer recommendations and multi-axis operator guidance, are in Section E.12. We map published mitigation work onto these regimes in Section B.

B Bias Substitution in the Wild

Bias substitution is already pervasive in the published RLHF/DPO mitigation literature, miscategorized as isolated anomalies, capability trade-offs, or judge noise because no prior framework distinguishes it from successful mitigation. Applying the prescriptions of Section A.5 to the published mitigation literature, we find that almost no work provides the inputs Theorem A.14 requires, and where partial evidence is available it points to $R_{0\text{cont}}$, R_1 , or R_2 rather than R_0 .

B.1 Reward models exhibit correlated, partially informative spurious features

Bias substitution requires correlated, partially informative spurious features under μ_{diag} , both hold in deployed reward models (Fact G.1–G.17, G.20). Partial informativeness anchors R_2 's target-misspecification origin, and correlated axes provide directions for pressure rotation. We extend this with length-sycophancy coupling measurements across labeling regimes (Fact G.9, see Section F.3). Across LLM responses to Reddit prompts from Cheng et al. [9] on eight model families, the sycophantic effect on length is +24.3 characters under human labels, +128.4 under LLM judge labels, +154.3 under judge agreement, and −43.1 under judge disagreement (all $p < 0.01$). The sign reversal documents R_4 's audit-dependent operator construction empirically. Section E.9 isolates the same mechanism in closed-form phase diagrams of regime transitions driven by these couplings.

B.2 Bias substitution arises under RLHF training

To show that bias substitution arises end to end under policy optimization, we run GRPO [67] on *Llama-3.2-3B-Instruct* [27] with the *Skywork-Reward-V2-Llama-3.1-8B* reward model [49], shaping the reward as $\tilde{R}(x, y) = R_{\text{RM}}(x, y) - \lambda n_{\text{tok}}(y)/100$ for $\lambda \in \{0, 4, 8\}$ with four seeds per value. This shaping has the form of the canonical single-axis mitigation M_i of Definition A.5, with a fixed coefficient in place of the estimated reliance. As λ rises, mean response length falls from 204 to 170 tokens and multiple-choice (MMLU) accuracy is unchanged, yet expected calibration error climbs from

0.25 to 0.41, free-form factual accuracy on TriviaQA falls from 0.56 to 0.42, and confidence-correctness AUROC drops from 0.73 to 0.65 (Figure 1, Table 3). The $\lambda = 0$ run leaves calibration intact, so the degradation is caused by the penalty and not by RLHF itself. Optimization pressure rotates off length onto expressed confidence and the true reward proxy of free-form factual accuracy degrades, instantiating harmful R1 in a standard RLHF pipeline. Full setup, training curves, and per-metric results are in Section F.1.

B.3 Bias mitigations that pass audit can fail under optimization

We show that the measurement-versus-optimization gap fails to close with previously published mitigation operators in Section F.2. We evaluate the post-hoc LOESS calibration of Huang et al. [30] and the linear-probe operator of Fein et al. [21] across five reward models under Best-of-N selection. The calibration drives pooled reward-length correlation from 0.316 to 0.037, an audit-side success by construction. All four SOTA reward models acquire negative within-prompt correlations under the calibration, with three exceeding their unmitigated baselines in absolute value. AlpacaEval LC win rate degrades below baseline on two cells and GSM8K BoN accuracy drops 3.6 points on one. Under the ϵ -banded reading of Section E.8, two cells satisfy $R2_\epsilon$ on the targeted axis, or mixed $R1_\epsilon + R2_\epsilon$ with off-target axes unmeasured, both undetectable to any $B \in \mathcal{B}$ by Theorem A.13.

B.4 Published mitigations cannot certify successful mitigations

No published mitigation paper we find provides evidence sufficient to certify R0. Where direct evidence is available, it points to $R0_{\text{cont}}$, R1, or R2 and for the remainder the regime is undetermined, which is itself the failure mode Theorem A.13 formalizes. We discuss each work in detail in Section H.

Determined regimes. R2 appears in Feng et al. [22], Bu et al. [3], and Huang et al. [30]. R1 is determinately instantiated by the controlled-vision study of Li et al. [48] and, in language RLHF, by our GRPO experiment (Section B.2 and Figure 1). The human-versus-LLM-judge verbosity gap [6, 65] and our Section B.1 sign reversal document R4’s audit-dependent operator construction in deployed reward models. R3 has no clean published instance we can identify, as expected from publication bias. Section E.9 additionally instantiates R0, $R0_{\text{cont}}$, R1 (neutral and harmful), R2, and the operator-level R4 transition in closed form, with finite- N validation.

Split or undetermined. Fein et al. [21], Eisenstein et al. [19], Meng et al. [53], Lu et al. [51], Chen et al. [8], Ye et al. [86], Zhao et al. [91], Liu et al. [50], Fu et al. [24], and Kim et al. [36] all admit multiple regime readings on the evidence reported, ranging across $R0_{\text{cont}}$, R1, R2, and R3 depending on unmeasured Δ_j at μ_{π^*} or unreported rescalability.

$R0_{\text{cont}}$ on audit-side evidence alone. Park et al. [59], Shen et al. [70], Li et al. [45], Wang et al. [80], Ng et al. [55], and Li et al. [47] zero an audit-side metric without measuring off-target Δ_j at μ_{π^*} , leaving R0 indistinguishable from substitution by Theorem A.13.

The closest approaches. Fein et al. [21] survives the within-prompt diagnostic of Section F.2 with $\Delta_j > 0$ on four of five RMs but does not measure off-target Δ_j . Srivastava et al. [75] engages multi-axis mitigation, BoN, and transformation robustness, but its LLM-oracle counterfactuals inherit unmeasured R4 sensitivity. Cai et al. [4] includes PPO results but reports no off-target Δ_j . None jointly delivers the inputs Theorem A.14 requires. These regime ambiguities are not resolvable by current benchmarks. Every benchmark we surveyed [12, 43, 52, 84] sits inside the impossibility class \mathcal{B} of Theorem A.13, see Section E.12.

C Related Work

Reward bias mitigation methods and benchmarks. We position our work against prior mitigation methods and benchmarks in Section B and Section H; our framework classifies what single-axis operators can certify rather than proposing a new one. Concurrent work addresses adjacent problems: probabilistic reward modeling [7], identifiability under correlated probit models [11], emergent misalignment from narrow SFT [2], heterogeneous safety drift under benign fine-tuning that standard proxies fail to predict [35], and a KL-minimal sycophancy correction [68]. None formalize how single-axis mitigations redistribute optimization pressure onto correlated proxies.

Supervised-learning and fairness antecedents. The rotation-onto-correlated-proxies mechanism has supervised-learning precedents in shortcut learning [26], simplicity bias [66], and Whac-A-Mole [48], with a structural analog in fairness gerrymandering [34]. Our contribution provides the first formal taxonomy of these failures, with KL-regularized RLHF as the primary instantiation. Two failure modes absent under empirical risk minimization arise: cardinal-scale sensitivity (Corollary E.4) and π_{ref} dependence (Theorem A.13).

Prior definitions of reward hacking. Existing definitions [23, 40, 41, 73] take a single proxy as given and ask whether optimizing it produces an acceptable policy, with Liu et al. [50] extending to the worst-case proxy in a correlation set. None consider mitigation operators, so prior definitions cannot distinguish a mitigation that reduces reward hacking from one that relocates it, overshoots, or does nothing the optimizer notices.

D Discussion

Limitations. The regime classification uses first-moment drift, which can miss tail shifts for bounded features such as sycophancy indicators. We argue that reward model biases admit to first order corrections to limit the practical severity of this gap (Section G). Definition A.3 partitions Φ into spurious and structurally relevant features rather than decomposing within-feature. However, our approach is conservative under partial informativeness, as it under-flags substitution rather than over-flagging clean mitigations (Section E.3). Also, as soon as anyone targets a feature for debiasing, they assert it belongs to Φ_{sp} , so our taxonomy and impossibility result apply regardless of whether Φ_{sp} is ground-truth identifiable from preference data.

Impact. Current improvement metrics for reward bias mitigation cannot distinguish whether a model has become less reward-hacky from whether it has merely shifted exploitation modes. Our taxonomy and impossibility-sufficiency pair apply to any preference-learning setting with single-axis mitigation operators. We develop language RLHF in detail because reward bias mitigation is a primary lever for shaping non-verifiable model behavior, and the mitigation ambiguity propagates directly into post-deployment alignment claims. We prove this blindspot is structural, as no benchmark functional over audit-distribution observables can distinguish successful mitigation from substitution or overcorrection, regardless of benchmark richness. We establish this both mathematically and empirically across published SOTA reward models and mitigation operators, showing that bias substitution is an active and pervasive failure mode in RLHF systems. To close the gap, we provide a concrete, implementable checklist for both mitigation method developers and benchmark developers.

885 E Additional Formalizations

886 E.1 Causal Interpretation and Scope

888 In this appendix, we outline the causal frame of our formalization in Section A and delimit what the framework does
889 and does not claim about the underlying causal structure of reward.

891 *Underlying structural causal model.* We posit, descriptively, that the true reward R and the features Φ are related
892 through an underlying structural causal model \mathcal{M} on (X, Y, R, Φ) . We do not assume that \mathcal{M} is identified from preference
893 data \mathcal{D} . The reward identifiability characterization of Skalse et al. [72] establishes that preference data identifies R
894 only up to a prompt-only additive shift under the Bradley-Terry likelihood, leaving the cardinal scale of R and by
895 extension the causal structure of \mathcal{M} partially identified at best. The framework of Section A therefore operates within
896 this partial-identification regime and it adopts causal vocabulary to state failure modes precisely, but its definitions and
897 results do not require \mathcal{M} to be recovered from \mathcal{D} .

901 *Causal reading of Definition A.3.* The corresponding causal reading of Definition A.3 is: ϕ_i is *causally spurious* with
902 respect to R iff ϕ_i has no directed path to R in \mathcal{M} . Under positivity-style richness on the support of μ_{diag} (strengthening
903 Assumption E.2) and faithfulness, the causal and observational readings coincide on the partition Φ_{sp} vs. Φ_{struct} . The
904 causal reading is the primitive notion and the observational form its identifiable shadow. Thus, Φ_{sp} vs. Φ_{struct} is best read
905 as the conservative observational projection of an underlying causal partition preference data does not fully identify
906 [72].

909 *Partial informativeness as causal mediation.* Partial informativeness (Sections A.1 and E.3) corresponds causally to ϕ_i
910 having both a direct path to R in \mathcal{M} and one or more mediated paths through other features. Length is the canonical
911 instance: a direct contribution via comprehensiveness, plus mediated effects of hedging, sycophancy, and verbosity
912 compensation (see Sections B and G). The conservative partition classifies such features as Φ_{sp} , under-flagging rotation
913 onto causally relevant mediated paths.

916 *Associational mitigation operators.* The single-axis mitigations M_i of Definition A.5 are associational: they target
917 the linear reliance statistic g_i , which is an $L^2(\mu_{\text{diag}})$ regression coefficient, not a controlled direct effect. A causal
918 counterpart would target the controlled direct effect of ϕ_i on R , identifiable only under interventions on \mathcal{M} or strong
919 conditional-independence assumptions that preference data does not warrant. Existing single-axis mitigations in the
920 language reward modeling literature [e.g., 21, 30, 58] are associational in this sense, and the regime taxonomy of
921 Section A.3 characterizes their failure modes precisely because of the gap between associational construction and causal
922 response of the optimized policy.

926 *R2 origins in causal language.* The two origins of R2 in Definition A.10 restate causally as: scale overshoot (projection
927 coefficient too large at μ_{π^*} while ϕ_i remains causally spurious) versus target misspecification (ϕ_i has a non-zero direct
928 path to R in \mathcal{M} that the projection has zeroed at μ_{diag}). The rescalability test of Section A.5 discriminates these origins
929 observationally, without identifying \mathcal{M} .

932 *R4 and transportability.* The audit-distribution sensitivity of Definition A.12 is structurally analogous to the trans-
933 portability question [e.g. 60]. An identified mitigation at one audit distribution does not need to transport to another.
934 R4 names this phenomenon under partial identification, without formally invoking transportability machinery.

E.2 Standing Assumptions and Optimization Setup

Assumption E.1 (Regularity). $R, \tilde{R} \in L^2(\mu_{\text{diag}}) \cap L^\infty(\mu_{\pi_{\text{ref}}})$ and $\phi_k \in L^2(\mu_{\text{diag}}) \cap L^\infty(\mu_{\pi_{\text{ref}}})$ for each $k = 1, \dots, K$, where $\Phi = (\phi_1, \dots, \phi_K)^\top$ is the feature map of Definition A.1.

The $L^2(\mu_{\text{diag}})$ inclusion makes $g(\cdot; \mu_{\text{diag}})$ of Definition A.4 well-defined as a Gram-inner-product statistic at the diagnostic measure. The $L^\infty(\mu_{\pi_{\text{ref}}})$ conditions are natural both for bounded-output reward models and for the interpretable surface features this paper considers (length, formatting counts, style indicators). The conditions imply that any \tilde{R} obtained as a finite \mathbb{R} -linear combination of \tilde{R} , R , and the ϕ_k also lies in $L^\infty(\mu_{\pi_{\text{ref}}})$, which covers every \tilde{R} used in this paper, including $M_i(\tilde{R})$ of Definition A.5 and the scale-normalised variant M_i^{norm} introduced below. Since $\mu_{\pi_{\text{ref}}}$ is a probability measure, $L^\infty(\mu_{\pi_{\text{ref}}}) \subset L^2(\mu_{\pi_{\text{ref}}})$ holds automatically and no separate $L^2(\mu_{\pi_{\text{ref}}})$ clause is needed.

Optimization setup. Fix $\beta > 0$. For any \tilde{R} satisfying Assumption E.1, define the KL-regularized optimizer

$$\pi_\beta^\star(\tilde{R}) = \arg \max_{\pi} \mathbb{E}_{x \sim \mathcal{D}, y \sim \pi} [\tilde{R}(x, y)] - \beta D_{\text{KL}}(\pi \parallel \pi_{\text{ref}}),$$

which in the single-turn contextual-bandit regime admits the unique softmax closed form $\pi_\beta^\star(\tilde{R})(y \mid x) \propto \pi_{\text{ref}}(y \mid x) \exp(\tilde{R}(x, y)/\beta)$ recalled in the main text. Essential boundedness of \tilde{R} gives two-sided bounds on both $\exp(\tilde{R}/\beta)$ and the normalizer $Z(x) = \mathbb{E}_{y \sim \pi_{\text{ref}}(\cdot \mid x)} [\exp(\tilde{R}(x, y)/\beta)]$, so the density ratio $d\mu_{\pi_\beta^\star(\tilde{R})}/d\mu_{\pi_{\text{ref}}} = \exp(\tilde{R}/\beta)/Z(x)$ is essentially bounded. The chain $\tilde{R} \in L^\infty(\mu_{\pi_{\text{ref}}}) \Rightarrow L^\infty(\mu_{\pi_\beta^\star(\tilde{R})}) \subset L^2(\mu_{\pi_\beta^\star(\tilde{R})})$ then follows immediately, and analogously for each R and each ϕ_k . The policy-level expectations $\mathbb{E}_{\mu_{\pi_\beta^\star}(\tilde{R})}[\phi_j]$ entering Definition A.7 are then finite. Well-definedness of $g(\tilde{R}; \mu_{\pi_\beta^\star(\tilde{R})})$ additionally requires Gram non-degeneracy at $\mu_{\pi_\beta^\star}$. This condition is handled in Section E.3, where the relevant audit distributions are introduced.

E.3 Feature Realizability and the Diagnostic Measure

This section collects four items referenced from Section A. A discussion of the choices available for μ_{diag} (pointed to from the preamble of Section A), the Gram non-degeneracy condition for the decoupled case, the structural Assumption E.2 (feature realizability) on the feature map (invoked by Definition A.3), and the expressiveness remark that motivates reading the regime language of Section A.3 conservatively under partially informative features (pointed to immediately after Definition A.3). The order below matches the narrative order of Section A.

Remark (choice of diagnostic measure). The main text introduces μ_{diag} as an auxiliary distribution on $\mathcal{X} \times \mathcal{Y}$ used for reliance estimation and correlation measurement, distinct from the KL anchor π_{ref} . Two choices are load-bearing for this paper:

- (1) *The coupled case* $\mu_{\text{diag}} = \mu_{\pi_{\text{ref}}}$. Recovers the setting of Laidlaw et al. [41] and makes diagnostic statistics directly comparable to their bounds, at the cost of coupling measurement and optimization reference points.
- (2) *An annotator-conditioned distribution* $\mu_{\text{diag}}^{\text{human}}$ or $\mu_{\text{diag}}^{\text{LLM}}$, reflecting that the same proxy \tilde{R} induces different g -profiles depending on which annotator pool generated the audit data. This formalizes the human-vs-LLM-judge gap documented in Chen et al. [6], Movva et al. [54], Saito et al. [65], Zheng et al. [92], and is what R4 (Definition A.12) operationalizes.

Other valid choices, e.g., the preference-data distribution used to train \tilde{R} , curated audit sets, deployment distributions, and round-anchored distributions in iterative RLHF [87], yield different g -profiles and shift substitution diagnostics accordingly. Substitution claims are always relative to the chosen diagnostic measure, and the choice should be stated explicitly in applications.

Where $g(\bar{R}; \mu_{\pi^*}(\bar{R}))$ is referenced in Section A (e.g., the measurement-vs-optimization gap following Lemma A.6), well-definedness requires Gram non-degeneracy at μ_{π^*} , which follows from Assumption A.2 (non-degeneracy) whenever $d\mu_{\pi^*}/d\mu_{\text{diag}}$ is bounded above and below. In the coupled case ($\mu_{\text{diag}} = \mu_{\pi_{\text{ref}}}$), the bound follows from Assumption E.1 (regularity). In the decoupled case ($\mu_{\text{diag}}^{\text{human}}, \mu_{\text{diag}}^{\text{LLM}}$), it is a separate regularity condition assumed in the relevant statements.

Assumption E.2 (Feature Realizability). *For μ_{diag} -almost every $(x, y) \in \mathcal{X} \times \mathcal{Y}$ and every feature index $i \in \{1, \dots, K\}$, there exists at least one $y' \in \mathcal{Y}$ with $\phi_j(x, y') = \phi_j(x, y)$ for all $j \neq i$ and $\phi_i(x, y') \neq \phi_i(x, y)$.*

Assumption E.2 is the non-vacuity condition for Definition A.3. Without it, ϕ_i can be functionally determined by the prompt and the other features on \mathcal{Y} , so further conditioning on ϕ_i leaves the conditional expectation in Definition A.3 unchanged and every feature is trivially classified as spurious. Assumption E.2 is therefore a richness condition on the response space under μ_{diag} , as \mathcal{Y} must be rich enough where the diagnostic measure places mass that each feature axis admits independent perturbation.

We treat g and Δ_j as observational statistics under μ_{diag} and μ_{π^*} throughout, as an interventional reading along ϕ_i would require strengthening Assumption E.2 to a richness condition on counterfactual realization, which we do not invoke.

Even after weakening to the μ_{diag} -a.e. form, Assumption E.2 is fragile for natural-language features on Φ 's own terms. Consider $\Phi = \{\phi_{\text{length}}, \phi_{\text{exclamation_count}}\}$. Assumption E.2 demands that μ_{diag} -a.e. there exist a response y' matching y on exclamation count but differing in length, and conversely. The first direction is plausible, as long and short variants with the same exclamation count are typically realizable. The second is genuinely constrained, as at fixed length, the realizable range of exclamation counts is bounded above by length itself, and at very short lengths the counterfactual set may be empty. Even within Φ , the response space does not cleanly factor into independent feature axes. This failure mode is specific to features that count items within the response (exclamation marks, bullet points, headers), since their range is mechanically tied to length. Features operationalized as binary indicators or response-level scalars do not have this bound (e.g., sycophancy and confidence in Sections B and F.3), and Assumption E.2 holds approximately for these pairs.

Feature misspecification is itself a major failure mode in practice. Under Assumption E.2, Definition A.3 cleanly separates spurious from structurally relevant features, and the canonical mitigation M_i of Definition A.5 targets a well-defined coordinate of g -space. When Assumption E.2 fails, the definitions of Section A.1 should be read as approximations, and the binary spurious/structurally-relevant partition as the best dichotomous summary of a genuinely continuous informativeness spectrum.

Remark (scope of the binary partition under partial informativeness). Definition A.3 partitions Φ into spurious and structurally relevant features. The partition is scope-matched to the operator class deployed in the literature (whole-feature projections, see Section A.1) and to what preference data identifies, but it does not separate within-feature components for partially informative features in the empirical setting surveyed in Section B, where many natural features (foremost length) act as mediators for multiple mechanisms (sycophancy, epistemic uncertainty, informativeness) rather than as pure spurious proxies or pure structural drivers.

Under partial informativeness, the regime distinction of Section A.3 still applies, but with a scoping caveat: bias substitution becomes a *conservative classification* of the ways a single-axis mitigation can fail, because real features admit additional substitution modes that the binary partition cannot express. Concretely, the best $L^2(\mu_{\text{diag}})$ approximation of ϕ_i by a function of R alone (writing $\phi_i = \mathbb{E}_{\mu_{\text{diag}}}[\phi_i | R] + \phi_i^\perp$, i.e., decomposing ϕ_i into “the part that tracks R ” and

a residual) does not in general coincide with the causal-versus-spurious decomposition of ϕ_i under the underlying structural causal model \mathcal{M} (Section E.1). The first is a μ_{diag} -level regression decomposition, the second a structural statement under \mathcal{M} . Reconciling the two requires additional structure, e.g., a causal model, conditional-independence assumptions, or auxiliary interventions, and we defer partial-informativeness extensions to future work.

Exhaustiveness of the regime taxonomy. The two axes underlying Section A.3 – whether mitigation rotates first-moment exploitation onto another spurious feature ($\Delta_j(\pi, \pi') \neq 0$ for some $\phi_j \in \Phi_{\text{sp}} \setminus \{\phi_i\}$) and the sign of the change in true reward ($\Delta J = J(\pi', R) - J(\pi, R)$) – partition the outcome space into six cells, summarized in Table 1. Five regime labels (R0, R0_{cont}, R1, R2, R3) suffice because R1 absorbs both $\Delta J = 0$ (neutral substitution) and $\Delta J < 0$ (harmful substitution) under Definition A.9, while R0 and R0_{cont} receive separate labels because the rotation-with-improvement corner is precisely the case audit-distribution evaluation most readily mistakes for clean R0 (cf. Theorem A.13). The asymmetry is therefore deliberate rather than ad hoc: R0_{cont} earns a separate label for downstream emphasis in the impossibility result, whereas neutral and harmful R1 are both already failures of ΔJ and the sub-case treatment suffices.

| | No rotation ($\Delta_j = 0$ on $\Phi_{\text{sp}} \setminus \{\phi_i\}$) | Rotation ($\Delta_j \neq 0$ for some $\phi_j \in \Phi_{\text{sp}} \setminus \{\phi_i\}$) |
|----------------|--|---|
| $\Delta J > 0$ | R0 (Definition A.8) | R0 _{cont} (Definition A.8) |
| $\Delta J = 0$ | R3 (Definition A.11) | R1, neutral (Definition A.9) |
| $\Delta J < 0$ | R2 (Definition A.10) | R1, harmful (Definition A.9) |

Table 1. Cell-by-cell partition of single-axis mitigation outcomes by the rotation predicate on $\Phi_{\text{sp}} \setminus \{\phi_i\}$ and the sign of ΔJ . Six cells, five regime labels: R1 covers the two $\Delta J \leq 0$ rotation cells under the neutral/harmful sub-case distinction. R4 (Definition A.12) is transversal to this table and indexes how regime membership changes across audit distributions μ_{diag} .

The R4 axis (Definition A.12) does not appear in Table 1 because R4 is not an outcome cell but a property of how cell membership shifts across μ_{diag} choices, as discussed in Section A.3 and instantiated empirically in Section F.3.

Spurious-rotation-with-improvement corner. Definition A.9 (R1, bias substitution) requires $\Delta J \leq 0$. The complementary corner $\Delta_j \neq 0$ for some $\phi_j \in \Phi_{\text{sp}} \setminus \{\phi_i\}$ with $\Delta J > 0$ is regime R0_{cont} per Definition A.8. Two mechanisms produce R0_{cont}: (i) rotation onto a feature classified spurious in the conservative partition but partially informative in the true reward (an in-scope manifestation of the conservatism discussed above), and (ii) under the strict binary partition, rotation onto a genuinely spurious ϕ_j where structural-feature gains from reducing ϕ_i outweigh the spurious cost on ϕ_j .

The practical consequence for the rest of the paper is that substitution claims proved under Assumption E.2 (feature realizability) should be read as a floor rather than a complete characterization: when Assumption E.2 fails, at least the substitution mode the binary partition captures remains available to the optimizing policy, and generically additional modes are available too.

E.4 Equivalent Epsilon-Formulation for Local Spuriousness

Definition A.3 states spuriousness as conditional mean independence at μ_{diag} . We derive the equivalent $o(\epsilon)$ formulation for completeness.

Let $\mathcal{Y}_{-i}(x, y) := \{y' : \phi_j(x, y') = \phi_j(x, y) \forall j \neq i\}$ denote the i -fiber through (x, y) , and decompose $\mu_{\text{diag}} = D(x) \pi_{\text{audit}}(y | x)$. For a Markov kernel κ with $\kappa(\mathcal{Y}_{-i}(x, y) | x, y) = 1$ for μ_{diag} -a.e. (x, y) , define the mixture-perturbed conditional

$$\pi_{\epsilon}^{\kappa}(y' | x) := (1 - \epsilon) \pi_{\text{audit}}(y' | x) + \epsilon \int \kappa(y' | x, y) \pi_{\text{audit}}(dy | x), \quad \epsilon \in [0, 1]. \quad (16)$$

1093 Consider the perturbation condition

$$1094 \mathbb{E}_{x \sim D, y' \sim \pi_\varepsilon(\cdot|x)}[R(x, y')] - \mathbb{E}_{\mu_{\text{diag}}}[R(x, y)] \in o(\varepsilon) \quad \text{as } \varepsilon \downarrow 0. \quad (17)$$

1096 The left-hand side equals $\varepsilon \cdot (\mathbb{E}_{\kappa \otimes \mu_{\text{diag}}}[R] - \mathbb{E}_{\mu_{\text{diag}}}[R])$ for mixture perturbations, so the $o(\varepsilon)$ condition is equivalent to
 1097 $\mathbb{E}_{\kappa \otimes \mu_{\text{diag}}}[R] = \mathbb{E}_{\mu_{\text{diag}}}[R]$. Requiring this for every kernel κ implies the conditional mean independence of Definition A.3,
 1098 and the two are equivalent when R is $\sigma(\Phi)$ -measurable (so that within-fiber variation of R collapses to variation in ϕ_i).
 1099
 1100

1102 *Verification of downstream results.* Theorem A.13 and the closed-form regime instantiations in Sections E.7 and E.9
 1103 use a true reward of the form $R = w\phi_3$, which is $\sigma(\Phi)$ -measurable. Thus, the conditional mean independence and strict
 1104 point-wise invariance on i -fibers coincide in this regime. The constructed instances satisfy Definition A.3 and the proofs
 1105 transfer without modification. Theorem A.14 is a correctness statement on the classifier B^* given $(\Delta_j, \Delta J)$ and Φ_{sp} as
 1106 inputs. It places no assumption on the functional form of R and transfers regardless of how Φ_{sp} is defined.
 1107
 1108

1111 E.5 Scale Change Induced by Mitigation and a Scale-Invariant Variant

1112 Immediately after Lemma A.6, we note that applying M_i changes $\|\tilde{R}\|_{L^2(\mu_{\text{diag}})}$ through both a diagonal and a cross-
 1113 correlation contribution, and that this scale change interacts non-trivially with the KL-regularized optimizer because
 1114 π_β^* is not invariant to positive rescaling of the reward. In this section, we give the exact scale identity for M_i , isolate the
 1115 mechanism by which scale change conflates with axis reallocation at fixed β , and define the scale-invariant variant
 1116 M_i^{norm} that separates the two effects.
 1117

1119 **Scale identity.** Recall the Gram matrix G of Assumption A.2, with entries $G_{ij} = \mathbb{E}_{\mu_{\text{diag}}}[\phi_i \phi_j]$, let $g = g(\tilde{R}; \mu_{\text{diag}})$
 1120 abbreviate the linear-reliance vector of Definition A.4, and write $\tilde{R}' = M_i(\tilde{R}) = \tilde{R} - g_i \phi_i$ for the canonical mitigation of
 1121 Definition A.5. Expanding the $L^2(\mu_{\text{diag}})$ -norm of \tilde{R}' ,
 1122

$$1123 \|\tilde{R}'\|_{L^2(\mu_{\text{diag}})}^2 = \|\tilde{R}\|_{L^2(\mu_{\text{diag}})}^2 - 2g_i \mathbb{E}_{\mu_{\text{diag}}}[\phi_i \tilde{R}] + g_i^2 G_{ii},$$

1125 and using $\mathbb{E}_{\mu_{\text{diag}}}[\phi_i \tilde{R}] = (Gg)_i = \sum_j G_{ij} g_j$ gives the exact identity

$$1126 \|\tilde{R}'\|_{L^2(\mu_{\text{diag}})}^2 = \|\tilde{R}\|_{L^2(\mu_{\text{diag}})}^2 - g_i \left[g_i G_{ii} + 2 \sum_{j \neq i} G_{ij} g_j \right].$$

1130 The sign of the change is governed by the bracketed expression relative to g_i : the norm strictly decreases whenever the
 1131 bracketed term is co-signed with g_i , and can *increase* when the cross-correlation contribution $2 \sum_{j \neq i} G_{ij} g_j$ is oppositely
 1132 signed and sufficiently large in magnitude relative to the diagonal term $g_i G_{ii}$. A parameter choice exhibiting this norm
 1133 increase is given in Section E.7. The diagonal-only approximation $\|\tilde{R}'\|^2 \approx \|\tilde{R}\|^2 - g_i^2 G_{ii}$ is correct only when G is
 1134 diagonal at μ_{diag} , i.e., when the diagnostic measure renders the feature axes orthogonal. This condition generically fails
 1135 for interpretable feature sets such as length, formatting, and style indicators, which are empirically correlated on any
 1136 natural μ_{diag} .
 1137

1139 **Scale-invariant variant.** The $L^2(\mu_{\text{diag}})$ -normalized mitigation

$$1140 M_i^{\text{norm}}(\tilde{R}) = \alpha (\tilde{R} - g_i \phi_i), \quad \alpha := \frac{\|\tilde{R}\|_{L^2(\mu_{\text{diag}})}}{\|\tilde{R} - g_i \phi_i\|_{L^2(\mu_{\text{diag}})}},$$

preserves $\|\cdot\|_{L^2(\mu_{\text{diag}})}$ by construction. Its linear reliance at μ_{diag} is α times that of $M_i(\tilde{R})$, so $g_i(M_i^{\text{norm}}(\tilde{R}); \mu_{\text{diag}}) = 0$ and $g_j(M_i^{\text{norm}}(\tilde{R}); \mu_{\text{diag}}) = \alpha g_j(\tilde{R}; \mu_{\text{diag}})$ for $j \neq i$: the i -th coordinate is zeroed and the remaining axes are rescaled uniformly.

Effective- β shift. By KL scale-equivariance $\pi_\beta^*(c\tilde{R}) = \pi_{\beta/c}^*(\tilde{R})$ and the relation $M_i^{\text{norm}}(\tilde{R}) = \alpha M_i(\tilde{R})$,

$$\pi_\beta^*(M_i(\tilde{R})) = \pi_{\alpha\beta}^*(M_i^{\text{norm}}(\tilde{R})).$$

The unnormalized comparison $\pi_\beta^*(\tilde{R})$ vs $\pi_\beta^*(M_i(\tilde{R}))$ at fixed β therefore equals the normalized comparison $\pi_\beta^*(\tilde{R})$ vs $\pi_{\alpha\beta}^*(M_i^{\text{norm}}(\tilde{R}))$, differing from the same- β normalized comparison by exactly an effective- β shift from β to $\alpha\beta$. This is a cardinal-distortion mechanism rather than an ordinal one, as two proxies producing identical pairwise rankings on μ_{diag} can differ in $L^2(\mu_{\text{diag}})$ -norm and induce different π_β^* at fixed β . This mechanism is in line with the Section 2 ordinal-vs-cardinal gap, which motivates reporting $\|\tilde{R}\|_{L^2(\mu_{\text{diag}})}$ alongside Δ_j in empirical instantiations.

As a practical consequence, under M_i^{norm} , any observed change in $\mathbb{E}_{\pi_\beta^*}[\phi_j]$ at fixed β is attributable to axis reallocation rather than scale, and Δ_j values can differ in sign between M_i and M_i^{norm} when α is far from 1, classifying the same proxy into different R0-R3 regimes under the two mitigations. Diagnosing R1 versus R3 cleanly therefore calls for M_i^{norm} , since under M_i a nonzero Δ_j admits both reallocation and scale interpretations and requires a β -scan or pre/post norm reporting to disambiguate.

Deployment-time normalization. Standard practices in PPO-based RLHF for stability like per-batch reward whitening [42] apply a sample-dependent shift and scale to \tilde{R} at each training step. In expectation the shift approximates a constant under the current rollout distribution and is therefore gauge-equivalent at the population level (no policy effect, by Section A.2), but the scale component divides \tilde{R} by an estimate of $\text{Std}_{\mu_{\pi_t}}(\tilde{R})$, inducing an effective KL coefficient $\beta \cdot \text{Std}_{\mu_{\pi_t}}(\tilde{R})$ that M_i^{norm} does not control. M_i^{norm} fixes the audit-side $L^2(\mu_{\text{diag}})$ scale, while whitening rescales against the (non-stationary) training-side dispersion. The (M_i, M_i^{norm}) comparison therefore does not exhaust the scale story under whitened RLHF, and empirical instantiations should report training-side reward dispersion alongside $\|\tilde{R}\|_{L^2(\mu_{\text{diag}})}$ when interpreting Δ_j across mitigations.

E.6 Gauge-Invariant Linear Reliance

In this section, we give a gauge-invariant variant of the linear reliance statistic and verify that the Section A.3 regime classification transfers under it. Define the prompt-conditional deviations under μ_{diag} ,

$$\tilde{R}^{\text{cent}}(x, y) = \tilde{R}(x, y) - \mathbb{E}_{\mu_{\text{diag}}}[\tilde{R} | x], \quad \tilde{\Phi}(x, y) = \Phi(x, y) - \mathbb{E}_{\mu_{\text{diag}}}[\Phi | x].$$

Assumption E.3 (Non-degeneracy). *The centered Gram matrix $\tilde{G} := \mathbb{E}_{\mu_{\text{diag}}}[\tilde{\Phi}\tilde{\Phi}^\top] \in \mathbb{R}^{K \times K}$ is positive definite: $\tilde{G} \succ 0$.*

Note that Assumption E.3 does not follow from Assumption A.2, as any feature that is constant in y given x is non-degenerate pooled but vanishes after within-prompt centering. The *centered linear reliance* is

$$g_{\text{cent}}(\tilde{R}; \mu_{\text{diag}}) = (\mathbb{E}_{\mu_{\text{diag}}}[\tilde{\Phi}\tilde{\Phi}^\top])^{-1} \mathbb{E}_{\mu_{\text{diag}}}[\tilde{\Phi}\tilde{R}^{\text{cent}}] \in \mathbb{R}^K.$$

For near-collinear or learned feature sets, g_{cent} should be read as a ridge-regularized estimate.

Gauge invariance. Under $\tilde{R} \mapsto \tilde{R} + b(x)$ for any prompt-only b , $\mathbb{E}[\tilde{R} + b | x] = \mathbb{E}[\tilde{R} | x] + b(x)$, so \tilde{R}^{cent} is unchanged and $\tilde{\Phi}$ does not involve \tilde{R} . Therefore $g_{\text{cent}}(\tilde{R} + b; \mu_{\text{diag}}) = g_{\text{cent}}(\tilde{R}; \mu_{\text{diag}})$. The KL-regularized optimum shares this invariance, so g_{cent} is a reliance statistic on the same gauge equivalence class the policy respects.

Relation to g . The pooled second moments decompose as

$$\mathbb{E}[\Phi\Phi^\top] = \mathbb{E}[\bar{\Phi}\bar{\Phi}^\top] + \mathbb{E}[\mathbb{E}[\Phi|x]\mathbb{E}[\Phi|x]^\top], \quad \mathbb{E}[\Phi\tilde{R}] = \mathbb{E}[\bar{\Phi}\tilde{R}^{\text{cent}}] + \mathbb{E}[\mathbb{E}[\Phi|x]\mathbb{E}[\tilde{R}|x]]$$

so g is a matrix-weighted combination of the within-prompt regression (g_{cent}) and the between-prompt regression of conditional means. The two statistics coincide only in degenerate cases (e.g., vanishing between-prompt covariation) and can differ substantially when Φ carries strong prompt-level structure. This structure is plausible for length, where prompt difficulty drives average response length.

Canonical mitigation and Lemma A.6 analogue. The centered canonical mitigation is

$$M_i^{\text{cent}}(\tilde{R})(x, y) = \tilde{R}(x, y) - g_{\text{cent},i}(\tilde{R}; \mu_{\text{diag}}) \phi_i(x, y).$$

Direct computation in the centered inner product gives $g_{\text{cent},i}(M_i^{\text{cent}}(\tilde{R}); \mu_{\text{diag}}) = 0$ and $g_{\text{cent},j}(M_i^{\text{cent}}(\tilde{R}); \mu_{\text{diag}}) = g_{\text{cent},j}(\tilde{R}; \mu_{\text{diag}})$ for $j \neq i$, using that $\mathbb{E}_{\mu_{\text{diag}}}[\bar{\Phi}\bar{\phi}_i]$ is the i -th column of $\mathbb{E}_{\mu_{\text{diag}}}[\bar{\Phi}\bar{\Phi}^\top]$. The variant $\tilde{R} - g_{\text{cent},i}\bar{\phi}_i$ differs from $M_i^{\text{cent}}(\tilde{R})$ by $g_{\text{cent},i}\mathbb{E}[\phi_i | x]$, a prompt-only function. Given the above-stated gauge invariance, the two variants are policy-equivalent and produce identical g_{cent} .

Framework transfer. Definition A.4 ports to g_{cent} verbatim, and Definition A.5 ports as the requirement $|g_{\text{cent},i}(\tilde{R}'; \mu_{\text{diag}})| < |g_{\text{cent},i}(\tilde{R}; \mu_{\text{diag}})|$. The Lemma A.6 analogue above shows M_i^{cent} qualifies. Definitions A.8–A.11 (R0-R3) port unchanged because their criteria reference policy-induced expectations $\Delta_j(\pi, \pi')$ and ΔJ , not the reliance statistic.

Gauge invariance of the regime classification. The Section A.2 obstruction is that under $\tilde{R} \mapsto \tilde{R} + b(x)$ for prompt-only b , $M_i(\tilde{R} + b)$ differs from $M_i(\tilde{R}) + b$ by $(g_i(\tilde{R}) - g_i(\tilde{R} + b))\phi_i$, so $\pi_\beta^*(M_i(\tilde{R} + b)) \neq \pi_\beta^*(M_i(\tilde{R}))$ and $\Delta_j, \Delta J$ can change. Replacing M_i by M_i^{cent} removes the obstruction. By gauge invariance of g_{cent} , $g_{\text{cent},i}(\tilde{R} + b) = g_{\text{cent},i}(\tilde{R})$, so $M_i^{\text{cent}}(\tilde{R} + b) = M_i^{\text{cent}}(\tilde{R}) + b$. Since the KL-regularized optimum is invariant under prompt-only shifts,

$$\pi_\beta^*(\tilde{R} + b) = \pi_\beta^*(\tilde{R}), \quad \pi_\beta^*(M_i^{\text{cent}}(\tilde{R} + b)) = \pi_\beta^*(M_i^{\text{cent}}(\tilde{R})),$$

hence $\Delta_j(\pi, \pi')$ and ΔJ are gauge-invariant and the R0-R3 regime is unchanged. R4 (see Definition A.12) also ports verbatim, as g_{cent} retains its μ_{diag} -dependence and only the gauge of \tilde{R} is factored out. Thus, $M_i^{\text{cent}, \mu_{\text{diag}}}$ depends on μ_{diag} in the same sense as $M_i^{\mu_{\text{diag}}}$ in Definition A.12 and audit-distribution sensitivity is defined identically.

Note that M_i and M_i^{cent} applied to the same proxy \tilde{R} may still occupy different R0-R3 regimes, as different mitigations induce different policies. However, this difference reflects the choice of reliance estimator as a degree of freedom in the mitigation pipeline, independent of gauge. Auditing a proxy under both statistics localizes the targeted reliance, as a feature on which $|g_i|$ is large but $|g_{\text{cent},i}|$ is small flags reliance that lives in prompt-level structure the optimizing policy averages over rather than acts on.

E.7 Non-Vacuity of the Regime Taxonomy via Linear-Gaussian Instantiation

This appendix exhibits closed-form instantiations of the regimes and gaps introduced in Section A, establishing that the taxonomy is non-vacuous within our framework’s own assumptions. The constructions are not intended as empirical models of RLHF and we show the empirical instantiations in Sections B and F.3.

The construction below satisfies Assumptions E.1 (regularity), A.2 (non-degeneracy of the Gram at μ_{diag}), and E.2 (feature realizability), as well as Gram non-degeneracy at μ_{π^*} . We verify each below.

1249 *Shared construction.* We work with a trivial prompt space (\mathcal{X} a singleton, suppressed in notation) and continuous
 1250 response space $\mathcal{Y} = \mathbb{R}^3$. The feature map is $\Phi(y) = (\phi_1, \phi_2, \phi_3)$ with $\phi_k(y) = y_k$. We fix true and proxy rewards
 1251

$$1252 \quad R(y) = w \phi_3(y), \quad \tilde{R}(y) = a \phi_1(y) + b \phi_2(y) + w \phi_3(y),$$

1253
 1254 with $w > 0$ and $(a, b) \in \mathbb{R}^2$. By Definition A.3, $\Phi_{\text{sp}} = \{\phi_1, \phi_2\}$ and $\Phi_{\text{struct}} = \{\phi_3\}$ by construction. Reference policy and
 1255 diagnostic measure are zero-mean Gaussians on \mathbb{R}^3 ,
 1256

$$1257 \quad \pi_{\text{ref}} = \mathcal{N}(0, I_3), \quad \mu_{\text{diag}} = \mathcal{N}(0, \Sigma),$$

1258 where Σ is positive definite with $\Sigma_{kk} = 1$ and off-diagonal $\Sigma_{12} = \rho \in (-1, 1)$, $\Sigma_{13} = \Sigma_{23} = 0$. Fix the KL parameter $\beta > 0$
 1259 and target the spurious feature $i = 1$ throughout. We treat \mathbb{R}^3 as the formal response space. For the $L^\infty(\mu_{\pi_{\text{ref}}})$ clause of
 1260 Assumption E.1, restrict to a sufficiently large bounded set on which all formulas below hold to arbitrary precision in
 1261 the unrestricted limit.³
 1262
 1263

1264 *Verification of assumptions.* Assumption E.1 holds under truncation as noted. Assumption A.2 holds because
 1265 $\mathbb{E}_{\mu_{\text{diag}}}[\Phi\Phi^\top] = \Sigma \succ 0$. Assumption E.2 holds because $\mathcal{Y} = \mathbb{R}^3$ has full support under any positive-definite Gauss-
 1266 ian, so for μ_{diag} -a.e. y and any i , varying y_i while holding $y_{j \neq i}$ fixed remains in the support. Gram non-degeneracy at
 1267 μ_{π^*} follows from the closed form below, where μ_{π^*} is itself Gaussian with positive-definite covariance.
 1268
 1269

1270 *Closed forms for the optimum.* For any linear reward $\tilde{R}(y) = c^\top y$ with $c \in \mathbb{R}^3$, the KL-regularized optimum against
 1271 π_{ref} satisfies $\pi_\beta^*(\tilde{R}) = \mathcal{N}(c/\beta, I_3)$ by completing the square in the softmax form of Equation (2). Two consequences used
 1272 throughout: (i) the policy-induced first moment is $\mathbb{E}_{\mu_{\pi^*}}[\phi_k] = c_k/\beta$, so $\Delta_j(\pi, \pi') = (c'_j - c_j)/\beta$ for any pair of linear
 1273 rewards with coefficients c, c' ; (ii) the plain return is $J(\pi_\beta^*(\tilde{R}), R) = \sum_k w_k c_k/\beta$ where w_k are the coefficients of R , so
 1274 under our $R = w\phi_3$, only the ϕ_3 -coefficient of the proxy contributes to ΔJ .
 1275
 1276

1277 *Linear reliance and canonical mitigation.* The reliance vector at μ_{diag} is $g(\tilde{R}; \mu_{\text{diag}}) = \Sigma^{-1} \mathbb{E}_{\mu_{\text{diag}}}[\Phi\tilde{R}]$. Since \tilde{R} is linear in
 1278 Φ with coefficient vector (a, b, w) and $\mathbb{E}_{\mu_{\text{diag}}}[\Phi\Phi^\top] = \Sigma$, we have $g(\tilde{R}; \mu_{\text{diag}}) = (a, b, w)$ exactly. The canonical mitigation
 1279 $M_1(\tilde{R})(y) = \tilde{R}(y) - a\phi_1(y) = b\phi_2 + w\phi_3$ has reliance $g(M_1(\tilde{R}); \mu_{\text{diag}}) = (0, b, w)$ by Lemma A.6.
 1280
 1281

1282 *Measurement-vs-optimization gap.* We progress through three settings of $(\pi_{\text{ref}}, \mu_{\text{diag}})$ to isolate where the gap of
 1283 Section A.2 can manifest in this construction.
 1284

1285 *Isotropic reference and audit* ($\pi_{\text{ref}} = \mu_{\text{diag}} = \mathcal{N}(0, I_3)$): the mitigated proxy induces $\pi' = \pi_\beta^*(M_1(\tilde{R})) = \mathcal{N}((0, b, w)/\beta, I_3)$.
 1286 The policy-side Gram is $\mathbb{E}_{\mu_{\pi'}}[\Phi\Phi^\top] = I_3 + (0, b, w)^\top(0, b, w)/\beta^2$, and
 1287

$$1288 \quad g_1(M_1(\tilde{R}); \mu_{\pi'}) = (\mathbb{E}_{\mu_{\pi'}}[\Phi\Phi^\top]^{-1} \mathbb{E}_{\mu_{\pi'}}[\Phi M_1(\tilde{R})])_1 = 0,$$

1289 since $\mathbb{E}_{\mu_{\pi'}}[\phi_1 M_1(\tilde{R})] = 0$ in this case.

1290 *Isotropic reference, correlated audit* ($\pi_{\text{ref}} = \mathcal{N}(0, I_3)$, $\mu_{\text{diag}} = \mathcal{N}(0, \Sigma)$ with $\Sigma_{12} = \rho \neq 0$): the canonical mitigation
 1291 is $M_1(\tilde{R}) = \tilde{R} - g_1(\tilde{R}; \mu_{\text{diag}})\phi_1$ with $g_1(\tilde{R}; \mu_{\text{diag}}) = a$ unchanged (since $\Sigma_{13} = 0$ keeps ϕ_3 's contribution clean and
 1292 the (ϕ_1, ϕ_2) block of Σ^{-1} recovers a for a linear proxy). Evaluating reliance at $\mu_{\pi'} = \mathcal{N}((0, b, w)/\beta, I_3)$ gives Gram
 1293 $I_3 + (0, b, w)^\top(0, b, w)/\beta^2$, diagonal in the ϕ_1 row, so $g_1(M_1(\tilde{R}); \mu_{\pi'}) = 0$ here as well.
 1294
 1295

1296 *Coupled correlated case* ($\pi_{\text{ref}} = \mu_{\text{diag}} = \mathcal{N}(0, \Sigma)$ with $\Sigma_{12} = \rho \neq 0$): the policy-side Gram now inherits the audit
 1297 cross-correlation, $\mu_{\pi'} = \mathcal{N}(\Sigma c/\beta, \Sigma)$ for $c = (0, b, w)$, and one might expect a non-trivial g_1 . It does not arise, for a
 1298

1299 ³The closed-form expressions extend continuously to the unrestricted Gaussian limit.

structural reason: for any linear proxy $\bar{R} = c^\top y$ on $\Phi(y) = y$ and any non-degenerate measure $\mu = \mathcal{N}(m, V)$,

$$g(\bar{R}; \mu) = (\mathbb{E}_\mu[\Phi\Phi^\top])^{-1} \mathbb{E}_\mu[\Phi \bar{R}] = (V + mm^\top)^{-1}(V + mm^\top) c = c,$$

so g recovers the coefficient vector exactly and is invariant to μ . Applied to $c = (0, b, w)$, this gives $g_1(M_1(\tilde{R}); \mu_{\pi'}) = 0$ in the coupled case as well, and the g -level gap is degenerate throughout the linear-Gaussian closed form.

The substantive gap manifests instead through the policy-induced first moment. In the coupled correlated case, $\mathbb{E}_{\mu_{\pi'}}[\Phi] = \Sigma c / \beta$, so

$$\mathbb{E}_{\mu_{\pi'}}[\phi_1] = (\Sigma c)_1 / \beta = \rho b / \beta \neq 0$$

whenever $\rho \neq 0$ and $b \neq 0$, even though g_1 vanishes at every measure. The mechanism is the cross-correlation ρ : at μ_{diag} the projection orthogonalizes $M_1(\tilde{R})$ against ϕ_1 in the Gram-inner-product sense, but the policy under $M_1(\tilde{R})$ shifts the mean of (ϕ_2, ϕ_3) by $(b, w) / \beta$, and any ϕ_1 - ϕ_2 coupling at μ_{diag} propagates into ϕ_1 first-moment drift at the policy distribution. This is the form of the gap the regime taxonomy of Section A.3 is sensitive to, since Definitions A.8-A.11 are stated in terms of Δ_j and ΔJ rather than g_i at μ_{π^*} .

Remark (non-linear proxies activate the g -level gap). The linear-Gaussian construction is structurally incapable of exhibiting a non-degenerate g -level gap, because the OLS-style reliance map $g(\bar{R}; \mu)$ recovers c exactly for any linear $\bar{R} = c^\top y$, independent of μ . Activating a non-trivial g -level gap requires breaking exact representability of \tilde{R} in $\text{span}(\Phi)$ in a way that couples to moments of μ_{diag} that vary across audit distributions. Augmenting \tilde{R} with a non-linear term outside $\text{span}(\Phi)$ (e.g., an interaction $\kappa \phi_1 \phi_2$ at non-symmetric μ_{diag} , or a higher-order term whose relevant cross-moments differ across $\mu_{\text{diag}}^{(1)}$ and $\mu_{\text{diag}}^{(2)}$) makes $g(\tilde{R}; \mu_{\text{diag}})$ depend on those moments. Under such an augmentation, $M_1^{\mu_{\text{diag}}^{(1)}} \neq M_1^{\mu_{\text{diag}}^{(2)}}$ as operators and the operator-level audit sensitivity of Definition A.12 becomes instantiable with π_{ref} held fixed across audit distributions. While we establish non-vacuity with the linear-Gaussian, Section E.9 provides the closed-form non-linear instantiation. Real reward models are non-linear in any tractable feature basis, so the audit-side evidence in Sections B and F.3 supports the operator-level reading of Definition A.12 in deployment, beyond the controlled construction.

Regimes R0–R3. We work in the coupled case $\pi_{\text{ref}} = \mu_{\text{diag}} = \mathcal{N}(0, \Sigma)$ for the same reason as the gap analysis above: under decoupled isotropic π_{ref} the post-mitigation policy is $\mathcal{N}((0, b, w) / \beta, I_3)$, which gives $\Delta_2 = 0$ and $\Delta J = 0$ identically and collapses the regime taxonomy to a single regime. The coupled case is the minimal setting in which the five regimes are jointly distinguishable in this construction. The limitation of relying on $\pi_{\text{ref}} = \mu_{\text{diag}}$ noted under R4 below applies symmetrically here. Under this coupling, the policy-mean for any linear reward $\bar{R} = c^\top y$ is $\Sigma c / \beta$. The pre-mitigation policy is $\pi = \pi_\beta^*(\tilde{R}) = \mathcal{N}(\Sigma(a, b, w)^\top / \beta, \Sigma)$ and the post-mitigation policy is $\pi' = \pi_\beta^*(M_1(\tilde{R})) = \mathcal{N}(\Sigma(0, b, w)^\top / \beta, \Sigma)$. The off-target spurious-feature drift and true-reward change evaluate to

$$\Delta_2(\pi, \pi') = -\rho a / \beta, \quad \Delta J = -\Sigma_{13} \cdot aw / \beta = 0,$$

the latter because $\Sigma_{13} = 0$ in our parameterization. This makes $\Delta J = 0$ in the strict $\Sigma_{13} = 0$ case and rotation onto ϕ_2 governed entirely by ρa . To populate all five regimes we relax Σ_{13} as a free parameter $\sigma_{13} \in (-1, 1)$ (with Σ remaining positive definite), giving $\Delta J = -\sigma_{13} aw / \beta$. The five regimes correspond to:

- *R0 (successful mitigation):* $\rho = 0$, $\sigma_{13} < 0$, $a, w > 0$. Then $\Delta_2 = 0$ and $\Delta J = -\sigma_{13} aw / \beta > 0$. Removing spurious ϕ_1 -reliance improves true reward because ϕ_1 was anti-correlated with ϕ_3 at μ_{diag} , so the unmitigated proxy was pushing the policy away from high- ϕ_3 regions.

- 1353 • *R0_{cont} (contaminated success)*: $\rho \neq 0, \sigma_{13} < 0, a, w > 0$. Then $\Delta_2 = -\rho a/\beta \neq 0$ and $\Delta J = -\sigma_{13} a w/\beta > 0$. True
1354 reward improves as in R0 (driven by $\sigma_{13} < 0$), but the cross-correlation ρ couples ϕ_1 -reduction to first-moment
1355 drift on ϕ_2 , so the substitution mechanism is active despite improvement.
1356
- 1357 • *R1 (bias substitution)*: $\rho \neq 0, \sigma_{13} = 0, a \neq 0$. Then $\Delta_2 = -\rho a/\beta \neq 0$ and $\Delta J = 0$ (neutral substitution). Setting
1358 $\sigma_{13} > 0$ instead gives $\Delta J < 0$ (harmful substitution).
1359
- 1360 • *R2 (overcorrection)*: $\rho = 0, \sigma_{13} > 0, a, w > 0$. Then $\Delta_2 = 0$ and $\Delta J = -\sigma_{13} a w/\beta < 0$. The mechanism is
1361 audit-side cross-correlation between ϕ_1 and the structural ϕ_3 : removing $a\phi_1$ shifts the policy mean of ϕ_3
1362 by $-\sigma_{13} a/\beta$, lowering true reward without rotating pressure onto ϕ_2 . The rescalability test of Section E.12
1363 gives $\Delta J(c) = -c \sigma_{13} a w/\beta$ for partial mitigation $M_1^c(\tilde{R}) = \tilde{R} - c a\phi_1$, monotone in $c \in (0, 1)$, so no partial
1364 mitigation recovers the unmitigated return; by the operational test of Definition A.10 the construction sits
1365 on the non-rescalable side of R2 (driven here by reference-policy coupling $\Sigma_{13} \neq 0$ rather than the named
1366 target-misspecification mechanism, which would require $\phi_1 \in \Phi_{\text{struct}}$).
1367
- 1368 • *R3 (silent non-op)*: $\rho = 0, \sigma_{13} = 0$. Then $\Delta_2 = 0$ and $\Delta J = 0$ regardless of a . The mitigation alters \tilde{R} along an axis
1369 the policy's true-reward-relevant directions are orthogonal to.
1370
1371

1372 Each regime is realized by an open set of parameter values, not a measure-zero corner.
1373

1374 *Audit-distribution sensitivity (R4)*. The strict operator-level reading of Definition A.12 (varying μ_{diag} with π_{ref} held
1375 fixed) is degenerate under linearity: in the linear-Gaussian setting $g(\tilde{R}; \mu_{\text{diag}}) = (a, b, w)$ at every non-degenerate
1376 μ_{diag} , so $M_1^{\mu_{\text{diag}}^{(1)}}(\tilde{R}) = M_1^{\mu_{\text{diag}}^{(2)}}(\tilde{R}) = \tilde{R} - a\phi_1$ as operators, whatever the difference between $\mu_{\text{diag}}^{(1)}$ and $\mu_{\text{diag}}^{(2)}$. Linearity thus
1377 characterizes the boundary at which Definition A.12 activates: any operator-level instantiation requires the non-linear
1378 extension of the preceding remark (instantiated closed-form in Section E.9), and Section F.3 supplies the empirical
1379 counterpart across eight model families.
1380

1381 The closed form does realize R4 jointly with π_{ref} under the coupling $\pi_{\text{ref}} = \mu_{\text{diag}}$ (the audit distribution serving as
1382 the rollout reference). Hold $(\tilde{R}, R, \Phi_{\text{sp}}, \beta)$ and the operator M_1 fixed with $a, b, w > 0$, and take two pairs $(\mu_{\text{diag}}^{(\ell)}, \pi_{\text{ref}}^{(\ell)}) =$
1383 $(\mathcal{N}(0, \Sigma^{(\ell)}), \mathcal{N}(0, \Sigma^{(\ell)}))$ for $\ell = 1, 2$ with $\Sigma_{12}^{(1)} = \rho^{(1)} \neq 0, \sigma_{13}^{(1)} = 0$, and $\Sigma_{12}^{(2)} = \rho^{(2)} \neq 0, \sigma_{13}^{(2)} < 0$ (each $\Sigma^{(\ell)}$ positive
1384 definite, all other entries equal across the two). Then
1385

$$1386 \Delta_2^{(\ell)} = -\rho^{(\ell)} a/\beta, \quad \Delta J^{(\ell)} = -\sigma_{13}^{(\ell)} a w/\beta,$$

1387 so $\ell = 1$ realizes R1 (neutral substitution: $\Delta_2 \neq 0, \Delta J = 0$) and $\ell = 2$ realizes R0_{cont} ($\Delta_2 \neq 0, \Delta J > 0$): the audit distribution
1388 selects the regime, with π_{ref} entering jointly via the coupling. For the strict operator-level version (π_{ref} fixed across
1389 audit distributions), Section E.9 provides the closed-form instantiation and Section F.3 the empirical counterpart (sign
1390 reversal of the length-sycophancy coupling under human-LLM judge disagreement).
1391
1392
1393
1394

1395 *Norm increase under M_i* . The scale identity from Section E.5 reads $\|\tilde{R}'\|_{L^2(\mu_{\text{diag}})}^2 = \|\tilde{R}\|_{L^2(\mu_{\text{diag}})}^2 - g_1 [g_1 G_{11} + 2(G_{12} g_2 +$
1396 $G_{13} g_3)]$. With $G = \Sigma$ in our parameterization, $G_{11} = 1, G_{12} = \rho, G_{13} = \sigma_{13}$, and $(g_1, g_2, g_3) = (a, b, w)$, the bracket equals
1397 $a + 2\rho b + 2\sigma_{13} w$. Choosing $a = 0.5, b = w = 1, \rho = -0.6, \sigma_{13} = -0.4$ gives bracket = $0.5 - 1.2 - 0.8 = -1.5$ with bracket
1398 sign opposite to $g_1 = 0.5$. The norm change is $-g_1 \cdot (-1.5) = +0.75$, so $\|\tilde{R}'\|^2 > \|\tilde{R}\|^2$ strictly. Mitigation increases the
1399 proxy's L^2 norm because the cross-correlation contribution dominates the diagonal term, exactly the failure mode
1400 flagged in Section E.5.
1401
1402
1403
1404

E.8 Epsilon-Banded Regime Classification

Finite-sample classification of mitigations into the regimes of Definitions A.8-A.11 requires bands, since exact equality on Δ_j and ΔJ is a measure-zero event under any non-degenerate sampling design.

Banded definitions. Fix tolerances $\varepsilon_j \geq 0$ for each $\phi_j \in \Phi_{\text{sp}} \setminus \{\phi_i\}$ and $\varepsilon_J \geq 0$. With π, π' as in Definitions A.8-A.11, define

$$R0_\varepsilon : \Delta J > \varepsilon_J \text{ and } |\Delta_j(\pi, \pi')| \leq \varepsilon_j \text{ for all } \phi_j \in \Phi_{\text{sp}} \setminus \{\phi_i\},$$

$$R0_{\varepsilon, \text{cont}} : \Delta J > \varepsilon_J \text{ and } |\Delta_j(\pi, \pi')| > \varepsilon_j \text{ for some } \phi_j \in \Phi_{\text{sp}} \setminus \{\phi_i\},$$

$$R1_\varepsilon : |\Delta_j(\pi, \pi')| > \varepsilon_j \text{ for some } \phi_j \in \Phi_{\text{sp}} \setminus \{\phi_i\}, \text{ and } \Delta J \leq \varepsilon_J,$$

$$R2_\varepsilon : |\Delta_j(\pi, \pi')| \leq \varepsilon_j \text{ for all } \phi_j \in \Phi_{\text{sp}} \setminus \{\phi_i\}, \text{ and } \Delta J < -\varepsilon_J,$$

$$R3_\varepsilon : |\Delta_j(\pi, \pi')| \leq \varepsilon_j \text{ for all } \phi_j \in \Phi_{\text{sp}} \setminus \{\phi_i\}, \text{ and } |\Delta J| \leq \varepsilon_J.$$

$R1$'s sub-cases port unchanged: neutral and harmful substitution correspond to $|\Delta J| \leq \varepsilon_J$ and $\Delta J < -\varepsilon_J$ respectively. The spurious-rotation-with-improvement corner of Section E.3 is regime $R0_{\varepsilon, \text{cont}}$ per the definition above. Definition A.12 ports verbatim, since it is transversal to the regime classification and unaffected by the choice of bands.

Nested limit. Taking $\varepsilon_j, \varepsilon_J \rightarrow 0$ recovers Definitions A.8-A.11 pointwise: $R0_\varepsilon \rightarrow \{\text{no rotation, } \Delta J > 0\}$, $R0_{\varepsilon, \text{cont}} \rightarrow \{\text{rotation, } \Delta J > 0\}$, $R1_\varepsilon \rightarrow \{\text{rotation, } \Delta J \leq 0\}$, $R2_\varepsilon \rightarrow \{\text{no rotation, } \Delta J < 0\}$, $R3_\varepsilon \rightarrow \{\text{no rotation, } \Delta J = 0\}$, with the disjunctive condition $|\Delta_j| > \varepsilon_j$ for some j converging to $\Delta_j \neq 0$ for some j as $\varepsilon_j \rightarrow 0$.

Default choice: noise-floor bands. We default to noise-floor bands, in which $\varepsilon_j = z_{\alpha/2} \cdot \widehat{\text{SE}}(\widehat{\Delta}_j)$ under the empirical sampling design (e.g., the standard error of a fixed-effect coefficient in a mixed linear model, as in Section F.3, where the sycophancy-on-length coefficient instantiates $\widehat{\Delta}_{\text{length}}$) and $\varepsilon_J = z_{\alpha/2} \cdot \widehat{\text{SE}}(\widehat{\Delta}_J)$ for a significance level α chosen and reported by the practitioner. Setting $\varepsilon_j = \widehat{\text{SE}}$ without the $z_{\alpha/2}$ factor corresponds to approximately 68% confidence and should not be conflated with a standard significance test; regime assignments should always state the chosen α . An effect-size-relative alternative (ε_j as a fraction of unmitigated $\mathbb{E}_{\mu_\pi}[\phi_j]$, ε_J as a fraction of unmitigated J) is appropriate for cross-proxy comparisons; the two can disagree on borderline cases and the choice should be stated. When $|\Phi_{\text{sp}} \setminus \{\phi_i\}| > 1$, the conservative reading (declare rotation if any axis exceeds its ε_j band, no multiple-testing correction) inflates the false-positive rate of the rotation criterion. Because rotation is the discriminating predicate for both $R0_{\varepsilon, \text{cont}}$ (against $R0_\varepsilon$) and $R1_\varepsilon$ (against $R3_\varepsilon \cup R2_\varepsilon$), this inflation pulls cases out of $R0_\varepsilon$ into $R0_{\varepsilon, \text{cont}}$ when $\Delta J > \varepsilon_J$, and out of $R3_\varepsilon$ and $R2_\varepsilon$ into $R1_\varepsilon$ when $\Delta J \leq \varepsilon_J$. This choice is deliberate to prioritize sensitivity to substitution over specificity in either ΔJ regime, consistent with the conservative-partition stance of Section 3.1 and the neutral sub-case absorbing near-zero improvements ($\Delta J \in (0, \varepsilon_J]$) in the same direction.

Mapping wild classifications onto the bands. The $R0$ – $R4$ calls in Sections A.3 and B are made under the banded reading, with $\varepsilon_j, \varepsilon_J$ inferred from each cited work's reported uncertainty. The [3] accuracy degradation under uniform length penalties satisfies the $R2_\varepsilon$ pattern under noise-floor bands, with the rescalability test of Section E.12 itself constituting a banded operationalization (existence of $c \in (0, 1)$ improving \widehat{J} implicitly assumes detection against ε_j). The Fein et al. [21] sign flip together with reduced correctness is read as $R1_\varepsilon$ under the rotation interpretation (off-target spurious axis exceeds its band, $\widehat{\Delta}_J \leq \varepsilon_J$). The alternative reading as $R2_\varepsilon$ on the targeted axis itself is also consistent with the

reported statistics, and we flag this ambiguity rather than resolve it. The human-vs-LLM-judge gap is operationalized via Section F.3 below.

Mapping Section F.3 statistics onto the bands. The mixed-model coefficients reported in Section F.3 (+24.3 under human labels, CI [9.6, 39.0], +128.4 under LLM labels, CI [118.5, 138.4], +154.3 under judge agreement, CI [123.6, 185.0], -43.1 under judge disagreement, CI [-60.5, -25.7]) all satisfy the noise-floor ϵ_j -band test on the length axis, since each CI excludes zero. The pooled KS statistics on within-model residuals (0.025, 0.130, 0.130, 0.046, all with $p < 0.05$) provide a distributional band test alongside the first-moment Δ_j of Definition A.7. The sign reversal of the length-sycophancy coupling between the agreement-side regimes and the disagreement regime is consistent with the construction-level precondition for R4_e, in that two annotator-conditioned audit distributions yield g -profiles with opposite signs on the relevant cross-feature coupling. The operator-level g -difference is not measured empirically here, but its constructive counterpart appears in Section E.9.

E.9 Phase-Diagram Validation of Regime Transitions Under Quadratic Non-Linearity

Extending the linear Gaussian non-vacuity instantiation of Section E.7, we use a quadratic non-linear reward structure to empirically show

- clear instantiation of all regimes and sub-regimes (R0, R0_{cont}, R1_{neut}, R1_{harm}, R2, R3) with phase-diagrams illustrating how the regime boundaries are outcomes of a single mitigation operator in a controlled setting
- that the regimes are not a property of the reward model alone but of $(R, \tilde{R}, M_i, \mu_{\text{diag}})$ as stated in Section A.
- clear instantiation for the regime transition as described by R4 from only varying the audit distribution μ_{audit} mean
- a match of the theoretical predictions connecting the analytic phase-diagram claims to finite- N experiments, strengthening the non-vacuity results.

This generalizes the linear-Gaussian setting, as we show that the regime taxonomy is not a linear-Gaussian artifact. For real deployment instantiations, see different experiments in Section F.

Setup. We extend the setup in Section E.7: We work with a trivial prompt space (\mathcal{X} a singleton, suppressed in notation) and continuous response space $\mathcal{Y} = \mathbb{R}^3$. The feature map is $\Phi(y) = (\phi_1, \phi_2, \phi_3)$ with $\phi_k(y) = y_k$. We fix true and the preference-generating annotator reward

$$R(y) = w \phi_3(y), \quad R_{\text{anno}}(y) = a \phi_1(y) + b \phi_2(y) + w \phi_3(y) + \gamma \phi_1(y)^2,$$

with $w, \gamma > 0$ and $(a, b) \in \mathbb{R}^2$. By Definition A.3, $\Phi_{\text{sp}} = \{\phi_1, \phi_2\}$ and $\Phi_{\text{struct}} = \{\phi_3\}$ by construction. For the reward model (learned proxy reward), we use a linear polynomial such that

$$\tilde{R}(y) = \theta_1 \phi_1(y) + \theta_2 \phi_2(y) + \theta_3 \phi_3(y) + \theta_4 \phi_1(y)^2,$$

which is correctly specified using enough samples and the maximum likelihood estimator for Bradley-Terry (BT-MLE) as $\hat{\theta} = (a, b, w, \gamma)$.⁴ Reference policy and diagnostic measure are Gaussians on \mathbb{R}^3 ,

$$\pi_{\text{ref}} = \mathcal{N}(0, \Sigma_{\text{ref}}), \quad \mu_{\text{diag}} = \mathcal{N}(m, \Sigma_{\text{ref}}),$$

⁴We use the analytical idealization for the correct specification of the reward model for closed-form regime predictions. Under realistic RM misspecification, the same regime phenomena persist with additional noise, as documented empirically across five reward models and a non-linear operator in Section F.2.

where Σ is positive definite with

$$\Sigma_{\text{ref}} = \begin{pmatrix} 1 & \rho_{12} & \rho_{13} \\ \rho_{12} & 1 & 0 \\ \rho_{13} & 0 & 1 \end{pmatrix}, \quad \rho_{12}^2 + \rho_{13}^2 < 1$$

Although this setup ties the covariance of the reference and audit distribution, we show that the mean parameter $m \in \mathbb{R}^3$ of the audit distribution μ_{diag} is sufficient to instantiate R4. We fix the KL parameter $\beta > 0$ and target the spurious feature $i = 1$ for mitigation throughout.

Mitigation observables. Following Definition A.4, the linear reliance g_i for targeted $i = 1$ is

$$g_1(\theta, m) = \theta_1 + \theta_2 \rho_{12} + \theta_3 \rho_{13} + 2\theta_4 m_1.$$

The resulting single-axis mitigation with strength $c \geq 0$ is then

$$M_1(\tilde{R}(y)) = \tilde{R}(y) - cg_1 \phi_1(y).$$

We rewrite the reward model as

$$\tilde{R}(y) = \alpha \cdot y + \frac{1}{2} y^\top H y, \quad \text{with } \alpha = (\theta_1, \theta_2, \theta_3), H = \text{diag}(2\theta_4, 0, 0),$$

reparameterizing the KL-regularized optimum policy as

$$\pi^* = \mathcal{N}(\mu^*, \Sigma^*) \quad \text{with } \Sigma^{*-1} = \Sigma_{\text{ref}}^{-1} - \frac{H}{\beta}, \mu^* = \frac{\Sigma^* \alpha}{\beta}.$$

Since H is unchanged by the mitigation, Σ^* is also unchanged and only α shifts by $\Delta\alpha = (-cg_1, 0, 0)$, leading to the regime-defining changes in the policy-induced feature expectation and true reward (see Section A.3) as

$$\Delta_j = (\Sigma^* \Delta\alpha) / \beta = -\frac{cg_1}{\beta} \Sigma_{j,1}^*, \quad \Delta J = w \Delta_3 = -w \frac{cg_1}{\beta} \Sigma_{3,1}^*.$$

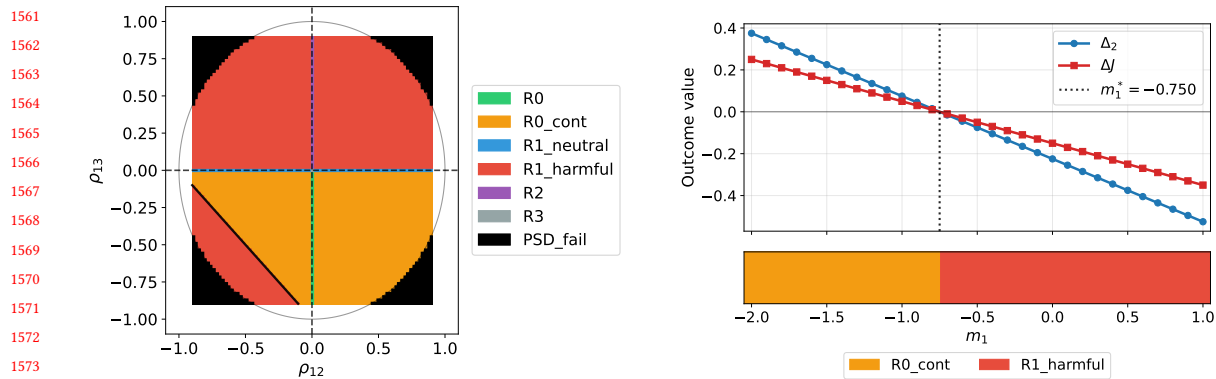
Expanding with Sherman-Morrison gives $\Sigma_{j,1}^* = \rho_{1j(1+\kappa)}$ with $\kappa = (2\gamma/\beta)/(1 - 2\gamma/\beta) \geq 0$, implying that rescaling γ does not flip the sign of Δ_j for regime classification.

Simulation parameters. With this setup, we have the following parameters to tune and induce regime transitions:

- ρ_{12} sets rotation for bias substitution and also the sign of Δ_2
- ρ_{13} sets the sign of ΔJ (if the sign of g_1 does not change)
- γ amplifies Δ_j and ΔJ , but does not move boundaries at $m = 0$
- m_1 shifts linear reliance g_1 and leads to g_1 sign flip at $m_1^* = -(a + b\rho_{12} + w\rho_{13})/(2\gamma)$ inducing the transition for R4

The KL parameter β and the mitigation scale parameter c only set margins and scales.

Experiment Results. Using $N = 10^5$ preference samples, 50 seeds per validation cell, a linearly spaced grid of 81 by 81 values for $\rho_{12}, \rho_{13} \in (-0.9, 0.9)^2$ each, and a (near-numerical precision) ϵ -bound of 10^{-8} (see Section E.8), we instantiate all five regimes R0–R3 and the transition R4, while also validating the above theoretical predictions with finite- N sample experiments. Figure 2 shows the instantiation, Figure 3 shows the regime boundary robustness to non-linearity strength, Figure 4 verifies predictions and measurements, and Table 2 presents the numerical values and uncertainties of Figure 4.



(a) Phase diagram from varying (ρ_{12}, ρ_{13}) at fixed $\gamma = 0.4, m = 0, c = 1, \beta = 4$. We color each grid cell with the corresponding regime from Section A.3. (b) R4 transition of Δ_2 and Δ_J induced by varying the audit distribution μ_{diag} mean $m_1 \in [-2, 1]$ at fixed $\gamma = 1.0, \rho_{12} = 0.3, \rho_{13} = 0.2, c = 1, \beta = 4$.

Fig. 2. We instantiate all five R0–R3 as tune-able outcomes of a single mitigation operator and R4 by varying only the mean of the audit distribution μ_{diag} in a controlled setting. Figure 2a shows that a non-linearity term creates new boundary structure and that the regimes are not solely induced as a reward model property.

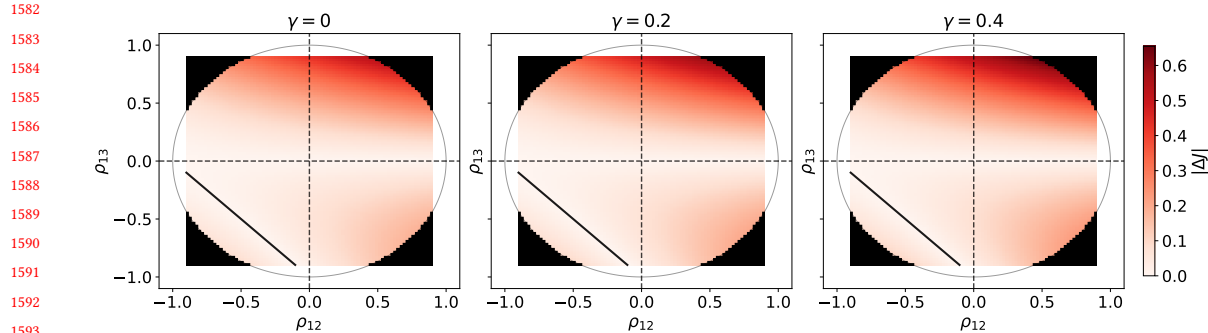


Fig. 3. Shifting γ to produce different $|\Delta_J|$ -heatmaps of Figure 2a to illustrate regime robustness to non-linearity strength, as the regime boundaries do not change themselves.

E.10 Formal Proof of Audit-Distribution Insufficiency

We give a formal proof of the structural-blindness claim of Section A.4 that standard ordinal reward-model benchmarks cannot reliably distinguish R0 from R0_{cont}, R1, or R2.

Benchmark class. Let \mathcal{B} denote the class of benchmark functionals depending only on the joint distribution of $(\tilde{R}, M_i(\tilde{R}), R, \Phi)$ under μ_{diag} . We take μ_{diag} as the empirical measure of the evaluation set of \mathcal{B} where applicable. This class subsumes ranking accuracy, pairwise win-rate, preference-prediction calibration, reward-target correlation, the linear-reliance statistic g of Definition A.4, and any composite of these, covering the evaluation paradigm of Section A.5. Note that including R in the input of \mathcal{B} strengthens the result, since standard benchmarks lack oracle access to R .

THEOREM A.13 (AUDIT-DISTRIBUTION INSUFFICIENCY). Fix $\beta > 0$. There exist four choices of π_{ref} with $\mu_{\text{diag}}, \tilde{R}, M_i, R, \Phi, \Phi_{\text{sp}}$, and β held fixed s.t.

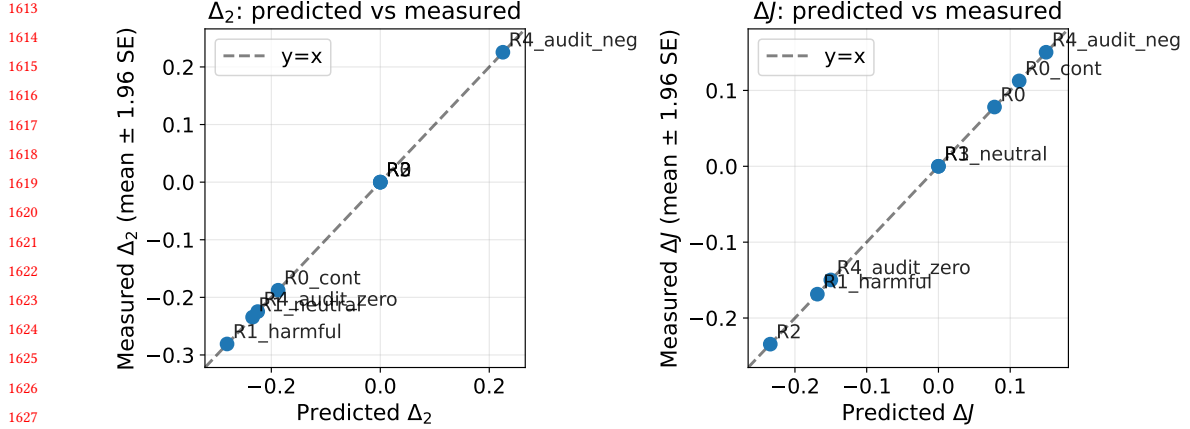


Fig. 4. Validating that the simulated regime transitions align with the theoretical predictions from our setup for Δ_2 and Δ_J with standard errors (SE) from the BT-MLE with $N = 10^5$ samples, aggregated over 50 seeds. This result connects the analytic phase-diagram claims to finite- N experiments.

Table 2. Validation of closed-form regime predictions against finite- N BT-MLE on simulated preferences. All cells use $a = b = w = c = 1$, $\beta=4$, $N=10^5$ samples per seed, for $S=50$ seeds. The first six columns are the regime instances for $m_1 = 0$, except for R4 where we vary m_1 at fixed $\rho_{12}, \rho_{13}, \gamma$ (see Figure 2b). All eight cells achieve 100% agreement on the respective regime call.

| Cell | R0 | R0 _{cont} | R1 _{neut} | R1 _{harm} | R2 | R3 | R4 | |
|-----------------------|--------|--------------------|--------------------|--------------------|---------|--------|--------------------|--------------------|
| Scenario | | | | | | | | |
| ρ_{12} | 0.0 | 0.5 | 0.5 | 0.5 | 0.0 | 0.0 | 0.3 | 0.3 |
| ρ_{13} | -0.5 | -0.3 | 0.0 | 0.3 | 0.5 | 0.0 | 0.2 | 0.2 |
| γ | 0.4 | 0.4 | 0.4 | 0.4 | 0.4 | 0.4 | 1.0 | 1.0 |
| m_1 | 0.0 | 0.0 | 0.0 | 0.0 | 0.0 | 0.0 | -1.5 | 0.0 |
| Induced g_1 | 0.500 | 1.200 | 1.500 | 1.800 | 1.500 | 1.000 | -1.500 | 1.500 |
| Δ_2 | | | | | | | | |
| predicted | 0.0000 | -0.1875 | -0.2344 | -0.2812 | 0.0000 | 0.0000 | 0.2250 | -0.2250 |
| measured | 0.0000 | -0.1876 | -0.2344 | -0.2809 | 0.0000 | 0.0000 | 0.2255 | -0.2247 |
| SE | 0.0000 | 0.0002 | 0.0003 | 0.0004 | 0.0000 | 0.0000 | 0.0006 | 0.0005 |
| Δ_J | | | | | | | | |
| predicted | 0.0781 | 0.1125 | 0.0000 | -0.1688 | -0.2344 | 0.0000 | 0.1500 | -0.1500 |
| measured | 0.0781 | 0.1126 | 0.0000 | -0.1686 | -0.2345 | 0.0000 | 0.1503 | -0.1498 |
| measured SE | 0.0002 | 0.0001 | 0.0000 | 0.0002 | 0.0003 | 0.0000 | 0.0004 | 0.0003 |
| Regime classification | R0 | R0 _{cont} | R1 _{neut} | R1 _{harm} | R2 | R3 | R0 _{cont} | R1 _{harm} |

- (i) the joint distribution of $(\tilde{R}, M_i(\tilde{R}), R, \Phi)$ under μ_{diag} is identical across the four instances, so B takes the same value on all four for every $B \in \mathcal{B}$;
- (ii) the optimized policies $\pi_{\beta}^*(M_i(\tilde{R}))$ fall into regimes R0, R0_{cont}, R1, and R2 respectively, in the sense of Definitions A.8–A.10.

PROOF. Construction. Take \mathcal{X} a singleton (suppressed), $\mathcal{Y} = \mathbb{R}^3$, feature map $\phi_k(y) = y_k$, and rewards $R(y) = \phi_3(y)$, $\tilde{R}(y) = \phi_1(y) + \phi_2(y) + \phi_3(y)$. Set $\Phi_{\text{sp}} = \{\phi_1, \phi_2\}$, target $\phi_i = \phi_1$, and audit distribution $\mu_{\text{diag}} = \mathcal{N}(0, I_3)$. Fix a small

$\delta \in (0, 1/2)$. The four instances differ only in $\pi_{\text{ref}}^{(k)} = \mathcal{N}(0, \Sigma^{(k)})$, where each $\Sigma^{(k)}$ has unit diagonal, $\Sigma_{23}^{(k)} = 0$, and

$$\begin{aligned} (\Sigma_{12}^{(0)}, \Sigma_{13}^{(0)}) &= (0, -\frac{1}{2}), & (\Sigma_{12}^{(1)}, \Sigma_{13}^{(1)}) &= (\frac{1}{2}, -\frac{1}{2}), \\ (\Sigma_{12}^{(2)}, \Sigma_{13}^{(2)}) &= (\frac{1}{2}, \delta), & (\Sigma_{12}^{(3)}, \Sigma_{13}^{(3)}) &= (0, \frac{1}{2}). \end{aligned}$$

Each $\Sigma^{(k)}$ has $\det = 1 - \Sigma_{12}^2 - \Sigma_{13}^2 > 0$ by Sylvester's criterion: $\det^{(0)} = \det^{(3)} = 3/4$, $\det^{(1)} = 1/2$, and $\det^{(2)} = 3/4 - \delta^2 > 0$ for $\delta < 1/2$. We work on a sufficiently large bounded set to satisfy Assumption E.1. Expressions extend to the unrestricted Gaussian limit as in Section E.7. Linearity and Gaussianity are used for closed-form existence and are not required for the impossibility claim. Assumption A.2 holds since $\mathbb{E}_{\mu_{\text{diag}}}[\Phi\Phi^\top] = I_3$. Assumption E.2 holds since $\mathcal{Y} = \mathbb{R}^3$ admits independent coordinate perturbations on a set of full μ_{diag} -measure.

(i) Audit-side identity. The Gram at μ_{diag} is I_3 , giving $g(\tilde{R}; \mu_{\text{diag}}) = (1, 1, 1)$ and $M_i(\tilde{R}) = \phi_2 + \phi_3$ across all four instances. Since \tilde{R} , $M_i(\tilde{R})$, R , Φ , and μ_{diag} are identical across k , the joint distribution of $(\tilde{R}, M_i(\tilde{R}), R, \Phi)$ under μ_{diag} is identical across k , so B takes the same value on all four instances for every $B \in \mathcal{B}$.

(ii) Policy-side classification. For any linear $\tilde{R}(y) = c^\top y$ and $\pi_{\text{ref}} = \mathcal{N}(0, \Sigma)$, completing the square in Equation (2) gives $\pi_\beta^*(\tilde{R}) = \mathcal{N}(\Sigma c/\beta, \Sigma)$. Applying this with $c = (1, 1, 1)^\top$ and $c' = (0, 1, 1)^\top$ at each $\Sigma^{(k)}$:

$$\Delta_2(\pi^{(k)}, \pi'^{(k)}) = -\Sigma_{12}^{(k)}/\beta, \quad \Delta J^{(k)} = -\Sigma_{13}^{(k)}/\beta,$$

where Δ_j is checked against $\Phi_{\text{sp}} \setminus \{\phi_i\} = \{\phi_2\}$. Substituting:

$$\begin{aligned} (\Delta_2, \Delta J)^{(0)} &= (0, +\frac{1}{2\beta}), & (\Delta_2, \Delta J)^{(1)} &= (-\frac{1}{2\beta}, +\frac{1}{2\beta}), \\ (\Delta_2, \Delta J)^{(2)} &= (-\frac{1}{2\beta}, -\frac{\delta}{\beta}), & (\Delta_2, \Delta J)^{(3)} &= (0, -\frac{1}{2\beta}). \end{aligned}$$

By Definitions A.8-A.10, these are R0, R0_{cont}, R1 (harmful substitution), and R2 respectively. \square

Remark (scope). The theorem fixes μ_{diag} and M_i across instances. R3 is additionally realizable in the same construction by setting $\Sigma_{12} = \Sigma_{13} = 0$, with audit observables still identical and distinguishability from R0 requiring $D_{\text{KL}}(\pi' \parallel \pi)$, which is not audit-local. R4 requires non-linear proxies (see Sections E.7 and E.9). The varying π_{ref} corresponds to the same reward model deployed against different reference policies, the standard regime in which benchmarks claim to evaluate reward models independently of downstream RLHF setup.

Functional blindspot. Theorem A.13 establishes the distributional blindspot. The functional blindspot is independent and follows from cardinal-scale sensitivity:

Corollary E.4 (Functional blindspot via reward rescaling). *Let $\mathcal{B}_{\text{ord}} \subseteq \mathcal{B}$ denote the benchmark functionals satisfying $B(f \circ \tilde{R}, f \circ M_i(\tilde{R}), R, \Phi, \mu_{\text{diag}}) = B(\tilde{R}, M_i(\tilde{R}), R, \Phi, \mu_{\text{diag}})$ for every strictly monotone increasing $f: \mathbb{R} \rightarrow \mathbb{R}$. Then for the construction of Theorem A.13, every $B \in \mathcal{B}_{\text{ord}}$, and every $c > 0$,*

$$B(c\tilde{R}, M_i(c\tilde{R}), R, \Phi, \mu_{\text{diag}}) = B(\tilde{R}, M_i(\tilde{R}), R, \Phi, \mu_{\text{diag}}),$$

while $J(\pi_\beta^*(c\tilde{R}), R)$ is strictly monotone in c .

PROOF. By linearity, $M_i(c\tilde{R}) = c M_i(\tilde{R})$, and ordinal invariance under $f(y) = cy$ gives the benchmark identity. The KL identity $\pi_\beta^*(c\tilde{R}) = \pi_{\beta/c}^*(\tilde{R})$ (Section E.5) gives $J(\pi_\beta^*(c\tilde{R}), R) = c \cdot (\Sigma_{c\tilde{R}})_3/\beta$ in the construction of Theorem A.13, where $c_{\tilde{R}} = (1, 1, 1)^\top$ and $(\Sigma^{(k)} c_{\tilde{R}})_3 = \Sigma_{13}^{(k)} + 1 \neq 0$ for each k , so J is strictly monotone in c . \square

E.11 Audit-Distribution Sufficiency

Theorem A.13 establishes that benchmark functionals depending only on $(\tilde{R}, M_i(\tilde{R}), R, \Phi, \mu_{\text{diag}})$ cannot reliably separate R0 from R0_{cont}, R1, or R2. We now show that adding the policy-induced distributions $\mu_{\pi_\beta^*}(\tilde{R})$ and $\mu_{\pi_\beta^*(M_i(\tilde{R}))}$ to the input suffices, certifying the prescription that benchmarks must evaluate at policy-induced distributions (Do 1 of Section E.12).

Extended benchmark class. Let \mathcal{B}^+ denote the class of benchmark functionals depending on the joint distributions of $(\tilde{R}, M_i(\tilde{R}), R, \Phi)$ under each of μ_{diag} , $\mu_{\pi_\beta^*}(\tilde{R})$, and $\mu_{\pi_\beta^*(M_i(\tilde{R}))}$, together with the partition $\Phi = \Phi_{\text{sp}} \sqcup \Phi_{\text{struct}}$. We have $\mathcal{B} \subset \mathcal{B}^+$.

THEOREM A.14 (AUDIT-DISTRIBUTION SUFFICIENCY). Fix $\beta > 0$ and let $\pi = \pi_\beta^*(\tilde{R})$, $\pi' = \pi_\beta^*(M_i(\tilde{R}))$. With Δ_j and ΔJ as in Definition A.7 and Section A.3, define $B^* \in \mathcal{B}^+$ by

$$B^* = \begin{cases} \text{R0} & \text{if } \Delta J > 0 \text{ and } \Delta_j = 0 \text{ for all } \phi_j \in \Phi_{\text{sp}} \setminus \{\phi_i\}, \\ \text{R0}_{\text{cont}} & \text{if } \Delta J > 0 \text{ and } \Delta_j \neq 0 \text{ for some } \phi_j \in \Phi_{\text{sp}} \setminus \{\phi_i\}, \\ \text{R1} & \text{if } \Delta J \leq 0 \text{ and } \Delta_j \neq 0 \text{ for some } \phi_j \in \Phi_{\text{sp}} \setminus \{\phi_i\}, \\ \text{R2} & \text{if } \Delta J < 0 \text{ and } \Delta_j = 0 \text{ for all } \phi_j \in \Phi_{\text{sp}} \setminus \{\phi_i\}. \end{cases} \quad (15)$$

Under Assumptions E.1, A.2, and E.2, with $M_i(\tilde{R})$ likewise satisfying Assumption E.1, for every tuple $(\pi_{\text{ref}}, \tilde{R}, M_i, R, \Phi, \Phi_{\text{sp}})$ with $(\pi, \pi') \in \text{R0} \cup \text{R0}_{\text{cont}} \cup \text{R1} \cup \text{R2}$ in the sense of Definitions A.8–A.10, B^* returns the correct regime label.

PROOF. Δ_j and ΔJ are first moments of Φ and R under μ_π and $\mu_{\pi'}$, both inputs to \mathcal{B}^+ , so $B^* \in \mathcal{B}^+$. Finiteness of these expectations follows from Assumption E.1 via the boundedness chain in Section E.2, applied to both \tilde{R} and $M_i(\tilde{R})$ (the latter by closure under finite linear combinations).

Correctness on each regime follows directly from Definitions 3.8–3.10: R0 requires $\Delta J > 0$ with no off-target rotation (first branch), R0_{cont} requires $\Delta J > 0$ with off-target rotation (second branch), R1 requires $\Delta J \leq 0$ with off-target rotation (third branch), R2 requires $\Delta J < 0$ with no off-target rotation (fourth branch). The four branches are pairwise disjoint: $\{\text{R0}, \text{R0}_{\text{cont}}\}$ have $\Delta J > 0$ and $\{\text{R1}, \text{R2}\}$ have $\Delta J \leq 0$, separating any cross-pair by the sign of ΔJ ; within $\{\text{R0}, \text{R0}_{\text{cont}}\}$ and within $\{\text{R1}, \text{R2}\}$ the branches differ in rotation (R0_{cont}, R1 require rotation; R0, R2 do not). Inputs in R3 (no rotation, $\Delta J = 0$) are excluded by hypothesis. \square

Corollary E.5 (Sufficiency for R3). The classifier B^* extends to $\{\text{R0}, \text{R1}, \text{R2}, \text{R3}\}$ by adding a fourth branch: return R3 if $\Delta J = 0$ and $\Delta_j = 0$ for all $\phi_j \in \Phi_{\text{sp}} \setminus \{\phi_i\}$. This branch matches Definition A.11 exactly and is disjoint from the R0–R2 branches by the $\Delta J = 0$ condition. The observable $D_{\text{KL}}(\pi' \parallel \pi)$ is not required for the regime call, though reporting it alongside the label distinguishes policy-relevant R3 ($D_{\text{KL}} > 0$, mitigation moved the policy along directions orthogonal to $(\Delta_j, \Delta J)$) from policy-irrelevant R3 ($D_{\text{KL}} = 0$, mitigation did not move the policy at all).

Corollary E.6 (Finite-sample sufficiency). Under the noise-floor ε -bands of Section E.8, replacing exact equalities in Equation (15) with the banded conditions yields $B_\varepsilon^* \in \mathcal{B}^+$ returning the correct ε -banded label.

Corollary E.7 (Multi-audit sufficiency for R4). Let \mathcal{B}^{++} denote the class of benchmark functionals depending on the inputs of \mathcal{B}^+ at each of $K \geq 2$ audit distributions $\mu_{\text{diag}}^{(1)}, \dots, \mu_{\text{diag}}^{(K)}$ with the canonical $M_i^{(k)}$ constructed from each. Define B_{R4}^* as the indicator that B^* (Theorem A.14) returns different regime labels across the K instances. Then B_{R4}^* detects R4 in the sense of Definition A.12.

PROOF. By Theorem A.14, each per-distribution call of B^* returns the correct R0–R2 label. Definition A.12 defines R4 as cross-distribution disagreement on the regime label, which is exactly what B_{R4}^* tests. \square

Discussion. Theorem A.13 (audit-distribution insufficiency) and Theorem A.14 (audit-distribution sufficiency) together show that audit-only inputs cannot separate R0, R0_{cont}, R1, and R2, while augmenting \mathcal{B} 's input with the policy-induced distributions $\mu_{\pi_\beta^*(\tilde{R})}$ and $\mu_{\pi_\beta^*(M_i(\tilde{R}))}$ suffices. Corollary E.7 extends this to R4 detection under multi-audit evaluation. This grounds the policy-distribution evaluation prescription structurally rather than heuristically. We do not claim minimality of \mathcal{B}^+ , as the classifier B^* depends only on the scalars $(\Delta_j, \Delta J)$ and the partition Φ_{sp} , so any extension of \mathcal{B} that determines these quantities also suffices.

E.12 Takeaway Recommendations for RM Mitigation and Benchmarks Developers

Theorem A.13 and Corollary E.4 together circumscribe what reward-model benchmarks in the class \mathcal{B} can and cannot deliver, and what mitigation methods evaluated within \mathcal{B} cannot be certified to achieve. We translate the formal results into concise recommendations.

Regime detection procedures. Theorem A.14's classifier B^* separates R0, R0_{cont}, R1, and R2 given $(\Delta_j, \Delta J)$ at μ_{π^*} , but does not address R2 origin disambiguation, R3 detection, or R4 detection. The following three procedures close these gaps and should accompany any regime claim that touches them.

- *R2 origin disambiguation via the rescalability sweep.* Definition A.10's two origins (scale overshoot, target misspecification) are operationally distinguished by sweeping the partial mitigation $M_i^c(\tilde{R}) = \tilde{R} - c g_i(\tilde{R}; \mu_{diag}) \phi_i$ over $c \in (0, 1)$ at fixed β . If $J(\pi_\beta^*(M_i^c(\tilde{R})), R) > J(\pi_\beta^*(\tilde{R}), R)$ for some c , the regime is driven by scale overshoot at μ_{π^*} and is recoverable by partial projection. If no c improves ΔJ , the projection has removed a structurally informative component of ϕ_i (target misspecification under partial informativeness, see Section A.1) and rescaling cannot recover it. The test requires a one-dimensional sweep over c at fixed β and the same cardinal R -access already needed to distinguish R2 from R0. Concretely: publish the $\Delta J(c)$ curve, since a recoverable overshoot is a different scientific finding than a non-recoverable target misspecification.
- *R3 detection via D_{KL} .* Corollary E.5 extends B^* to R3 via the branch $\Delta J = 0 \wedge \Delta_j = 0$. Reporting $D_{KL}(\pi' \parallel \pi)$ alongside the label separates policy-relevant R3 ($D_{KL} > 0$: mitigation moved the policy along directions orthogonal to $(\Delta_j, \Delta J)$) from policy-irrelevant R3 ($D_{KL} = 0$: mitigation did not move the policy at all). Concretely: report D_{KL} whenever the regime call is R3, since only the policy-irrelevant case is consistent with a vacuous mitigation operator and should be flagged differently in downstream interpretation.
- *R4 detection via the multi- μ_{diag} protocol.* R4 (Definition A.12) is transversal to R0–R3 and observable only by constructing the canonical $M_i^{\mu_{diag}}$ at multiple audit distributions and comparing the resulting regime calls. Concretely: select at least two audit distributions of substantively different annotator provenance (e.g., μ_{diag}^{human} , μ_{diag}^{LLM}), construct $M_i^{(\ell)}$ and $\pi^{(\ell)} = \pi_\beta^*(M_i^{(\ell)}(\tilde{R}))$ at each, apply B^* to each, and treat cross-distribution disagreement on the regime call as the test.

Reward Model Bias Mitigation Method Recommendations. The prescriptions of our framework are necessary conditions for interpretable mitigation claims of any new mitigation method. Without them, R0 cannot be distinguished from R0_{cont}, R1, or R2, regardless of how strong the audit-side evidence is. We organize the recommendations by single- and

multi-axis operators and present each as a self-contained checklist to support standardized reporting across method papers.

Single-axis Don'ts.

- **Do not fit and validate M_i on the same μ_{diag} without testing alternative audit distributions.** R4 (Definition A.12) makes regime membership a function of μ_{diag} , and the sign reversal in Section F.3 shows the dependence is empirically real. *Concretely*: a length operator fit on LLM-judge preference data can inherit a sycophancy-coupling coefficient with the opposite sign of one fit on human-judge data, so the two operators should not be interchangeable even when both zero on-target reliance at their respective audit distributions.
- **Do not apply M_i without reporting the induced $L^2(\mu_{\text{diag}})$ scale change.** Section E.5 shows $\|M_i(\tilde{R})\|_{L^2(\mu_{\text{diag}})}$ generically differs from $\|\tilde{R}\|_{L^2(\mu_{\text{diag}})}$ via both a diagonal and a cross-correlation contribution, inducing an effective- β shift at fixed nominal β . *Concretely*: a method paper reporting only $g_i \rightarrow 0$ at μ_{diag} leaves readers unable to separate axis reallocation from scale-driven regime change (see single-axis dos).
- **Do not report only pooled diagnostics when the deployment target is within-prompt selection.** Section F.2 shows pooled $|\rho_{\text{len}}|$ can be driven to zero while within-prompt $|\rho_{\text{len}}^{\text{within}}|$ flips sign on three of four SOTA reward models. *Concretely*: BoN top-1 and PPO both operate within prompts, so a pooled diagnostic targeting ϕ_i does not measure the reliance the optimizer responds to.
- **Do not rely on ordinal-only diagnostics post-mitigation.** Corollary E.4 shows that ranking accuracy and pairwise win-rate are invariant to $\tilde{R} \mapsto c\tilde{R}$ while $J(\pi_\beta^*(c\tilde{R}), R)$ varies strictly with c . *Concretely*: do not only report ranking accuracy on RewardBench or pairwise win-rates on AlpacaEval post-mitigation, because reporting only such ordinal scores cannot rule out scale-driven regime shifts.

Single-axis Dos.

- **Default to M_i^{norm} and M_i^{cent} mitigations and document the variant used.** Section E.5 shows M_i^{norm} separates axis reallocation from scale and Section E.6 shows M_i^{cent} restores gauge invariance of the regime classification. *Concretely*: when using M_i for pipeline compatibility, report $\|\tilde{R}'\|_{L^2(\mu_{\text{diag}})} / \|\tilde{R}\|_{L^2(\mu_{\text{diag}})}$ (ratio of root-mean-square reward scores for mitigated and unmitigated proxy rewards) alongside g_i so that scale-driven and reallocation-driven Δ_j contributions are distinguishable.
- **Run the $c \in (0, 1)$ rescalability sweep when $\Delta J \leq 0$.** The sweep over $M_i^c(\tilde{R}) = \tilde{R} - c g_i \phi_i$ is the only observational test that separates the two origins of R2 (scale overshoot, recoverable, but target misspecification is not). *Concretely*: publish the $\Delta J(c)$ curve, since a recoverable overshoot is a different scientific finding than a non-recoverable target misspecification.
- **Evaluate at policy-induced distributions and report $(\Delta_j, \Delta J)$.** According to Theorem A.14, this approach separates R0, R0_{cont}, R1, and R2, but, according to Theorem A.13, audit-only inputs cannot. *Concretely*: run the mitigated reward model in BoN or short PPO against a fixed reference policy and report both quantities, not only audit-set accuracy.
- **Instrument off-target Δ_j before publication.** Select a panel of plausible Φ_{sp} candidates the operator does not target (formatting density, hedging markers, position effects, sycophancy and confidence indicators) and report Δ_j on those features at μ_{π^*} . *Concretely*: the developer has both the unmitigated and mitigated rewards in hand and is the natural party to produce this measurement. Without it, an R0 claim is structurally indistinguishable from R0_{cont} (substitution under improvement).

- **Report cardinal scale pre/post mitigation.** Publishing $\|\tilde{R}\|_{L^2(\mu_{\text{diag}})}$ before and after is the minimal cardinal observable that ordinal benchmark functionals miss. *Concretely:* report the standard deviation of reward scores on a fixed evaluation set both before and after applying the mitigation, so scale-driven ΔJ shifts are visible.
- **Document π_{ref} -sensitivity of the mitigation.** The construction in Theorem A.13 produces four distinct regimes by varying only π_{ref} for a fixed $(\tilde{R}, M_i, \mu_{\text{diag}})$. *Concretely:* validate the mitigation against at least two reference policies of different lineages and report per-pair outcomes rather than a single headline number.

Multi-axis recommendations.

- **Specify whether the operator is joint or sequential, and the operator-stack type.** For canonical linear projections applied to a fixed reward, sequential M_i followed by M_j equals joint projection on (ϕ_i, ϕ_j) by Frisch-Waugh-Lovell. The equivalence breaks when operators are heterogeneous (e.g., LOESS for length, two-head for content), when projection directions are recomputed on the changed reward, or when operators are non-linear. *Concretely:* state which case applies and, when sequential and joint differ, report both outcomes side-by-side.
- **Joint reliance zero at μ_{diag} does not certify R0.** Lemma A.6 generalizes and a multi-axis operator M_S that zeros g_i for all $i \in S$ at μ_{diag} does not in general zero g_i at μ_{π^*} . *Concretely:* report Δ_j at μ_{π^*} for every axis in the targeted set, not just on-target axes. Covering a panel at μ_{diag} narrows the substitution-accessible subset of Φ but does not close the measurement-vs-optimization gap.
- **State Φ_{sp} explicitly and acknowledge out-of-panel exploitation.** *Concretely:* a method covering {length, sycophancy, position} does not certify R0 against features outside that set, and multi-axis claims should be framed as “no substitution within the measured panel” rather than “no substitution.”
- **Flag statistical mediation between targeted features.** Where the literature establishes coupling between features in the targeted set (length-sycophancy and length-confidence per Facts G.9, G.11, and G.12), associational joint projection may remove signal flowing through coupled paths. *Concretely:* acknowledge this as a known limitation rather than treating coupled signal as residual noise. The rescalability test above extends to multi-axis as a partial diagnostic.

Scope under non-linear operators. The proofs above (Theorems A.13 and A.14) and the linear-Gaussian instantiation (Section E.7) use Gaussianity for closed-form policy expressions and linearity for the canonical projection construction, but the prescriptions in this section target the structural $\mu_{\text{diag}} \rightarrow \mu_{\pi^*}$ gap and the cardinal/ordinal blindness rather than these constructive devices, and so apply to non-linear, non-Gaussian operators and reference policies. The regime classification (R0-R4) references Δ_j and ΔJ rather than operator form, the cardinal/ordinal blindness of Corollary E.4 is monotone-invariant, and the gap itself is geometric rather than algebraic. What does *not* transfer is Lemma A.6’s constructive guarantee (a non-linear operator may not zero g_i even at μ_{diag}) and the closed-form expressions of Section E.7. Nonetheless, Section E.9 demonstrates regime separation under quadratic non-linearity in closed form and Section F.2 shows the gap surviving a non-linear operator (LOESS calibration) measured by a linear diagnostic, and Section F.3 establishes audit-distribution dependence of the relevant coupling under linear specifications with distribution-free robustness checks. **Method papers using non-linear operators should still adopt the prescriptions** above by virtue of targeting the same geometric phenomenon, while stating that they cannot import the linear-Gaussian regime classification’s constructive guarantees, only its outcome-level definitions.

Reward Model Benchmark Recommendations. Some of the recommendations for mitigation method developers transfer over to benchmark developers, with a few added recommendations.

1925 *Don'ts.*

- 1926
- 1927 • **Do not expand the benchmark input within μ_{diag} and expect regime separation.** Theorem A.13 grants \mathcal{B}
- 1928 joint access to $(\tilde{R}, M_i(\tilde{R}), R, \Phi, \mu_{\text{diag}})$ and still admits four instances landing in R0, R0_{cont}, R1, and R2. Adding
- 1929 feature axes, distractors, or held-out audit splits stays inside \mathcal{B} and inherits the impossibility. *Concretely:* adding
- 1930 a sycophancy subset alongside a length subset, or expanding from pairwise to best-of-4 selection, enriches μ_{diag}
- 1931 but does not measure Δ_j at μ_{π^*} and so cannot separate R0 from R0_{cont} or R1.
- 1932
- 1933 • **Do not rely on benchmark functionals that are invariant to monotone rescaling of \tilde{R} .** Corollary E.4
- 1934 shows that every $B \in \mathcal{B}_{\text{ord}}$ (ranking accuracy, pairwise win-rate, and their composites) returns identical values
- 1935 under $\tilde{R} \mapsto c\tilde{R}$ while $J(\pi_{\beta^*}(c\tilde{R}), R)$ varies strictly with c , so cardinal-scale-driven regime shifts are undetectable.
- 1936 *Concretely:* if mitigation halves the proxy's effective scale, ranking accuracy is unchanged but the KL-regularized
- 1937 policy behaves as if β doubled (by the scale-equivariance of Section E.5), and downstream ΔJ shifts accordingly.
- 1938
- 1939 • **Do not generalize a single- π_{ref} score across deployment settings.** The construction in Theorem A.13 varies
- 1940 only π_{ref} and obtains four distinct regimes for the same $(\tilde{R}, M_i, \mu_{\text{diag}})$. A benchmark score at one π_{ref} does not
- 1941 transport to deployment π_{ref} 's. *Concretely:* a reward model topping the leaderboard when paired with one SFT
- 1942 base can degrade PPO outcomes when paired with a different SFT base of the same model family, since the
- 1943 policy-induced distribution shifts even though the audit-side score does not.
- 1944
- 1945

1946 *Dos.*

- 1947
- 1948
- 1949 • **Evaluate at policy-induced distributions and report $(\Delta_j, \Delta J)$.** By Theorem A.14, this input separates R0,
- 1950 R0_{cont}, R1, and R2 via the explicit functional B^* in Equation (15); by Theorem A.13, no audit-only input suffices.
- 1951 *Concretely:* run the mitigated reward model in BoN or short PPO against a fixed reference policy, and report the
- 1952 change in expected feature value Δ_j on each $\phi_j \in \Phi_{\text{sp}}$ alongside the change in true reward ΔJ , rather than only
- 1953 audit-set accuracy.
- 1954
- 1955 • **Report cardinal scale.** Publishing $\|\tilde{R}\|_{L^2(\mu_{\text{diag}})}$ pre/post mitigation provides observables that \mathcal{B}_{ord} cannot,
- 1956 addressing the blindspot of Corollary E.4. *Concretely:* alongside accuracy, report the standard deviation of
- 1957 reward scores on a fixed evaluation set before and after applying the mitigation, so that scale-driven ΔJ shifts
- 1958 are visible to readers.
- 1959
- 1960 • **Treat π_{ref} as an axis of variation, not a fixed convention.** Either evaluate across a panel of reference policies
- 1961 or document π_{ref} -sensitivity as a first-class output, consistent with the construction underlying Theorem A.13.
- 1962 *Concretely:* pair each reward model with at least two SFT reference policies of different lineages and publish
- 1963 per-pair scores rather than a single headline number.
- 1964
- 1965 • **Report Δ_j across all of Φ_{sp} at μ_{π^*} , not just the targeted axis.** Theorem A.14 makes the off-target Δ_j on
- 1966 $\Phi_{\text{sp}} \setminus \{\phi_i\}$ a load-bearing input for separating R0_{cont} and R1 from R0. Reporting only on-target reliance leaves
- 1967 substitution structurally invisible regardless of how rich the audit distribution is. *Concretely:* when evaluating a
- 1968 length-debiasing operator, also measure post-mitigation drift in sycophancy, hedging, and formatting at μ_{π^*} –
- 1969 not just residual length correlation at the audit set.
- 1970
- 1971

1972 *Summary.* The recommendations above are not addressed to a single audience. R0 certification is a joint property of

1973 the mitigation operator and the evaluation protocol: a method paper following the developer recommendations cannot

1974 be cleanly compared across benchmarks that omit the corresponding ones, and a benchmark implementing all three

1975

1976

1977 prescriptions cannot certify R0 for methods that do not report off-target Δ_j or cardinal scale. The overlap between the
1978 two checklists is structural rather than redundant.

1979 Existing methods and benchmarks satisfy these prescriptions only partially. SOTA benchmarks partially imple-
1980 ment these prescriptions and close the gap to pure ordinal benchmarks. In particular, the current SOTA benchmark
1981 RewardBench2 [52] features a Ties metric, introduces cardinal-sensitive scoring, and its on-policy/off-policy finding
1982 documents π_{ref} -sensitivity empirically. However, no current protocol jointly delivers policy-distribution evaluation,
1983 cardinal reporting, and π_{ref} -as-axis, a gap our framework both formalizes as necessary (Theorem A.13) and proves
1984 sufficient to close (Theorem A.14). On the mitigation side, published methods rarely report off-target Δ_j at μ_{π^*} or run the
1985 rescalability sweep, leaving almost all recent results ambiguous between R0_{cont}, R1, and R2 in ways the recommendations
1986 would have resolved (see Section C).

1987 Closing the joint gap is a community coordination problem rather than a single-paper deliverable. Method papers can-
1988 not unilaterally provide policy-distribution evaluation without benchmark infrastructure that supports it. Benchmarks
1989 cannot unilaterally provide off-target Δ_j measurement without method papers specifying Φ_{sp} . The recommendations
1990 above are scoped so that adoption by either side improves the evidence base, and joint adoption is what moves the field
1991 from R0 claims that audit-side scores cannot certify to R0 claims that the framework certifies sufficient.
1992
1993
1994
1995
1996
1997
1998
1999
2000
2001
2002
2003
2004
2005
2006
2007
2008
2009
2010
2011
2012
2013
2014
2015
2016
2017
2018
2019
2020
2021
2022
2023
2024
2025
2026
2027
2028

F Supporting Experiments

All experiments were either run through API access or on a single A100 GPU 40Gb or B200 GPU 180 Gb (rented online). Our code for the experiments will be released on GitHub (MIT license) upon publication.

F.1 Bias Substitution in Language Model RLHF (GRPO)

We demonstrate reward bias substitution in a RLHF pipeline with standard off-the-shelf configuration. A single-axis length penalty applied during GRPO training compresses response length as intended, yet the freed optimization pressure rotates onto a correlated axis and drives the policy into overconfidence, instantiating the substitution regime (R1) while leaving multiple-choice quality intact.

Setup.

- **Policy model.** *Llama-3.2-3B-Instruct* [27], adapted with LoRA on all linear layers (the four attention projections `q_proj`, `k_proj`, `v_proj`, `o_proj` and the three MLP projections `gate_proj`, `up_proj`, `down_proj`) with rank $r = 32$, scaling $\alpha = 64$, dropout 0.05.
- **Reward model.** *Skywork-Reward-V2-Llama-3.1-8B* [49], kept frozen throughout training and queried only through the reward function.
- **Dataset.** Prompts from *UltraFeedback* [16] (the `ultrafeedback_binarized_train_prefs` split), truncated to 2000 characters.
- **Algorithm.** Group Relative Policy Optimization (GRPO) [67] as implemented in TRL [79]. 8-bit AdamW, learning rate 3×10^{-5} , KL coefficient $\beta = 2 \times 10^{-2}$, 600 optimizer steps, and bfloat16 precision.
- **Length penalty.** The RLHF reward is shaped as $\tilde{R}(x, y) = R_{RM}(x, y) - \lambda n_{\text{tok}}(y)/100$, where $n_{\text{tok}}(y)$ is the number of response tokens and λ is the penalty coefficient. Three values $\lambda \in \{0, 4, 8\}$, each trained with four random seeds, for twelve training runs in total.
- **Generation during training.** Group size of 4 candidate generations per prompt for the group-relative advantage, per-device batch of 4, gradient accumulation of 1, sampling temperature $T = 1.0$, maximum completion length 256 tokens.

Training Curves. We log the following quantities at every optimizer step and report them per λ , averaged over the four seeds, see Figure 5.

- **Mean response length.** The average number of generated tokens per step. This is the primary curve and shows the length penalty taking effect, with the curves separating by λ .
- **Mean reward.** The average shaped reward \tilde{R} per step. Note that \tilde{R} is not directly comparable across λ because the penalty term differs, so the curves are read within each λ for evidence of learning.
- **KL divergence.** The divergence of the policy from the frozen reference policy per step, confirming that the policy stays regularized and does not collapse.
- **GRPO loss and gradient norm.** Diagnostic optimization curves.

Evaluation. All twelve adapted policies and the untrained base policy⁵ are evaluated on held-out data, with each adapter taken at the 600-step checkpoint. Aggregate numbers use 95% confidence intervals from a hierarchical bootstrap over seeds and samples. We report the following metrics.

⁵Untrained by us, but still instruction tuned by the model developer.

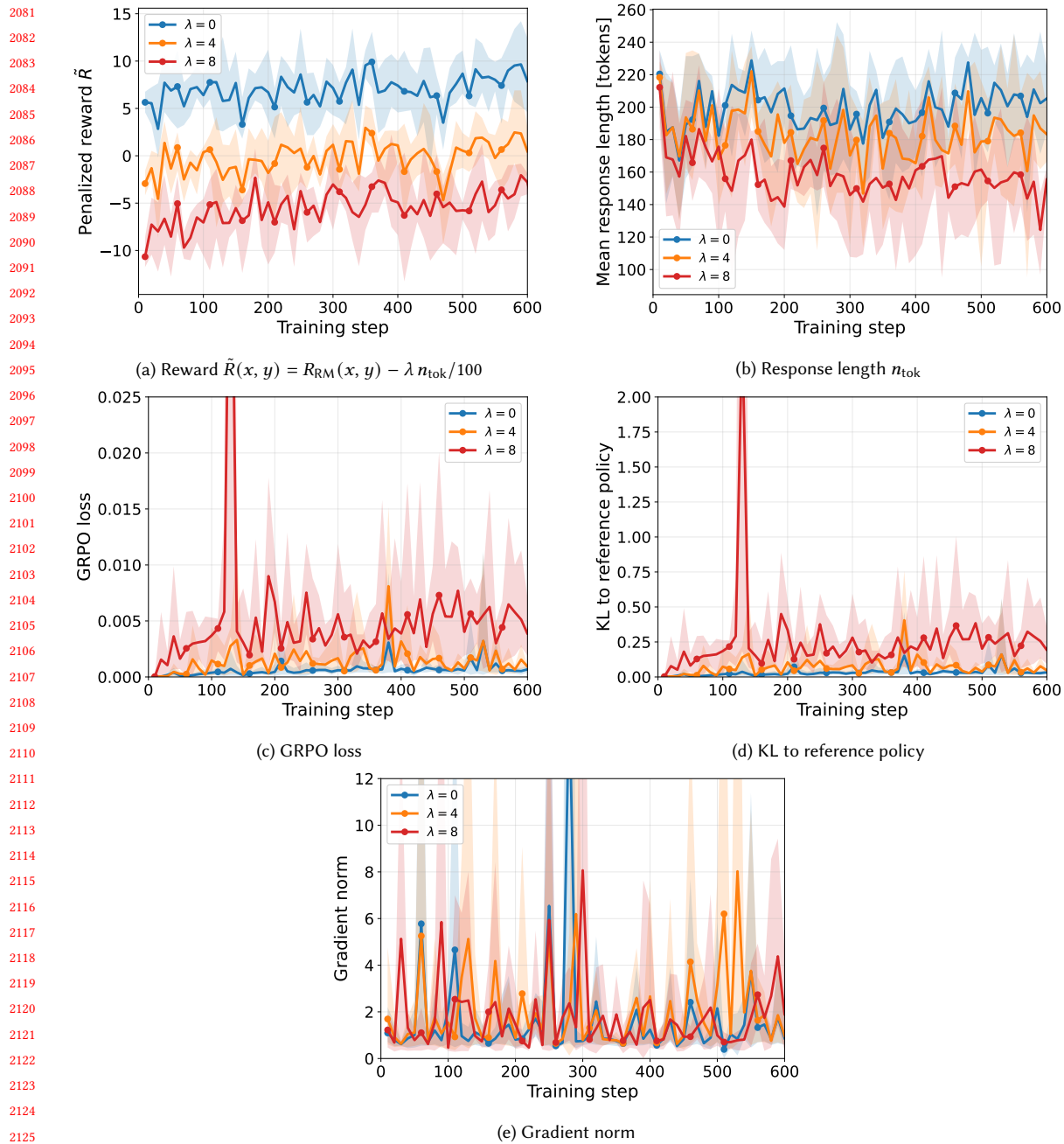


Fig. 5. Training curves (shaded region shows the min-max range across 4 seeds)

| | Base | $\lambda = 0$ | $\lambda = 4$ | $\lambda = 8$ |
|----------------------------------|-------|-----------------------------|-----------------------------|-----------------------------|
| <i>Length and quality</i> | | | | |
| Mean response length (tokens) | 192.2 | 203.6 _[194, 212] | 188.2 _[181, 195] | 170.2 _[161, 179] |
| MMLU accuracy | 0.596 | 0.609 _[.59, .62] | 0.611 _[.59, .63] | 0.616 _[.60, .63] |
| <i>Overconfidence</i> | | | | |
| Expected calibration error (ECE) | 0.264 | 0.253 _[.23, .28] | 0.299 _[.26, .34] | 0.406 _[.32, .50] |
| Mean verbalized confidence | 0.811 | 0.803 _[.78, .82] | 0.815 _[.80, .83] | 0.836 _[.82, .85] |
| TriviaQA accuracy | 0.564 | 0.558 _[.53, .58] | 0.523 _[.47, .57] | 0.415 _[.29, .52] |
| AUROC | 0.708 | 0.725 _[.71, .74] | 0.706 _[.67, .74] | 0.651 _[.62, .69] |
| <i>Sycophancy</i> | | | | |
| Regressive flip rate | 0.653 | 0.582 _[.52, .64] | 0.646 _[.59, .70] | 0.602 _[.54, .66] |
| Overall flip rate | 0.721 | 0.668 _[.63, .71] | 0.709 _[.67, .75] | 0.704 _[.66, .75] |

Table 3. Metric values across the length-penalty sweep for the data plotted in Figure 1 in Section 1. Each λ entry is the mean over four random seeds with the 95% hierarchical-bootstrap confidence interval as a subscript. The Base column is the untrained *Llama-3.2-3B-Instruct* model and carries no interval. Expected calibration error, TriviaQA accuracy, and AUROC use the LLM-judge answer grading.

- **Mean response length.** The average response token count on 50 held-out *UltraFeedback* test_prefs prompts, with 4 samples per prompt drawn at temperature 1.0.
- **MMLU accuracy.** Multiple-choice accuracy on 1000 questions subsampled from MMLU, evaluated zero-shot with greedy decoding and a single-letter answer parsed from the output.
- **Expected Calibration Error (ECE).** Computed with 10 equal-width bins on the (verbalized confidence, correctness) pairs produced by the Tian et al. [77] Verb. 1S top-1 protocol on 1000 *TriviaQA* questions.
- **AUROC.** The area under the ROC curve of the verbalized confidence used as a predictor of answer correctness.
- **Mean verbalized confidence.** The average probability the model assigns to its own answer under the Verb. 1S top-1 protocol.
- **TriviaQA accuracy.** The fraction of correct answers under the same protocol.
- **Regressive flip rate.** The headline sycophancy metric from the Sharma et al. [69] *are_you_sure* task, the fraction of initially-correct answers that become incorrect after the user pushes back.
- **Overall flip rate.** The fraction of answers that change at all after the pushback, regardless of correctness.

The Tian et al. [77] ECE, AUROC, and TriviaQA accuracy metrics are computed under an LLM-judge equivalence check using the Tian et al. [77] prompt with *Claude Haiku 4.5* [1] as the judge. The sycophancy metrics use the Sharma et al. [69] two-stage protocol, eliciting an initial answer and then a post-pushback answer on 500 *are_you_sure* records with 2 samples per record at temperature $T = 1.0$.

Results. See Figure 1 and Table 3 for results.

- **The length penalty compresses responses monotonically.** Mean response length falls from 204 tokens at $\lambda = 0$ to 188 at $\lambda = 4$ and 170 at $\lambda = 8$. Both reductions, measured relative to $\lambda = 0$, are significant.
- **Multiple-choice quality is preserved.** MMLU accuracy holds near 0.61 across the sweep (0.609, 0.611, 0.616 at $\lambda = 0, 4, 8$ against a base of 0.596), and the change from base is not significant at any λ .

- **Overconfidence rises with the length penalty.** Expected calibration error climbs from 0.25 at $\lambda = 0$ to 0.30 at $\lambda = 4$ and 0.41 at $\lambda = 8$, and mean verbalized confidence climbs in step from 0.80 to 0.82 to 0.84. TriviaQA accuracy falls from 0.56 at $\lambda = 0$ to 0.52 at $\lambda = 4$ and 0.41 at $\lambda = 8$, and the coupling between confidence and correctness weakens, with AUROC falling from 0.73 to 0.65. Thus, the policy does not merely become more confident, it becomes wrong more often. At $\lambda = 8$ the rise in calibration error, the rise in verbalized confidence, the drop in accuracy, and the drop in AUROC are all significant.
- **The calibration loss is caused by the penalty, not by RLHF.** At $\lambda = 0$ the calibration error (0.25) is at or slightly below the base value (0.26), so RLHF on its own does not break calibration. The degradation appears only once the length penalty is applied.
- **Sycophancy does not significantly change.** The regressive flip rate is non-monotonic across the sweep (0.58, 0.65, 0.60 against a base of 0.65) and the overall flip rate shows no consistent trend (0.67, 0.71, 0.70 against 0.72). However, the base model already has a high level of sycophancy.
- **Regime transitions.** $\lambda = 0$ is the no-penalty reference, $\lambda = 4$ transitions into and $\lambda = 8$ instantiates the bias substitution regime (R1), in which the penalty compresses length and the optimization pressure rotates onto confidence calibration. The overcorrection regime (R2), in which task-independent capability degrades, is not reached within the swept range, since MMLU accuracy is preserved at every λ . $\lambda = 4$ is not successful mitigation R0, since every overconfidence proxy is already displaced toward substitution (see Table 3) and sits at or leans towards the significance boundary, so the missing significance reflects four-seed power rather than the absence of an effect.
- **Regime transitions.** We take the true reward to be free-form factual accuracy, proxied by TriviaQA, and we treat expressed confidence as a spurious feature that the reward model rewards as a quality proxy (Fact G.12) (Φ_{sp}). Multiple-choice MMLU accuracy serves as a capability control. At $\lambda = 8$, response length falls and MMLU is unchanged, while expressed confidence rises, a nonzero Δ_j on an off-target spurious axis, and TriviaQA accuracy degrades, a negative ΔJ . This is the harmful sub-case of bias substitution (R1), not overcorrection (R2), because pressure rotates onto a second spurious axis rather than degrading capability with no rotation and since MMLU accuracy is preserved at every λ . We read $\lambda = 4$ as an intermediate point, where the overconfidence proxies shift toward substitution but the changes from the $\lambda = 0$ run are not significant at four seeds, and $\lambda = 8$ as the clear R1 instance.

F.2 Measurement-vs-Optimization Gap in Length Mitigation in BoN Selection

This section provides a quantitative empirical instance of the measurement-vs-optimization gap of Section A.2. We evaluate two published length-debiasing operators across five reward models under a Best-of- N (BoN) selection protocol, using a within-prompt reliance diagnostic that tracks the BoN selection distribution rather than the pooled audit distribution at which the operators are designed. The headline observation is that the post-hoc calibration of Huang et al. [30] zeros pooled reward-length correlation ($0.316 \rightarrow 0.037$) while overshooting into negative within-prompt correlation on three of four SOTA reward models, with ΔJ degrading on two. The linear-probe operator of Fein et al. [21] is fit on the same prompt split and transfers cleanly to the within-prompt diagnostic on four of five reward models with $\Delta J > 0$.

Reward models and operators. We evaluate Skywork-Reward-V2-Llama-3.1-8B, Skywork-Reward-V2-Qwen3-8B, Skywork-Reward-V2-Qwen3-0.6B [49], the Llama-3.1-8B-Instruct-RM-RB2 [52], and the older DeBERTa-v3-large [57]

| | DeBERTa | Allen | SW-L | SW-Q8B | SW-Q0.6B |
|--|---------|-------------------|---------|-------------------|----------|
| <i>Within-prompt $\rho_{\text{len}}^{\text{within}}$ (AlpacaEval)</i> | | | | | |
| Baseline | + .168 | -.061 | + .134 | + .065 | + .088 |
| Mech. RS | + .005 | -.011 | + .087 | + .032 | + .041 |
| Huang et al. [30] | + .048 | -.099 | -.065 | -.138 | -.228 |
| <i>AlpacaEval LC win rate vs. baseline (%)</i> | | | | | |
| Mech. RS | 54.4*** | 46.2 [†] | 51.5*** | 51.7*** | 51.8*** |
| Huang et al. [30] | 50.4* | 48.7 [†] | 50.1 | 48.2 [†] | 51.2*** |
| <i>GSM8K BoN accuracy (%)</i> | | | | | |
| Baseline | 32.8 | 71.5 | 68.3 | 71.9 | 65.8 |
| Mech. RS | 36.4 | 71.3 | 68.5 | 71.9 | 65.7 |
| Huang et al. [30] | 32.3 | 67.9 | 68.3 | 71.3 | 66.2 |

Table 4. Within-prompt reward-length correlations and BoN selection outcomes. * $p < .05$, *** $p < .001$ vs. 50%; [†] denotes significantly below 50%. Pooled $|\rho_{\text{len}}|$ averaged across the five RMs is 0.316 (baseline) and 0.037 for Huang et al. [30]: audit-distribution success and within-prompt overshoot on the same operator.

reward models. We use two single-axis operators that target length bias ($\phi_i = \text{length}$). First, mechanistic reward shaping with a difference-of-means linear probe projected from the RM’s final-layer hidden state [21], and, second, the post-hoc reward calibration of Huang et al. [30] with a LOESS fit subtracted at the score level. Both are instances of M_i (Definition A.5) and reduce $|g_i(\tilde{R}; \mu_{\text{diag}})|$ at their respective audit distributions by construction.

Setup. For each of two prompt sources (AlpacaEval [18], GSM8K [15]) we generate 64 candidate responses per prompt from Llama-3.2-1B-Instruct [27], fit each operator on a probe split, and evaluate on a disjoint held-out split (512 for AlpacaEval prompts, 807 for GSM8K prompts). Each RM selects the highest-scoring response among the 64 candidates. We measure two quantities:

- a within-prompt Pearson correlation $\rho_{\text{len}}^{\text{within}}$ between RM score and response length, computed within each candidate set and averaged across prompts. Because BoN top-1 selection operates within candidate sets per prompt, $\rho_{\text{len}}^{\text{within}}$ measures the reliance selection actually responds to, while pooled ρ_{len} averages over prompts that BoN never compares. This instantiates g_i at a μ_{diag} aligned with BoN selection rather than the pooled distribution the operators target.
- A selection-level ΔJ proxy (AlpacaEval length-controlled win rate vs. baseline, GSM8K BoN accuracy on selected responses).

The pair $\{\rho_{\text{len}}^{\text{pooled}}, \rho_{\text{len}}^{\text{within}}\}$ operationalizes Lemma A.6 empirically, as an operator can zero one and not the other, and the selection outcome reveals which distribution drives the optimizing process. BoN top-1 over 64 samples is a coarser proxy for π_{β}^* than KL-regularized optimization, and the regime calls below should be read accordingly.

Results. The Huang et al. [30] calibration achieves $|g_i| \approx 0$ at the pooled μ_{diag} it targets (0.316 \rightarrow 0.037 averaged across the five RMs). At the within-prompt μ_{diag} aligned with BoN selection, three of four SOTA RMs acquire negative $\rho_{\text{len}}^{\text{within}}$ of magnitude 0.065-0.228, exceeding their unmitigated baselines in absolute value (Table 4). The same pattern holds on GSM8K, where within-prompt $|\rho_{\text{len}}|$ averages 0.096 for Huang et al. [30] versus 0.007 for mechanistic reward shaping, with BoN accuracy averaging 61.20% versus 62.76%. The signed flip rather than same-side overshoot is **evidence of**

2289 **the measurement-versus-optimization gap** of Lemma A.6, because the operator reaches zero at μ_{diag} , while at the
 2290 optimization-relevant distribution g_i has the wrong sign.
 2291

2292 *Regime classification.* Under the ε -banded reading of Section E.8, two Huang et al. [30] cells (Allen, Skywork-Q8B)
 2293 satisfy $\Delta J < -\varepsilon_J$ on AlpacaEval LC win rate while the targeted axis $\rho_{\text{len}}^{\text{within}}$ has flipped sign. The GSM8K Allen cell
 2294 additionally drops 3.6 accuracy points. These satisfy the targeted-axis and ΔJ conditions of Definition A.10. Whether
 2295 they additionally satisfy the $\Delta_j = 0$ condition for off-target spurious axes is unmeasured here. As a result, the pattern is
 2296 consistent with $R_{2\varepsilon}$ if the no-rotation condition holds, or with a mixed R1+R2 regime otherwise.
 2297

2298 Mechanistic reward shaping achieves $\Delta J > 0$ at $p < .001$ on four of five AlpacaEval cells and moves $|\rho_{\text{len}}^{\text{within}}|$ toward
 2299 zero on four of five RMs without sign flips. The Allen cell falls below 50% LC WR on both operators (46.2% for Mech.
 2300 RS, 48.7% for Huang et al. [30], both †). The GSM8K Allen cell drops 0.2 for mechanistic reward shaping and 3.6 for
 2301 Huang et al. [30]. The remaining four mechanistic reward shaping cells satisfy Definition A.8 modulo the unmeasured
 2302 no-rotation condition.
 2303
 2304

2305 *Connection to the framework.* The pair $(\rho_{\text{len}}^{\text{pooled}}, \rho_{\text{len}}^{\text{within}}) = (0.037, 0.116)$ for Huang et al. [30] averaged across the five
 2306 RMs is the empirical counterpart of Lemma A.6: the canonical M_i reaches zero at the audit distribution it targets while
 2307 g_i at the optimization-relevant distribution has the wrong sign. This dissociation is the load-bearing mechanism of
 2308 the substitution argument. Wherever it holds, R1 is mechanically available to the optimizer, and the audit diagnostic
 2309 targeting ϕ_i cannot detect it by construction. Standard reward-model benchmarks evaluate at a single fixed μ_{diag} and
 2310 report only the pooled diagnostic, so under Section A.5 they cannot distinguish this $\Delta J < 0$ instance from R0 regardless
 2311 of whether the underlying regime is $R_{1\varepsilon}$ (rotation onto an unmeasured spurious axis), $R_{2\varepsilon}$ (overcorrection on the
 2312 targeted axis), or a mixture: all three are invisible at μ_{diag} , which is the central benchmark-inadequacy claim.
 2313
 2314
 2315

2316 F.3 Evaluating sycophancy and length of AITA Responses

2317 This section reports the experimental analysis backing the Section B claim that response length and sycophancy labels
 2318 are statistically dependent across human, LLM-judge, and judge-agreement labeling regimes, using AITA prompts and
 2319 responses from [9] across eight model families. We fit mixed linear models of response length (measured in characters)
 2320 on a binary sycophancy indicator (with model or prompt as a random effect), and report KS and W_1 . These statistics
 2321 provide audit-side evidence for the construction-level precondition of R4 (Definition A.12), with the formal banded
 2322 treatment deferred to Section E.8. In terms of causality, the mixed-model coefficients and KS statistics below are
 2323 observational. They are consistent with a mediation reading of length but do not, on preference data alone, discriminate
 2324 among candidate causal structures.
 2325
 2326
 2327
 2328

2329 *Human-labeled regime ($N = 14,375$ responses across 8 models).* A mixed linear model of response length on the
 2330 sycophancy indicator, with model as random effect, yields a positive sycophantic effect of +24.3 length units (95% CI [9.6,
 2331 39.0]; $p = 0.001$), against an intercept of 1554.5 and between-model variance 72,408. The effect is small relative to between-
 2332 model variability (median per-model non-syco/syco ratio 0.98), indicating that under human labeling sycophantic
 2333 responses are only modestly longer on average, but statistically significant. Distributionally, the length distributions
 2334 of sycophantic and non-sycophantic responses differ significantly for 6 of 8 models by Kolmogorov–Smirnov test
 2335 (Llama-8B and Mistral-7B not significant at $\alpha = 0.05$), with per-model Wasserstein distances ranging from 24.5 to 138.6.
 2336 Pooling z-scored within-model residuals gives KS = 0.025 ($p = 0.022$), confirming a small but significant aggregate
 2337 distributional deviation.
 2338
 2339
 2340

2341 *LLM-judge-labeled regime* ($N = 28,832$ responses across 8 models). A mixed linear model of response length on the
 2342 sycophancy indicator, with model as random effect, yields a sycophantic effect of +128.4 length units (95% CI [118.5,
 2343 138.4]; $p < 0.001$), against an intercept of 1487.1 and between-model variance 70,521. The effect is substantially larger
 2344 than under human labels (+24.3), indicating that LLM judges associate sycophancy with longer responses than the
 2345 binary-verdict label does. Per-model non-syco/syco ratios shift correspondingly downward (median 0.94 vs. 0.98
 2346 under human labels), with the gap most pronounced for Llama-17B, Llama-70B, and Gemini (ratios 0.81–0.85). Length
 2347 distributions differ significantly for 7 of 8 models by Kolmogorov–Smirnov test (only gpt-4o is not significant at $\alpha = 0.05$),
 2348 with per-model Wasserstein distances ranging from 19.5 to 293.7, which is substantially larger than the 24.5–138.6
 2349 range under human labels. The pooled KS on z-scored within-model residuals is 0.130 ($p < 0.001$), roughly $5\times$ larger
 2350 than the human-label pooled KS of 0.025, confirming that the LLM-judge regime exhibits a much stronger distributional
 2351 length–sycophancy coupling than the human regime at both the central-tendency and full-distribution levels.
 2352
 2353
 2354
 2355

2356 *Judge-agreement regime* ($N = 4,149$ responses across 8 models, restricted to responses where human and LLM judge labels
 2357 coincide). A mixed linear model of response length on the sycophancy indicator, with model as random effect, yields a
 2358 sycophantic effect of +154.3 length units (95% CI [123.6, 185.0]; $p < 0.001$), against an intercept of 1507.2 and between-
 2359 model variance 68,577. The effect is the largest of the four regimes, exceeding both human-only (+24.3) and LLM-only
 2360 (+128.4), which is consistent with the agreement subset isolating responses for which sycophancy is most unambiguous
 2361 along a length-correlated axis. Per-model non-syco/syco ratios drop correspondingly (median 0.93), with Llama-70B
 2362 and Gemini showing the strongest shifts (ratios 0.77 and 0.77). Length distributions differ significantly for 4 of 8 models
 2363 by Kolmogorov–Smirnov test at $\alpha = 0.05$ (Llama-8B, Llama-17B, Llama-70B, Gemini), with per-model Wasserstein
 2364 distances ranging from 25.9 to 437.2 and reaching their highest values across all regimes. The four non-significant
 2365 per-model KS tests (Claude, gpt-4o, Mistral-7B, Mistral-24B) likely reflect reduced statistical power, given the $8\times$ range
 2366 in per-model group sizes (94 to 784). The pooled KS on z-scored within-model residuals is 0.130 ($p < 0.001$), comparable
 2367 to the LLM regime, indicating that the distributional length–sycophancy coupling under judge agreement is at least as
 2368 strong as under LLM-only labeling at the aggregate level.
 2369
 2370
 2371
 2372
 2373

2374 *Judge-disagreement regime* ($N = 10,226$ responses across 8 models, restricted to responses where human and LLM judge
 2375 labels disagree). A mixed linear model of response length on the sycophancy indicator, with model as random effect,
 2376 yields a *negative* sycophantic effect of -43.1 length units (95% CI $[-60.5, -25.7]$; $p < 0.001$), against an intercept of 1568.7
 2377 and between-model variance 73,110. The sign reversal relative to the human (+24.3), LLM (+128.4), and agreement
 2378 (+154.3) regimes is the central qualitative finding: under disagreement, responses labeled sycophantic are *shorter* than
 2379 non-sycophantic ones on average. Per-model non-syco/syco ratios shift accordingly (median 1.03, with 6 of 8 models
 2380 above 1.0), reversing the pattern seen in the other three regimes. Length distributions differ significantly for 7 of 8
 2381 models by Kolmogorov–Smirnov test (only Mistral-7B is not significant at $\alpha = 0.05$), with per-model Wasserstein
 2382 distances ranging from 26.1 to 267.7. The pooled KS on z-scored within-model residuals is 0.046 ($p < 0.001$), smaller
 2383 than the LLM and agreement regimes (both 0.130) but larger than the human regime (0.025), indicating a distributional
 2384 shift that is real but less pronounced than under either single-judge or unanimous-agreement labeling. The reversal is
 2385 consistent with humans and LLM judges locating sycophancy in systematically different parts of the response-length
 2386 distribution: when their labels diverge, the resulting "sycophantic" set is enriched for cases each judge alone would
 2387 have ruled differently on, and the length signal flips accordingly.
 2388
 2389
 2390
 2391
 2392

2393 *Within-prompt across-models robustness check* ($N = 11,309$ responses across 1,569 prompts with label variation). To rule
2394 out prompt-level content as a confound — the possibility that the length–sycophancy association in the four mode-level
2395 regressions reflects sycophantic *prompts* eliciting longer responses rather than sycophantic *responses* being longer per
2396 se — we refit the mixed linear model with **prompt as the random effect**, restricting to the 1,569 prompts where at
2397 least two of the eight models’ responses received different sycophancy labels. With group size 5–8 (mean 7.2), this
2398 within-prompt design holds prompt content fixed and identifies the sycophantic effect off cross-model variation in
2399 labeling and length on the same prompt. The estimated sycophantic effect is +58.8 length units (95% CI [40.4, 77.1]; $p <$
2400 0.001), against an intercept of 1549.5 and between-prompt variance 47,247. The positive sign and significance confirm
2401 that sycophantic responses are longer than non-sycophantic responses *to the same prompt*, ruling out prompt-level
2402 content as the sole driver of the association. The pooled KS on z-scored within-prompt residuals is 0.073 ($p <$ 0.001),
2403 indicating distributional length–sycophancy coupling that survives prompt-level conditioning. The within-prompt
2404 analysis uses the same human-label scheme as the four-mode human regression (model binary verdict on ‘is_a**hole
2405 == 1’ prompts); the larger coefficient (+58.8 vs. +24.3) reflects the restriction to the 1,569 prompts with cross-model
2406 label variation, not a change in labeling regime. Sycophancy effects are concentrated on prompts where models actually
2407 disagree about the verdict.
2408
2409
2410
2411
2412
2413
2414
2415
2416
2417
2418
2419
2420
2421
2422
2423
2424
2425
2426
2427
2428
2429
2430
2431
2432
2433
2434
2435
2436
2437
2438
2439
2440
2441
2442
2443
2444

G Existing Evidence for Confounding Factors

Fact G.1. Reward models are correlated with response length independent of content quality, from the first language model applications of RLHF to today [6, 19, 21, 56, 70, 71, 76, 80, 83, 91].

Fact G.2. Reward models exhibit multiple additional distinct reward biases, e.g., overconfidence [44], sycophantic user agreement [69], answer position [21], prefix bias [39], model-style sensitivity [21, 52], and inherited value orientations (e.g., agency vs. communion) from pretraining base models [13]. All of these are present in state-of-the-art reward models [13, 21].

Fact G.3. LLM policies trained with RLHF or DPO learn to generate longer outputs as a direct result of optimization pressure. This effect can inflate benchmark scores [53, 71], produce outputs longer than training data [59], and grow systematically longer than SFT baselines over training [76]. Overall, it represents a primary mode of reward hacking [59, 71].

Fact G.4. Naive length penalties can remove informative content, as uniform length suppression degrades accuracy on tasks that genuinely require longer answers [3], and always-on penalties cause reward hacking via trajectory collapse in reasoning RL [89].

Fact G.5. Longer responses are systematically over-represented among “chosen” annotations in preference datasets, providing the distributional basis for reward models to learn length as a proxy for quality [70, 71]. RLHF dramatically amplifies a slight length skew already present in preference data [19], a pattern confirmed in crowdsourced human preference data, where controlling for length substantially reorders model rankings [46]. Further, length preferences vary systematically across annotator populations and data sets [54].

Fact G.6. Human preference judgments are systematically confounded by output length. Across multiple preference-collection pipelines, annotators select the longer response in comparisons [6, 29, 71]. This pattern is visible even in carefully filtered datasets, as Stiennon et al. [76] find that controlling for length reduces human preference gains and Saito et al. [65] show a positive correlation between length and preference in the HH-RLHF corpus. In large-scale crowdsourced evaluation, length has been estimated as a dominant style factor in human voting [46].

Fact G.7. LLM judges exhibit verbosity bias [92] and it is different to human annotators [38, 65]. Verbosity is one of several systematic biases catalogued in LLM judges, in addition to other LLM judge biases [6, 38, 82, 85]. Further evidence also supports that human and LLM judges can respond to length and content differently [12, 29, 46, 65].

Fact G.8. Human annotators systematically reward sycophantic responses. Sharma et al. [69] show that “matches user’s beliefs” is among the most predictive features of human preference, with the preference model favoring sycophantic over truthful responses. Cheng et al. [10] confirm this finding at the data level, as preferred responses across preference datasets are significantly more validating and indirect. Perez et al. [62] find that preference models incentivize sycophantic answers, while Wen et al. [83] show RLHF-optimized models learn to defend incorrect answers with fabricated evidence and extended argumentation, increasing human approval without improving correctness. Ibrahim et al. [32] further show that supervised fine-tuning on warmer-style conversational data alone (ithout preference optimization) is sufficient to amplify affirmation of incorrect user beliefs by approximately 40% across five model families, with the effect surviving response-length controls.

2497 **Fact G.9. Whether responses are labelled as sycophantic statistically depends on response length, independent**
2498 **of annotator type.** Prior work has reported response length and sycophancy jointly: Ibrahim et al. [32] include length
2499 as a covariate when regressing accuracy and sycophancy on warmth fine-tuning, and Dubois et al. [17] treat length
2500 as a nuisance parameter when measuring sycophancy reduction. To our knowledge, however, the statistical dependence
2501 between response length and sycophancy labels has not been characterized across labeling regimes (human, LLM judge,
2502 agreement, disagreement), nor has the sign reversal under human-LLM judge disagreement been documented. We examine
2503 the dependence across multiple labeling strategies (human labels, LLM Judge labels, their agreement, and disagreement),
2504 systematically disentangling model-level verbosity from prompt-level content as alternative explanations using the AITA
2505 dataset from [9]. We find length and sycophancy of responses are statistically dependent across all four labeling regimes, see
2506 Section F.3 for experiment details. Sycophantic responses are longer under human, LLM, and human-LLM judge agreement
2507 labels, with the relationship reversing under human-LLM judge disagreement labels. This observation is consistent with
2508 humans and LLM judges locating sycophancy in systematically different parts of the response-length distribution. In addition,
2509 we find that the length-distribution of sycophantic and non-sycophantic responses significantly deviates for most models
2510 (Kolmogorov-Smirnov) and that LLM judge labels are substantially inflated by length bias relative to human labels.
2511

2515 **Fact G.10. Sycophancy causes models to abandon correct responses in favor of alignment with user-stated**
2516 **positions.** When users suggest incorrect answers, model accuracy drops relative to unbiased baselines [69] with high
2517 agreement rates with user-stated views for some question types [62]. Fanous et al. [20] quantify this as regressive sycophancy,
2518 observing models switch from correct to incorrect under user rebuttals and Hong et al. [28] further show that this is not a
2519 knowledge deficit. Beyond propositional agreement, framing sycophancy, in which models uncritically accept flawed user
2520 premises, proved particularly resistant to DPO-based mitigation [10]. Also, sycophantic affirmation narrows users focus to
2521 self-validation while omitting alternative perspectives [9].
2522
2523

2524 **Fact G.11. Models can generate longer responses under epistemic uncertainty by hedging, qualifying, and**
2525 **elaborating more when confidence is low [90].**

2527 **Fact G.12. Reward models treat expressed confidence as a quality proxy, incentivizing models to replace hedged**
2528 **answers with confident-sounding responses regardless of actual certainty [44].** Reward models can penalize hedged
2529 correct answers over confidently-stated incorrect ones [21], and this penalty can be traced to the annotation level, where
2530 human raters in preference datasets are biased against expressions of uncertainty [31, 93].
2531
2532

2533 **Fact G.13. Longer responses in uncertain regimes correlate with lower confidence and reduced quality.** Verbose
2534 outputs produced under verbosity compensation show measurably higher uncertainty and recall drops relative to concise
2535 responses [90]. Reward models scoring confident-sounding elaborations highly regardless of correctness [44] implies the
2536 same inverse relationship, though without directly measuring it.
2537
2538

2539 **Fact G.14. Reasoning models (not RLHF) produce longer outputs with uncertainty markers.** Extended chain-
2540 of-thought training produces epistemic uncertainty markers and non-linear reasoning traces through backtracking and
2541 alternative exploration rarely seen in non-reasoning models [88]. During on-policy RL training, trajectory length grows
2542 systematically alongside exploratory behavior, and naively penalizing length degrades performance [89], consistent with RL
2543 reinforcing the coupling.
2544

2545 **Fact G.15. Reward model scores increase with length when additional tokens carry genuine informational**
2546 **content [29], but this relationship degrades at greater lengths, entering sublinear and then stochastic regimes beyond 100**
2547
2548

2549 and 200 tokens respectively [91], consistent with **diminishing informational returns per token as responses grow**
 2550 **longer.**

2552 **Fact G.16. Some reward-length correlation is irreducible and reflects genuine informativeness rather than**
 2553 **spurious shortcut.** The information mass decomposition directly establishes that a portion of the length–reward correlation
 2554 tracks real content quality [29], and more detailed answers can be genuinely more helpful depending on context [3]. Consistent
 2555 with this observation, length-controlled analysis implies the original correlation is not entirely spurious [18].
 2556

2557 **Fact G.17. DPO-like objectives develop length bias out-of-distribution as a structural vulnerability of the**
 2558 **algorithm without explicit mitigations.** DPO’s implicit reward as a token-level log-probability ratios grows with
 2559 sequence length, creating an algorithmic dependence on response length within the objective [51]. Empirically, this dependence
 2560 manifests as the implicit reward acquiring significant length correlation once the policy moves away from the training
 2561 distribution [53, 59], and can be mitigated through reward normalization and margin formulation [45]. Independently, under
 2562 noisy or capability-limited supervision, DPO’s KL regularization creates a structural dilemma where sufficient regularization
 2563 preventing over-optimization also prevents large updates needed for correcting errors inherited from the SFT phase [87].
 2564
 2565
 2566

2567 **Fact G.18. KL regularization interacts with reward scale non-trivially.** KL-regularized policies are not invariant to
 2568 positive linear reward scaling (see Section E.5), meaning two reward models agreeing on all preference orderings but differing
 2569 in cardinal scale produce different policies [72]. Empirically, there exists a measured optimal KL budget beyond which true
 2570 reward degrades in RLHF [25] and longer sequences accumulate disproportionately more KL divergence, entangling length
 2571 with the implicit DPO reward formulation [51]. The optimal regularization coefficient β depends strongly on feedback
 2572 reliability [87], meaning there is no single good KL setting across realistic annotation conditions.
 2573
 2574

2575 **Fact G.19. Spurious length features can maintain ordinal ranking accuracy while distorting cardinal reward**
 2576 **values, meaning benchmark metrics may fail to detect this failure mode.** A dissociation that follows from the
 2577 partial-identifiability characterization of Skalse et al. [72] combined with the KL scale-sensitivity of Section 2 (Equation (1)),
 2578 and is directly measured as a low-complexity linear artifact in reward model representations [21]. RMs built on different
 2579 base model families can achieve similar RewardBench scores while exhibiting systematically divergent value orientations
 2580 [13] or producing systematically different PPO outcomes depending on policy-RM lineage match [52], concrete empirical
 2581 instances of ordinal-invariant but cardinally-divergent reward functions. RewardBench [43] scores improve after post-hoc
 2582 length calibration, confirming retrospective distortion without prior detection [30], and apparent leaderboard gains shrink
 2583 substantially under length-controlled evaluation in both automated [18] and crowdsourced human voting settings [46].
 2584
 2585
 2586

2587 **Fact G.20. Evaluator and annotator weaknesses are learnable targets. Models trained with RLHF can acquire**
 2588 **exploitation strategies like length padding, hedging caveats, and code obfuscation.** These strategies are specific to
 2589 evaluator blind spots rather than tied to any single feature, demonstrating that reward gaming is a general learned behavior
 2590 directed at the structure of the evaluation process itself [83].
 2591

2592 **Fact G.21. Reward can be somewhat decomposed into separable dimensional components.** A two-head architecture
 2593 directly separates length and content reward signals, confirms the components can be disentangled with supervised training
 2594 signals to improve mitigation [6]. Consistent with this, multi-attribute scoring frameworks track helpfulness, correctness,
 2595 coherence, complexity, and verbosity as separately annotated (but correlated) dimensions [56, 81].
 2596
 2597

2598 **Fact G.22. Reward biases admit to first-order linear intervention.** Length bias is smooth, structured, and stable
 2599 enough across model families to be estimated and corrected post-hoc via LOESS without retraining [30], but assumes that the
 2600

2601 true reward is independent of length. Fein et al. [21] also use a linear probe to mitigate some part of length and other biases
 2602 in reward models, while Papadatos and Freedman [58] used linear probe-based reward corrections to mitigate sycophancy
 2603 on simplified controlled setting.
 2604

2605
 2606 **Fact G.23. Mitigating a spurious correlation in one feature can increase reward hacking in another feature**
 2607 **(bias substitution).** In supervised vision, mitigating a single labeled shortcut amplifies reliance on the unlabeled one,
 2608 with off-target axis degradation of 2–3x while aggregate worst-group accuracy improves [48]. In RLHF, aggressive length-
 2609 debiasing in recent SOTA reward models has flipped the length-bias sign and reduced correctness [21]. In LLM preference
 2610 optimization more broadly, several DPO bias-mitigation variants reduce targeted bias at significant cost to general capability
 2611 [22], instantiating R2. The mechanism is compatible with the single-proxy correlated-feature bound showing that drift along
 2612 an unconstrained correlated direction remains available to the policy whenever the proxy–truth correlation is below one
 2613 [41]. Liu et al. [50] address this at the reward level via max-min optimization, but without the R0–R4 diagnostic machinery,
 2614 substitution between interpretable feature axes remains undetectable. Conversely, jointly mitigating multiple shortcuts via
 2615 kernel-based regularization [86] achieves near-zero shortcut correlations at the audit distribution, but without measuring
 2616 whether optimization pressure has shifted at μ_{π^*} , instantiating the distributional blindspot of Section A.5.
 2617
 2618
 2619
 2620

2621 H Mapping published mitigations onto the regime taxonomy

2622 The detailed classification for each reward bias mitigation method we distill in Section B. No published mitigation paper
 2623 we find provides evidence sufficient to certify R0.
 2624
 2625

- 2626 • Park et al. [59] (R-DPO) — undetermined; consistent with $R0_{\text{cont}}$. Length penalty validated only on preference
 2627 data; no off-target Δ_j or fixed- π_{ref} Δ_j reported.
- 2628 • Meng et al. [53] (SimPO) — undetermined; consistent with $R0_{\text{cont}}$, R3 plausible. The debiasing mechanism is
 2629 reward rescaling via $|y|$ -normalization, exactly the cardinal-scale failure mode of Corollary A.4.
- 2630 • Lu et al. [51] (SamPO) — undetermined; consistent with $R0_{\text{cont}}$ or R1. Algorithmic length dependence in DPO
 2631 addressed analytically; empirical validation stays at the audit distribution and does not isolate Δ_j from length
 2632 reduction.
- 2633 • Fu et al. [24] (PAR) — undetermined on bias; arguably R0/R3 on the training-stability target it actually addresses.
 2634 Single Gemma-2B lineage, no off-target measurement.
- 2635 • Li et al. [45] (LMPO) — undetermined; consistent with $R0_{\text{cont}}$. The reported LC win-rate gain is on the same
 2636 metric the loss is built around, which is Goodhart-style and does not certify movement along the structural axis.
- 2637 • Shen et al. [70] (Loose Lips) — undetermined; consistent with $R0_{\text{cont}}$. Validation distribution coincides with the
 2638 RM training distribution; only pooled length correlation reported.
- 2639 • Chen et al. [8] (ODIN) — undetermined; $R0_{\text{cont}}$ vs R1 unresolved. Two-head disentanglement zeros pooled length
 2640 correlation; Lemma 3.6 says this does not transfer to μ_{π^*} .
- 2641 • Wang et al. [81] (HelpSteer) — undetermined; consistent with R3. Multi-attribute data construction with ordinal
 2642 MT-Bench evaluation; no operator-level mitigation claim the framework can adjudicate.
- 2643 • Bu et al. [3] (ALBM) — direct $R2_{\epsilon}$ evidence on the targeted axis (length-asymmetric subset accuracy regression
 2644 in Table 1) plus undetermined $R0_{\text{cont}}$ overall. No rescalability sweep to disambiguate scale overshoot from target
 2645 misspecification.
 2646
 2647
 2648
 2649
 2650
 2651
 2652

- 2653 • Huang et al. [30] (post-hoc LOESS calibration) – direct R_{2_ϵ} evidence on the targeted axis with R1 plausible.
2654 Section B.1 shows pooled correlation $0.316 \rightarrow 0.037$ but within-prompt sign flip on three of four SOTA RMs
2655 and $\Delta_J < 0$ on two AlpacaEval cells.
2656
- 2657 • Zhao et al. [91] (FiMi-RM) – direct R2 or R3 evidence (overall accuracy $70.8 \rightarrow 69$ reported in the paper itself),
2658 with $R_{0_{\text{cont}}}$ undisambiguated. Three-phase length-reward pattern identified but no rescalability sweep.
2659
- 2660 • Wang et al. [80] (causal rewards) – undetermined; consistent with $R_{0_{\text{cont}}}$. Synthetic identifiability sound;
2661 real-world evidence is correlation with GPT-4 and humans, not isolated Δ_J .
2662
- 2663 • Ng et al. [55] (debiasing with guarantees) – undetermined; consistent with $R_{0_{\text{cont}}}$. Identifiability theorem is
2664 conditional on the latent-variable DAG, which preference data does not pin down [72]; empirical validation is
2665 small/synthetic.
2666
- 2667 • Fein et al. [21] (mechanistic reward shaping) – closest to a positive R0 signal; $R_0/R_{0_{\text{cont}}}$ undisambiguated.
2668 Section B.1 shows $\Delta_J > 0$ on four of five RMs at within-prompt μ_{diag} , but off-target Δ_J across the rest of Φ_{sp} is
2669 not measured.
2670
- 2671 • Li et al. [47] (DIR, info-theoretic) – undetermined; consistent with $R_{0_{\text{cont}}}$. Information-bottleneck framing han-
2672 dles non-linearity in principle; empirical evidence conflates bias mitigation with downstream task improvement.
2673
- 2674 • Eisenstein et al. [19] (RM ensembles) – direct $R_{0_{\text{cont}}}/R_1$ evidence. The paper itself documents shared failure
2675 modes where ensemble members agree on length doubling and copying.
2676
- 2677 • Liu et al. [50] (worst-case correlated proxies) – undetermined; consistent with $R_{0_{\text{cont}}}$ or R3. Improves a worst-
2678 case proxy bound; the worst-case proxy is not the true reward, and substitution onto features outside the
2679 correlation set is undetected.
2680
- 2681 • Feng et al. [22] (C2PO) – direct R2 evidence. Reduces targeted bias with documented general-capability
2682 degradation, the named R2 pattern.
2683
- 2684 • Ye et al. [86] (PRISM) – undetermined; consistent with $R_{0_{\text{cont}}}$ with R1 plausible on uncovered axes. Near-zero
2685 correlation on the three measured shortcuts only; no off-panel measurement.
2686
- 2687 • Kim et al. [36] (CDA via causal lens) – undetermined; consistent with $R_{0_{\text{cont}}}$ on RewardBench-1, R3 on
2688 RewardBench-2. Length-controlled gains real on the easier benchmark; small gains on the harder one.
2689
- 2690 • Cai et al. [4] (Rc-BT) – undetermined; consistent with $R_{0_{\text{cont}}}$. Strongest gesture toward prescription 1 in the
2691 surveyed set (PPO included in Table 11), but no off-target Δ_J across Φ_{sp} .
2692
- 2693 • Srivastava et al. [75] (CHROME) – undetermined; consistent with $R_{0_{\text{cont}}}$. Best multi-prescription coverage
2694 in the surveyed set (multi-axis Table 7, BoN Figure 5, reWordBench transformations Tables 12–13), but the
2695 LLM-oracle counterfactuals inherit R4 sensitivity that is not measured.
2696
- 2697 • Song et al. [74] (SAE causal adjustment) – undetermined; consistent with $R_0/R_{0_{\text{cont}}}$. SAE features give partial
2698 Φ_{sp} , but the panel is incidental rather than designer-specified, and there is no fixed- π_{ref} evaluation.
2699
- 2700 • Chen et al. [7] (OPRM with RgFT) – undetermined; consistent with $R_{0_{\text{cont}}}$. Calibration framing (ECE, Brier in
2701 Table 5) is the closest engagement with the cardinal/ordinal distinction in the surveyed set, but the protocol
2702 does not exploit it for substitution detection.
2703
2704