

CAUSALLY-ENHANCED REINFORCEMENT POLICY OPTIMIZATION OF LARGE LANGUAGE MODELS

Anonymous authors

Paper under double-blind review

ABSTRACT

Large language models (LLMs) trained with reinforcement objectives often achieve superficially correct answers via shortcut strategies, pairing correct outputs with spurious or unfaithful reasoning and degrading under small causal perturbations. We introduce *Causally-Enhanced Policy Optimization* (CE-PO), a drop-in reward-shaping framework that augments policy optimization with a differentiable proxy for causal coherence along the generation pathway from prompt (Z) to rationale (X) to answer (Y). CE-PO estimates model-internal influence with Jacobian-based sensitivities, counterfactually hardens these signals to suppress nuisance cues, and fuses the resulting coherence score with task-accuracy feedback via a Minkowski (power-mean) combiner, exposing a single tunable between accuracy and coherence trade-off. The unified reward integrates with PPO/GRPO without architectural changes. Across reasoning benchmarks and causal stress tests, CE-PO reduces reward hacking and unfaithful chain-of-thought while improving robustness to correlation-causation flips and light counterfactual edits, all at near-parity accuracy. Experimental results across 4 datasets show that CE-PO improves accuracy over baselines by 5.49 % on average (up to 9.58 %), while improving robustness to correlation-causation flips and light counterfactual edits.

1 INTRODUCTION

Outcome-centric reinforcement training of large language models (LLMs)—whether via standard policy optimization or preference-based tuning (Christiano et al., 2017; Ouyang et al., 2022; Schulman et al., 2017; Shao et al., 2024) optimizes *what* answer is produced but places little pressure on *how* that answer is obtained (Sim et al., 2025). As a result, models routinely learn to “game” reward signals (e.g., by exploiting length, format, or lexical overlap) and to pair correct answers with unfaithful or spurious reasoning traces (Wen et al., 2024; Arjovsky et al., 2020). This failure mode is not merely theoretical: without intervention, we observe models can game both accuracy proxies and *gradient-based sensitivity* surrogates—e.g., by guessing final answer or maximizing output length as in section 3. This mismatch between scoring for outcomes and expecting coherent reasoning harms robustness, causing severe and threatening hallucination.

We argue that robust reasoning requires not only accurate outputs but also *causal coherence* along the generation pathway. By *coherence* we mean that (i) the prompt genuinely informs the rationale, (ii) the prompt along with rationale in turn genuinely informs the final answer, and (iii) small, causally relevant edits to the prompt or rationale induce correspondingly appropriate changes downstream, whereas irrelevant edits do not. Concretely, let $(Z \rightarrow X \rightarrow Y)$ denote the intended chain from prompt (Z), to rationale (X), to answer (Y) (Bao et al., 2024a). If the model’s internal influence respects this order, chain-of-thought (CoT) (Wei et al., 2022) traces are more likely to be right for the right reasons (Ouyang et al., 2022; Shao et al., 2024) and to generalize under distribution shift. However, current reinforcement objectives offer no internal signal that rewards such coherence, as previous literature on causal reward in LLM typically focus on causally indifferent to external factors instead of reasoning coherence (Wang et al., 2025). Related efforts to promote faithful reasoning typically rely on external supervision or decoding-time control: rationale supervision and “right-for-the-right-reasons” regularizers require human explanations and usually target only input and generation; concept-bottleneck (Koh and Liang, 2017) and invariance/IRM approaches (Kamath et al., 2021) enforce intermediate structure but do not address generative $\text{prompt} \rightarrow \text{rationale} \rightarrow \text{answer}$ flows; CoT verification and constrained decoding act *post hoc* (Chi et al., 2025; Heimersheim et al.,

2024) rather than shaping training. In contrast, we seek a *differentiable, annotation-free, training-time* signal that explicitly rewards influence coherence across $Z \rightarrow X \rightarrow Y$ and integrates into standard policy optimization.

Translating the above goal into a trainable signal raises three coupled issues. *Measurability*—existing faithfulness tools (saliency, raw input-gradient attributions) are fragile and can fail basic sanity checks or depend weakly on the learned model, limiting their usefulness as training signals (Adebayo et al., 2018); gradient regularization that makes models “right for the right reasons” typically relies on human-provided rationales and targets only input \rightarrow output links in discriminative settings, not the generative prompt \rightarrow rationale \rightarrow answer pathway (Ross et al., 2017), while causal-reasoning benchmarks show LLMs still conflate correlation and causation (Jin et al., 2024). This motivates a *differentiable, model-internal* proxy for directed influence along $Z \rightarrow X$, $X \rightarrow Y$, and $Z \rightarrow Y$ that is compatible with policy-gradient optimization. *Robustness to spurious sensitivity*—naïve gradient signals are easily inflated by non-semantic factors (token frequency, position, format), encouraging shortcut learning (Geirhos et al., 2020) and *reward hacking* under outcome-only RL objectives (Amodei et al., 2016); moreover, CoT can appear compelling yet remain unfaithful and brittle under small perturbations (Wei et al., 2022; Bao et al., 2024b; Chi et al., 2024). Hence we require *counterfactual hardening* that explicitly breaks semantic links and filters out nuisance directions rather than trusting raw sensitivities. *Objective balance*—robustness-oriented methods (e.g., invariance penalties) often trade off in-distribution accuracy and can underperform ERM when mis-specified (Arjovsky et al., 2020; Rosenfeld et al., 2020; Kamath et al., 2021), whereas optimizing only task reward preserves shortcut exploitation (Amodei et al., 2016). We therefore need a *tunable* mechanism to balance accuracy against coherence during optimization, avoiding both over-regularization and reward hacking.

Our approach. We introduce Causally-Enhanced Policy Optimization(CE-PO), a drop-in reward-shaping scheme that adds a *Jacobian-based causal-coherence signal* to standard policy optimization and fuses it with task accuracy via a *Minkowski (power-mean) combiner*. Concretely, we compute model-internal Jacobian influence signals for the *prompt \rightarrow rationale* and *rationale \rightarrow answer* links; harden them with a simple counterfactual procedure (removing the resulting nuisance directions by reshuffling) to reduce sensitivity to superficial cues; normalize and aggregate into a single *coherence score*; and combine this score with accuracy using a power-mean with tunable weights/exponent, yielding a single reward that trades off coherence and accuracy. The unified reward is plugged into PPO/GRPO without architectural changes. In summary, our contributions are threefold:

- We propose a differentiable, model-internal *DCE-proxy* (influence-coherence) reward along $Z \rightarrow X \rightarrow Y$, coupled with *counterfactual residualization* (see subsection 2.1) to suppress shortcut sensitivities and *reward hacking*.
- We introduce a *Minkowski (power-mean) combiner* (detailed in subsection 2.2) that exposes a tunable trade-off between task accuracy and causal-coherence signals, serving as a drop-in objective for PPO/GRPO with efficient reward computation.
- CE-PO demonstrates consistent performance improvements, supported by comprehensive ablation analyses and competitive out-of-distribution generalization. In addition, our study highlights notable generation patterns and offers insights into model behavior.

2 CAUSAL ENHANCED POLICY OPTIMIZATION

Overview. In this section, we give implementation details of *CE-PO* and formalize it. Given a rollout with prompt Z , rationale X , and answer Y , as shown in Figure 1. The method unfolds in three steps that map one-to-one to the subsections below. (i) **Raw influence signals** : we compute *Base Jacobians*, i.e., local sensitivities of downstream span likelihoods with respect to upstream token embeddings, as differentiable, model-internal proxies for $Z \rightarrow X$, $X \rightarrow Y$, and for $Z \rightarrow Y$ removing mediator X . (ii) **Counterfactual hardening** : Deploying raw Jacobian scores as rewards risks *reward hacking*, where spurious factors (e.g., length) dominate the signal Figure 2. To suppress such shortcut and formatting sensitivities, we break semantic links (reshuffle/mismatch) and remove the induced nuisance directions from the raw signals. (iii) **Reward normalization and fusion** We normalize Jacobian responses into a stable scalar *coherence score* and then fuse it with task accuracy using a Minkowski (power-mean) combiner, yielding a unified reward that exposes a tunable accuracy-coherence trade-off and is optimized with PPO/GRPO. This section defines each component, provides complexity notes, and clarifies how they compose into the training loop.

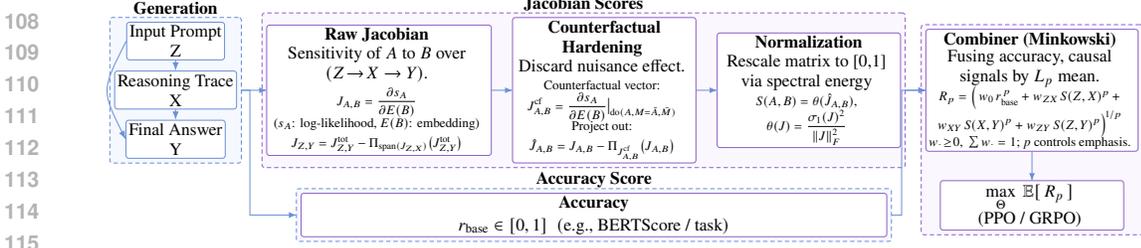


Figure 1: **CE-PO pipeline (left \rightarrow right)**. *Generation*: $Z \rightarrow X \rightarrow Y$ (X/Y split by a special token; see Fig. 14). *Base Jacobians*: local sensitivities for $Z \rightarrow X$, $X \rightarrow Y$, $Z \rightarrow Y$. *Counterfactual hardening*: build source-side counterfactuals and project their subspace from the base signals to obtain hardened residuals. *Normalization*: use spectral energy to convert Jacobian matrix into $[0,1]$ scalar. *Fusion*: combine $S(A, B)$ signals with accuracy r_{base} via a Minkowski combiner; optimize the unified reward with PPO/GRPO.

2.1 JACOBIAN SIGNALS, AND CAUSALLY-ENHANCED REWARD

We consider a differentiable language model Θ and our goal is, as previously described, to improve the accuracy on Y (i.e., r_{base}) yet simultaneously to ensure the alignment of influence with this order: Z should shape X , and X, Z should shape Y . Our proposed method is illustrated in Figure 1, providing a causally grounded reward signal R_p by combining r_{base} and Jacobian-based scores measuring causal effects along $Z \rightarrow X \rightarrow Y$.

Base Jacobians. While Y -accuracy rewards *what* answer is produced, it does not constrain *how* the answer is obtained. To shape the reasoning pathway, we introduce a differentiable, model-internal *Direct Causal Effect (DCE) proxy* (Pearl, 2009; VanderWeele, 2011) for coherence along $Z \rightarrow X \rightarrow Y$: small and causal changes in upstream tokens should induce appropriate changes downstream. Concretely, let s_X and s_Y denote the log-likelihoods of the rationale span X and the answer span Y posterior to Z and Z, X as in Equation 1, and let $E(A) \in \mathbb{R}^{T_A \times d}$ be the token-embedding matrix for a text segment A (length T_A , hidden size d) and $a_{<t}$ denotes the prefix strictly before index t .

$$s_X = \sum_{t \in X} \log p_{\Theta}(x_t | Z, x_{<t}), \quad s_Y = \sum_{t \in Y} \log p_{\Theta}(y_t | Z, X, y_{<t}) \quad (1)$$

We measure the influence from A to the score of B by the *block Jacobian* $J_{AB} \triangleq \partial s_B / \partial E(A)$. For the intended chain $Z \rightarrow X \rightarrow Y$, the *Base Jacobians* (raw, unhardened sensitivities) are

$$J_{ZX} = \frac{\partial s_X}{\partial E(Z)} \in \mathbb{R}^{T_Z \times d}, \quad J_{XY} = \frac{\partial s_Y}{\partial E(X)} \in \mathbb{R}^{T_X \times d}, \quad J_{ZY}^{total} = \frac{\partial s_Y}{\partial E(Z)} \in \mathbb{R}^{T_Z \times d}. \quad (2)$$

In addition, to approximate the Direct Causal Effect (DCE) from Z to Y , we need to isolate X effect out of J_{ZY}^{total} . We remove the mediated path via projection: forming a via- X template from J_{ZX} ($Z \rightarrow X$) and the principal direction of J_{XY} ($X \rightarrow Y$), then subtracting its projection from the total effect to obtain the direct effect:

$$J_{ZY} = J_{ZY}^{total} - \Pi_{J_{ZX} \odot \bar{u}}(J_{ZY}^{total}), \quad (3)$$

where $\bar{u} = \frac{1}{T_X} \sum_{t \in X} \frac{J_{XY}[t,:]}{\|J_{XY}[t,:]\|_2}$, and $\Pi_A(B) = \frac{\langle B, A \rangle_F}{\langle A, A \rangle_F} A$, with $\langle B, A \rangle_F$ indicating the Frobenius inner product between A and B , and \odot be Hadamard product.

We use the Jacobian as a local, direct causal effect *proxy*: for any upstream-downstream pair $(A, B) \in \{(Z, X), (X, Y), (Z, Y)\}$, with mediators held at their observed values, $J_{AB} = \partial s_B / \partial E(A)$ quantifies how an infinitesimal perturbation to A 's embeddings shifts the downstream log-likelihood s_B —i.e., which directions in A immediately move B , and by how much. Crucially, J_{AB} is not equal to a causal DCE in the mediation sense; it is the *first-order Taylor approximation* around the observed point (Pearl, 2001). They are already informative enough and avoid the cost and instability of higher-order (Hessian/influence-function) estimators in deep/LLM settings (Pearlmutter, 1994; Martens, 2010; Koh and Liang, 2017; Li et al., 2024).

We use J_{AB} to serve as the starting point for our *coherence* signal. In the next subsections we apply *counterfactual hardening* to remove nuisance directions and then normalize/aggregate the residual influence into a scalar coherence score that is fused with task accuracy.

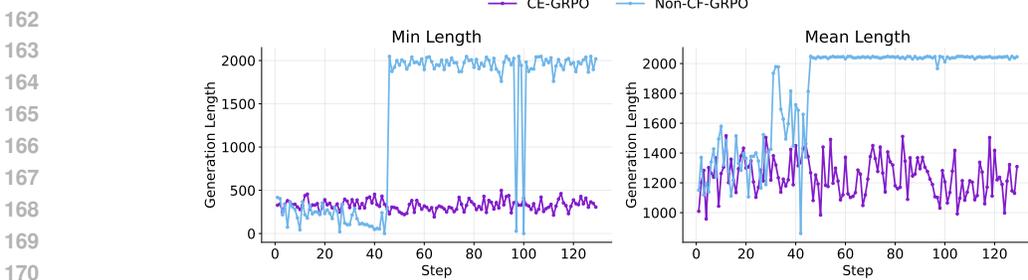


Figure 2: Training trajectories comparison: **Non-CF-GRPO** (GRPO with non-CF award) vs **CE-GRPO** (GRPO with our CE award) on Qwen-3-4B-Thinking with `max_tokens=2048`, in terms of minimum (left) and mean (right) generation length. Non-CF-GRPO curves plateau as generations approach the token limit, indicating length-driven reward hacking, while CE-GRPO doesn’t.

Counterfactual-hardened Jacobian score. Jacobian-based rewards, when computed directly on the base input, may conflate causal content with nuisance factors such as style, formatting and token-frequency statistics. As a result, innocuous paraphrases or vacuous elongation can inflate the gradient norm and concentrate its top singular direction, enabling *length-driven reward hacking* up to the token limit (Gao et al., 2023; Rafailov et al., 2024). Empirically, as shown in Figure 2 (blue), the length hacking existing in non-counterfactual Jacobian rewards (Non-CF-GRPO conflates).

To suppress shortcut-driven sensitivities in the *Base Jacobians*, we harden them with pair-specific counterfactuals. For each link $(A, B) \in \mathcal{L} := \{(Z, X), (X, Y), (Z, Y)\}$, let $M(A, B)$ denote mediators on $A \rightarrow B$ paths (empty except $M(Z, Y) = \{X\}$). We break the semantic alignment between A (and M) and B by permuting the tokens of A (and of M when present), yielding (\bar{A}, \bar{M}) that preserves surface statistics (length/frequency) while being approximately independent of B . Let s_B be the log-likelihood of span B and $E(A) \in \mathbb{R}^{T_A \times d}$ the embedding matrix of A . **Token reshuffling is used to construct nuisance dimensions because we target to destroy semantics while perfectly preserving low level nuisance cues, comparing to alternative methods like deleting or replacing, etc.** We compute the *counterfactual Jacobian*

$$J_{AB}^{\text{cf}} = \left. \frac{\partial s_B}{\partial E(A)} \right|_{(A, M) \rightarrow (\bar{A}, \bar{M})}. \quad (4)$$

We construct K *source-side* counterfactuals and calculate their Jacobians by Equation 4, which are then averaged to have \bar{J}_{AB}^{cf} (we set $K=4$, larger values generally improve results but incur higher computational cost).

Then, we extract the principal counterfactual subspace via the top- k left singular vectors U_{AB} of J_{AB}^{cf} , define the projector $\Pi_{AB}(J) = U_{AB}U_{AB}^\top J$, following (Ravfogel et al., 2020). Taking the J_{AB} in Equation 2 and Equation 3, we obtain the deconfounded (direct) block by removing these directions:

$$\hat{J}_{AB} = J_{AB} - \Pi_{AB}(J_{AB}). \quad (5)$$

Permuted signals expose a spurious nuisance subspace; projecting it out yields a deconfounded, causally aligned signal that resists surface correlations and reward hacking (Ravfogel et al., 2020).

As shown in Figure 2 (purple), the Non-CF-GRPO baseline drives the mean generation length to the token cap, revealing a length-based reward-hacking loop. With our counterfactual reshuffling and projection (CE-GRPO), the curve stays well below the cap and does not plateau, while reward improves. The mechanism is simple: source-side counterfactuals preserve surface statistics (length/frequency) while scrambling semantics (Ravfogel et al., 2020), so the counterfactual Jacobian spans a nuisance subspace; orthogonally projecting the base Jacobian off this subspace—akin to Neyman orthogonalization (Oprescu et al., 2019)—decouples the reward from length/format sensitivities and correlations (Chernozhukov et al., 2018). For the $Z \rightarrow Y$ link, we additionally subtract the via- X component before projection (a first-order direct-effect correction), further reducing spurious incentives.

Reward normalization and fusion Finally, we *scalarize* the block to a score via the scale-free spectral energy share $S(A, B) = \phi(\hat{J}_{AB})$, where $\phi(J) = \sigma_1(J)^2 / \|J\|_F^2 \in [0, 1]$ to squeeze and standardize the process (Vershynin, 2018). Concretely,

$$S(Z, X) = \phi(\hat{J}_{ZX}), \quad S(X, Y) = \phi(\hat{J}_{XY}), \quad S(Z, Y) = \phi(\hat{J}_{ZY}). \quad (6)$$

Let $r_{\text{base}} \in [0, 1]$ be a task accuracy score on Y (e.g., semantic correctness, measured by BERTScore (Zhang et al., 2019) in our setting). We define the new causally grounded reward by integrating accuracy and the three counterfactual-enhanced signals via a weighted power mean (Minkowski combiner) (Bullen, 2013),

$$R_p = \left(w_0 r_{\text{base}}^p + w_1 S(Z, X)^p + w_2 S(X, Y)^p + w_3 S(Z, Y)^p \right)^{1/p}, \quad (7)$$

with $p \in \mathbb{R}$, weights $w_i \geq 0$, and $\sum_{i=0}^3 w_i = 1$. As $p \rightarrow 0$, R_p approaches a weighted geometric mean (minimal optimization); larger p emphasizes the largest channel, which respectively matches deconfounded Jacobian causal scores and BERTScore as in subsection 2.2. Using this R_p in any common policy-gradient training of LLM (e.g., PPO, GRPO), the goal is to maximize $\mathbb{E}[R_p]$.

Note of r_{base} : We set r_{base} by BERTScore (i.e., the cosine similarity between Y and ground truth in BERT embedding space), as it is a continuous measure between 0 and 1 for accuracy, aligning with the scale of the causal term. In contrast, 0/1 judgments are often sensitive to judging prompts and are less compatible; mixing such binary signals with the causal term can induce oscillatory updates. Ablation studies on replacing BERTScore with 0/1 rewards are provided in section 3. In addition, we also provide TRPO-form (Schulman et al., 2015) lower bound for CE-PO with details in Appendix B and the discussion for training efficiency can be found in Appendix C.

Training objective Let $R_p(Z, X, Y) \in [0, 1]$ denote the unified reward from Eq. (7). We treat R_p as a scalar terminal reward for each rollout (Z, X, Y) : $r_T = R_p$, $r_t = 0$ for $t < T$, and **do not backpropagate through** R_p (standard policy-gradient; the Jacobian signals are used only to compute R_p).

CE-PPO. With behavior policy π_{old} , importance ratio $w_t(\theta) = \pi_\theta(a_t|s_t)/\pi_{\text{old}}(a_t|s_t)$, and GAE advantages \hat{A}_t computed from the sparse reward, the clipped objective is

$$\mathcal{J}_{\text{CE-PPO}}(\theta) = \mathbb{E} \left[\min(w_t(\theta)\hat{A}_t, \text{clip}(w_t(\theta), 1-\epsilon, 1+\epsilon)\hat{A}_t) - \beta \text{KL}(\pi_\theta \parallel \pi_{\text{ref}}) \right]. \quad (8)$$

We optionally standardize R_p per batch to stabilize advantages.

CE-GRPO. For each prompt we sample K candidates, compute $R_{p,k}$, and form group-relative advantages $\tilde{A}_k = (R_{p,k} - \mu_{\text{grp}})/(\sigma_{\text{grp}} + \epsilon)$, with $\mu_{\text{grp}}, \sigma_{\text{grp}}$ the group mean/std. The GRPO-style objective is

$$\mathcal{J}_{\text{CE-GRPO}}(\theta) = \mathbb{E} \left[\frac{1}{K} \sum_{k=1}^K \min(w_k(\theta)\tilde{A}_k, \text{clip}(w_k(\theta), 1-\epsilon, 1+\epsilon)\tilde{A}_k) - \beta \text{KL}(\pi_\theta \parallel \pi_{\text{ref}}) \right]. \quad (9)$$

2.2 REWARD SIGNALS TRADE-OFF

Our RL objective (maximizing $\mathbb{E}[R_p]$) blends two signals, semantic accuracy (e.g., BERTScore) and causal coherence (counterfactual-enhanced Jacobian). These signals misalign at times, mirroring RLHF’s trade-off between surface correctness and shortcut avoidance (Kalai et al., 2025; Lightman et al., 2023). By defining R_p via the Minkowski norm, we enable a dynamically tunable balance between the two signals.

Figure 3 illustrates how the Minkowski p -norm interpolation provides a smooth mechanism to mediate the trade-off between semantic accuracy and causal coherence. As $p \rightarrow +\infty$, the reward collapses to the larger signal, typically BERTScore, which is unprojected and numerically higher. This drives accuracy-focused optimization (as observed in all models). Conversely, as $p \rightarrow -\infty$, the reward collapses to the smaller signal (e.g., when $p = 0$ the interpolation equals geometric mean), which is usually the counterfactual Jacobian, lowered by residualization. This, in turn, emphasizes causal coherence. This multi-objective behavior parallels prior work on balancing conflicting RL rewards (Won et al., 2019; Roijers et al., 2013). To illustrate how CE-PO combine both signals can thus enhance LLM reasoning capability to be both causally capable yet not template-following, we provide an LLM generation example in Figure D.

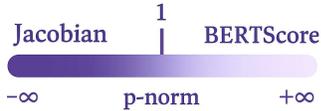


Figure 3: Tunable balance enabled by Minkowski p -norm interpolation in R_p .

3 EXPERIMENTS

In this section, we show performance of CE-PO on top of PPO and GRPO (noted as **CE-PPO** or **CE-GRPO**) against simple baselines, using BERTScore as the reward (denoted as PPO-BERTScore or GRPO-BERTScore) and vanilla LLMs, across representative LLMs and standard evaluation datasets.

Experiment Setting. Training is performed on a single node with four A100 or H100 GPUs. Unless stated otherwise, we use a rollout group size of 6 per prompt, KL penalty of $\lambda_{KL} = 0.001$ to stabilize optimization, and cap decoder outputs at 2048 tokens and input within 512 tokens. In determining R_p , extensive experiments indicate that a good configuration could be $p = 0.2$ and BERTScore, Jacobian scores weights, respectively, as $(w_0, w_1, w_2, w_3) = (0.7, 0.1, 0.1, 0.1)$, which is set as default for comparison unless otherwise specified in ablation studies. We set learning rate of 2×10^{-6} , batch size 32, and training for 10 epochs with temperature 0.6. Validation interval is set as 10 and checkpoints are saved per 100 steps. Experiments are based on Verl (Sheng et al., 2024) framework with mild modification for separate thread for reward calculation, and we also set warm-up step = 10 for critic.

To avoid random guesses and generate X and Y , all models are prompted to generate CoT reasoning (Wei et al., 2022). Evaluation is reported as accuracy, and the response correctness is determined via an LLM-as-Judge protocol (Zheng et al., 2023) using GPT-4o-mini (see prompt in Figure 15, which is designed to distinguish guess or contradictions especially for multi-choice questions).

We evaluate both standard instruction-tuned LLMs and “thinking” variants that emit explicit intermediate reasoning and instruction following. Although compute constraints preclude experiments on 70B-scale models, our evaluation covers compact to mid-sized systems that are widely used in research: Llama-3-8B-Instruct, Llama-3.2-3B-Instruct, Phi-3.5-mini-Instruct, Qwen3-4B-Thinking, and Qwen3-1.7B-Thinking (Meta AI, 2024b;a; Microsoft, 2025; Qwen Team, 2025b;a).¹

Training and Validation Dataset. We train and tune on four reasoning-heavy sets: *BBEHCausal*, the causal reasoning split from BIG-Bench Extra Hard dataset, including counterfactual deduction or cause inference questions; (Google DeepMind, 2025)); *CaseHOLD*, tasks centering on multiple-choice legal holdings given case citation context (CaseHOLD, 2022); *MATHHARD*, the level-5 subset targeting multi-step mathematical reasoning with multiple-choice questions (Hendrycks et al., 2021); and *IfQA*, an open-domain QA under counterfactual *if*-clause presuppositions requiring hypothetical reasoning (Yu et al., 2023). Note that all datasets are formatted as question answer pairs, examples of them are presented respectively in Table 6 and Table 7.

Testing Dataset. For same-field evaluation, we use *BBEHMATH* (challenging multistep arithmetic), where we filter shorter and more trivial questions and normalize alphabetic numerals to digits with unfolded calculation for comparability on this dataset (Google DeepMind, 2025). We further assess generalization with *CLadder*, covering association/intervention/counterfactual queries (Causal-NLP, 2023), *LegalBench* (case-understanding split) (Guha et al., 2023) analyzing the cause or reason of legal case, and *LogiQA*, a multiple-choice dataset emphasizing deductive logical reasoning sourcing from officer entrance exams (lucasmccabe, 2021).

Main Experiment Results. Under the training and validation dataset split detailed above, Figure 4 shows CE-PPO delivering smooth, near-monotonic validation gains across five backbones and four validation sets: rapid early improvement followed by mild saturation, which is indicative of stable and standard RL optimization plot rather than overfitting. In Table 1, across all five backbones, both causal variants outperform their non-causal baselines: CE-PPO/CE-GRPO improve the best non-CE baseline by 2.3 to 9.58 points. Comparing CE-GRPO with the vanilla baseline, we observe gains of 5.8–6.0 points on Llama-3-8B and Llama-3.2-3B; the largest improvement appears on Qwen-3-1.7B-Thinking (≈ 10 points), whereas Phi-3.5-mini-Instruct shows the most marginal uplift. These patterns suggest that architectural and scale differences modulate the benefits of CE-GRPO. CE-GRPO is *not* uniformly stronger than CE-PPO, while slightly higher on Qwen-1.7B/4B, Phi-3.5-mini, and Llama-3-8B, it lags behind CE-PPO on Llama-3.2-3B, indicating complementary strengths rather than strict dominance.

¹Note that all experiments use Instruct or Thinking models, as our outputs must be cleanly formatted into X and Y , which base models struggle to separate. Since Instruct models are typically harder to improve with RL, the observed gains further highlight the effectiveness of CE-PO.

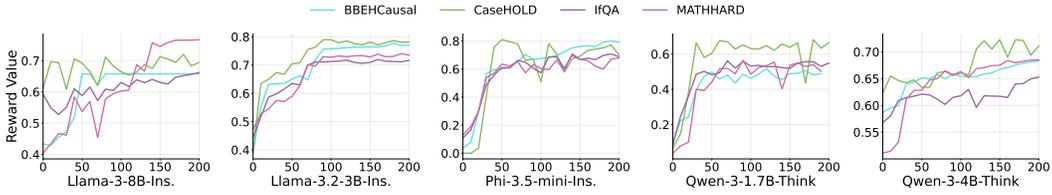


Figure 4: Reward curves of CE-PPO training on validation datasets across five models.

Table 1: **Main results across models, RL methods, and datasets.** We report LLM-as-Judge accuracy (%) under different RL variants: vanilla model, standard PPO/GRPO with BERTScore reward, and causal-enhanced (CE-PPO/CE-GRPO). We repeatedly generate 4 times of each trained model, and thus report mean and deviation below.

Model	RL Method	Dataset (%)				Average
		BBEHMATH	CLadder	LegalBench	LogiQA	
Qwen-3-1.7B-Thinking	Vanilla	1.36 ± 1.96	34.77 ± 0.68	57.80 ± 2.59	50.00 ± 0.45	35.98
	PPO-BERTScore	5.91 ± 0.00	37.27 ± 0.91	58.26 ± 1.95	52.50 ± 0.00	38.48
	GRPO-BERTScore	3.64 ± 1.29	35.00 ± 0.52	60.67 ± 0.94	57.04 ± 0.68	39.09
	CE-PPO	10.45 ± 0.52	49.32 ± 0.68	60.67 ± 1.05	58.05 ± 1.14	44.62
	CE-GRPO	10.00 ± 1.29	44.77 ± 0.00	70.64 ± 0.00	56.82 ± 0.45	45.56
Qwen-3-4B-Thinking	Vanilla	6.82 ± 0.45	42.95 ± 0.91	74.77 ± 1.95	80.45 ± 0.68	51.25
	PPO-BERTScore	5.00 ± 1.05	45.00 ± 0.00	73.64 ± 0.47	80.00 ± 0.45	50.91
	GRPO-BERTScore	6.82 ± 0.45	40.68 ± 0.68	76.61 ± 0.65	83.86 ± 0.91	51.99
	CE-PPO	9.55 ± 0.52	45.45 ± 0.68	79.34 ± 1.05	79.77 ± 0.91	53.53
	CE-GRPO	7.27 ± 0.00	49.09 ± 0.64	76.90 ± 0.00	84.09 ± 0.45	54.34
Phi-3.5-mini-Ins.	Vanilla	6.13 ± 0.52	40.23 ± 0.68	68.35 ± 1.95	54.09 ± 0.45	42.20
	PPO-BERTScore	6.59 ± 0.00	43.40 ± 0.87	66.97 ± 0.00	54.77 ± 0.00	42.93
	GRPO-BERTScore	5.45 ± 0.87	41.59 ± 0.64	68.81 ± 1.30	52.73 ± 0.45	42.15
	CE-PPO	8.41 ± 0.00	42.95 ± 0.91	70.18 ± 1.95	55.90 ± 0.68	44.36
	CE-GRPO	6.81 ± 0.64	45.45 ± 0.45	70.18 ± 0.65	57.95 ± 0.00	45.10
Llama-3.2-3B-Ins.	Vanilla	3.63 ± 0.91	30.00 ± 0.68	45.41 ± 0.47	41.59 ± 0.45	30.16
	PPO-BERTScore	6.59 ± 0.68	33.63 ± 0.64	45.87 ± 0.00	38.41 ± 0.45	31.12
	GRPO-BERTScore	6.50 ± 0.45	31.36 ± 0.91	42.73 ± 2.57	44.09 ± 0.00	31.17
	CE-PPO	10.67 ± 1.14	43.64 ± 0.68	55.05 ± 1.30	48.18 ± 0.91	39.38
	CE-GRPO	7.95 ± 0.91	35.23 ± 0.00	54.13 ± 2.59	47.27 ± 0.45	36.15
Llama-3-8B-Ins.	Vanilla	7.05 ± 0.45	39.54 ± 0.91	61.93 ± 1.95	65.68 ± 0.64	43.55
	PPO-BERTScore	5.27 ± 0.64	37.95 ± 0.00	61.93 ± 0.00	60.91 ± 0.45	41.52
	GRPO-BERTScore	10.91 ± 1.14	39.09 ± 0.91	57.80 ± 0.00	67.50 ± 0.00	43.83
	CE-PPO	14.86 ± 1.59	44.77 ± 0.45	67.35 ± 0.79	68.40 ± 0.91	48.84
	CE-GRPO	13.64 ± 0.91	45.00 ± 0.64	66.97 ± 0.00	71.82 ± 0.45	49.36

Out of Distribution(OOD) Evaluation. Trained exclusively on causal- and math-reasoning corpora, our CE-PO models are stress-tested for OOD generalization on TruthfulQA (Lin et al., 2021) (misconception resistant factuality), CodeMMLU (Manh et al., 2024) (programming knowledge), and SuperGPQA (Du et al., 2025) (expert level STEM QA), evaluating 200 uniformly random samples from each benchmark. As shown in Figure 5, they deliver consistent—but task-dependent improvements over baselines: effects are strongest on structure-heavy code reasoning, moderate on STEM-oriented SuperGPQA, and smaller yet steady on TruthfulQA, suggesting that causal regularization chiefly benefits multi-step reasoning while preserving overall robustness.

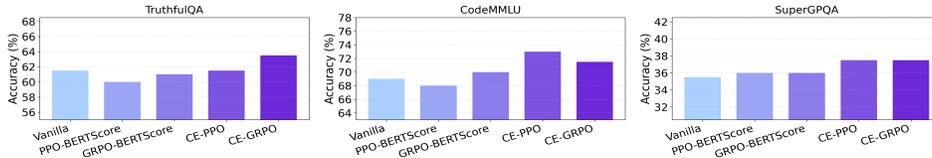


Figure 5: Qwen-3-4B-Thinking accuracy in OOD evaluation scenarios.

Ablation Studies. We test CE-PO’s composite objective under a fixed setup (Llama-3-8B-Instruct, GRPO, max_tokens=2048; datasets: BBEHMath, CLadder, LegalBench, LogiQA). The base reward is r_{base} (BERTScore); coherence uses hardened link scores $S(Z, X)$, $S(X, Y)$, $S(Z, Y)$. We ablate

Table 2: Ablation Studies on Llama-3-8B-Ins.

Ablation	BBEHMATH	CLadder	LegalBench	LogiQA	Average
$p = 10$	6.36	39.09	56.90	67.27	42.91
$p = -2$	12.73	43.63	65.52	70.00	47.97
$p = 1$	10.00	50.00	60.34	64.55	46.22
$w=[0, 1/3, 1/3, 1/3]$	10.91	44.54	61.93	68.18	46.89
$w=[4/5, 1/10, 1/10, 0]$	5.45	23.63	40.37	41.82	27.82
$w=[4/5, 1/10, 0, 1/10]$	3.63	30.00	45.87	49.09	32.15
$w=[4/5, 0, 1/10, 1/10]$	5.45	31.82	43.12	40.90	30.32
$w=[1/2, 1/6, 1/6, 1/6]$	8.18	41.82	67.24	66.36	45.90
Perturbed Reward	7.27	40.00	65.52	69.09	45.97
Judge Reward (binary)	6.36	39.09	60.34	66.36	43.04
CE-GRPO	13.64	47.27	66.97	72.73	50.65

the Minkowski exponent p and weights w , drop individual links, add Gaussian noise to r_{base} , and swap BERTScore for an LLM-as-Judge. Ablation results in Table 2 yield:

Insight 1: Reward hacking would emerge once we remove any of $S(Z, X)$, $S(Z, Y)$ and $S(X, Y)$. Models under such flawed RL can’t even excel the original models as observed in row 5 to 7 in Table 2. The likely cause is an incomplete constraint, i.e., the causal loop $Z \rightarrow X \rightarrow Y$ is not closed, so the optimization ceases to be strictly benign. With a detailed supporting plot provided in Figure 11 of the reward signal and generation length, we can observe that initial removal breaks the prior balance (reward dips), and subsequent updates discover a length shortcut (length conflating).

Insight 2: By selecting different p values, we can observe the performances of p as 10 identical to BERTScore only and p as -2 to Jacobian scores only (i.e. w as $[0, 1/3, 1/3, 1/3]$). The same case of p as 1, which reveals that tuning p can dynamically tune the optimization objective.

Insight 3: Under a perturbed reward (Gaussian noise $\sigma=0.05$), the model maintains satisfactory performance, indicating robustness of the training design. Using GPT - 4o - mini as an LLM-as-Judge (vs. BERTScore) yields performance similar to the original model, revealing what’s prescribed in the Note of r_{base} in subsection 2.1.

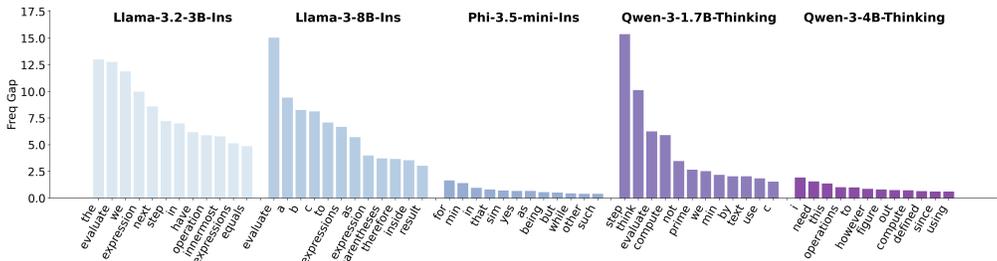


Figure 6: Frequency gap of generation per word where CE-GRPO exceeds BERTScore and Vanilla LLM (Top-12 per model).

Generation Pattern Analysis. Figure 6 shows that CE-GRPO does not induce uniform gains across all tokens but selectively amplifies reasoning or causal related words, suggesting that causal regularization guides the model toward tokens central to stepwise inference rather than surface templates. Complementarily, Figure 7 reveals the structural relation between metrics: $S(Z, X)$ is nearly independent of $S(X, Y)$ and $S(Z, Y)$, while the latter two are moderately correlated. This indicates that “evidence infusion” and “answer stability” provide orthogonal yet complementary signals. Their joint use, as in CE-GRPO, therefore supports coherence across the full $Z \rightarrow X \rightarrow Y$ pathway, rather than relying on any single proxy.

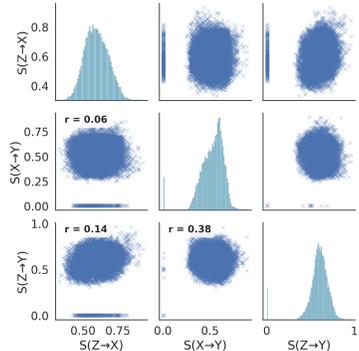


Figure 7: Pairwise scatter matrix of causal metrics on Qwen-3-4B-Thinking.

Turning to sequence-level behavior, Figure 9 reports the response lengths of CE-GRPO, GRPO-BERTScore, and the vanilla model on the training sets. The

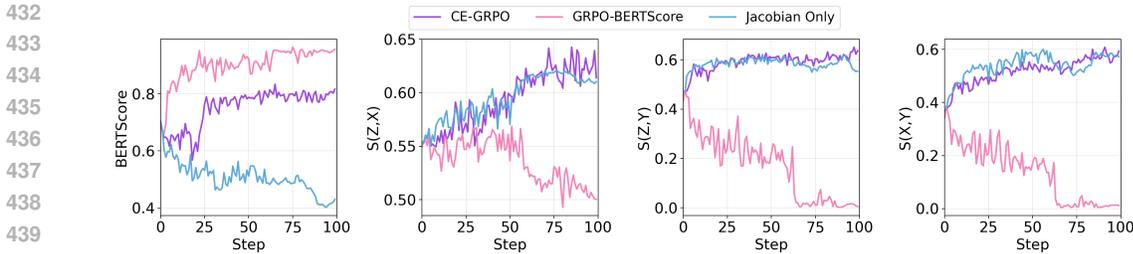


Figure 8: Training trajectories on Llama-3-8B-Instruct, reporting BERTScore, $S(Z, X)$, $S(Z, Y)$, $S(X, Y)$ (left to right). Comparison is made between CE-GRPO and optimization with only causal signal (Jacobian) or accuracy signal (BERTScore), supporting the claim of reward signal balance in subsection 2.2.

means are closely matched, indicating our optimization is largely *length-indifferent*. CE-GRPO achieves better task performance without increasing the token budget beyond the baseline, suggesting the efficiency of CE-GRPO.

Finally, the training reward trajectories in Figure 8 highlight three regimes: BERTScore-only optimization inflates BERTScore only while collapsing causal channels, Jacobian-only boosts sensitivities but remains unstable and harms task accuracy. This finding echoes prior work showing that proxy metric optimization (e.g., ROUGE) misaligns with human preferences (Stiennon et al., 2020), while self-rewarded training inflates output length, underscoring the tension between verifiable rewards and reward-gaming behaviors (Yuan et al., 2025). CE-GRPO achieves smooth improvements across $S(Z, X)$, $S(X, Y)$, and $S(Z, Y)$ while keeping BERTScore competitive. This pattern underscores CE-GRPO’s ability to sustain the causal chain from prompt to reasoning to answer, offering a practical safeguard against shortcut learning.

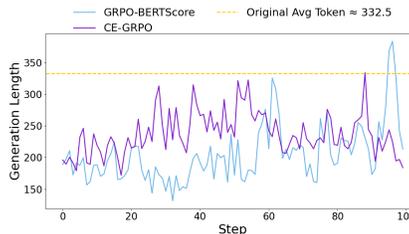


Figure 9: Response length comparison on Llama-3-8B-Ins

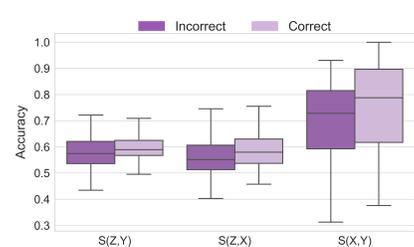


Figure 10: Boxplots of $S(Z, X)$, $S(Z, Y)$, and $S(X, Y)$ for LLM-as-Judge accuracy regarding their causal trajectory.

beyond surface outputs, and by statistically aligns with cognitive accuracy.

Jacobian Signals and Accuracy Figure 10 shows the distributions of Jacobian-based scores for samples annotated by correctness with LLM-as-Judge with strict judging criteria provided in Figure 15 conducted by GPT-4o. Annotations were collected on 1000 samples generated by Llama-3-8B-Ins. In total, 371 cases were labeled as consistent and 629 as inconsistent. To assess differences between the two groups, we conducted t-test with results reported in Table 5. We evaluate whether higher $S(A, B)$ scores align with factual correctness, as validity of these causal signals requires separation under strict accuracy criteria. As in Table 5, all three metrics show significantly higher means for the correct group ($p < 0.01$), confirming that $S(A, B)$ faithfully tracks reasoning–answer validity

Experiments on CE-PO’s Scalability for r_{base} Generally, the Minkowski combiner (Equation 7) benefits from a continuous $r_{base} \in [0, 1]$ that fuses smoothly with our dense causal signals (Appendix C). BERTScore offers such a graded, semantics-aware signal and avoids the oscillatory behavior of discrete 0/1 rewards (e.g., LLM-as-Judge in Table 2). To show that CE-PO fusion is not limited to BERTScore-style similarity rewards, we further replace r_{base} with a DeBERTa-v3-large reward model, representing classifier-style accuracy feedback. As summarized in Table 3, CE-PO maintains strong performance across both reward types using the same backbones and datasets as Table 1. Replacing BERTScore with pretrained reward model also doesn’t introduce very significant difference, suggesting these two signals are both reliable and feasible.

Model	Method	BBEHMATH	CLadder	LegalBench	LogiQA	Average
Qwen-3-4B-Thinking	GRPO-RM	6.61	44.44	77.33	81.03	00.00
	CE(RM)-GRPO	8.39	52.37	77.15	83.70	00.00
Phi-3.5-mini-Ins.	GRPO-RM	00.00	00.00	00.00	00.00	00.00
	CE(RM)-GRPO	00.00	00.00	00.00	00.00	00.00
Llama-3-8B-Ins.	GRPO-RM	8.39	39.54	63.00	69.58	00.00
	CE(RM)-GRPO	14.79	45.00	67.35	72.62	00.00

Table 3: Benchmarking comparison across RM (pretrained DeBERTa reward model) RL method and CE-PO employing RM as signal.

LLM as Judge Reliability To validate our choice of GPT-4o-mini as the LLM-as-a-judge, we conducted precisely this inter-rater reliability (IRR) analysis. We constructed an evaluation set by randomly sampling 200 generated examples. These examples were then assessed by three distinct sources: (1) human annotators, serving as the gold standard; (2) our selected judge, GPT-4o-mini; and (3) a stronger model, GPT-4o, for a comparative baseline. For the human annotation process, the 200 examples were distributed equally among four graduate student annotators, with each responsible for a unique batch of 50 examples. The results are summarized in Table 4.

Rater Pair Comparison	Agreement (%)	Cohen’s Kappa (κ)
GPT-4o-mini vs. Human	94.0%	0.88
GPT-4o vs. Human	96.5%	0.93
GPT-4o-mini vs. GPT-4o	98.5%	0.97

Table 4: Inter-rater reliability results of different annotators

4 CONCLUSION

We propose counterfactual-enhanced policy optimization that couples task accuracy with continuous causal signals, yielding smoother training and consistent gains across models and datasets, especially on Thinking backbones. The approach is a drop-in for PPO/GRPO and reduces reward-hacking behaviors, improving stability during training and robustness at convergence. Future work based on CE-PO could possibly work on scale interventions and automate the accuracy–causality balance; CE-PO opens new optimization avenues by incorporating richer causal signals (e.g., Hessian-based beyond simple Jacobian) and enabling adaptive schemes that tune this trade-off on the fly.

REFERENCES

- 540
541
542 Julius Adebayo, Justin Gilmer, Michael Muelly, Ian Goodfellow, Moritz Hardt, and Been Kim. Sanity
543 checks for saliency maps. *Advances in neural information processing systems*, 31, 2018.
- 544
545 Dario Amodei, Chris Olah, Jacob Steinhardt, Paul Christiano, John Schulman, and Dan Mané.
546 Concrete problems in ai safety. *arXiv preprint arXiv:1606.06565*, 2016.
- 547
548 Martin Arjovsky, Léon Bottou, Ishaan Gulrajani, and David Lopez-Paz. Invariant risk minimization.
549 In *International Conference on Learning Representations (ICLR)*, 2020.
- 550
551 Fan Bao et al. Non-causal reasoners in large language models. *arXiv preprint arXiv:2407.01489*,
552 2024a.
- 553
554 Guangsheng Bao, Hongbo Zhang, Linyi Yang, Cunxiang Wang, and Yue Zhang. Llms with chain-of-
555 thought are non-causal reasoners. *arXiv preprint arXiv:2402.16048*, 2024b.
- 556
557 Atilim Gunes Baydin, Barak A. Pearlmutter, Alexey A. Radul, and Jeffrey Mark Siskind. Automatic
558 differentiation in machine learning: A survey. *Journal of Machine Learning Research*, 18(153):
559 1–43, 2018. URL <https://jmlr.org/papers/v18/17-468.html>.
- 560
561 Peter S Bullen. *Handbook of means and their inequalities*, volume 560. Springer Science & Business
562 Media, 2013.
- 563
564 CaseHOLD. Casehold. Hugging Face dataset, 2022. URL [https://huggingface.co/
565 datasets/casehold/casehold](https://huggingface.co/datasets/casehold/casehold).
- 566
567 Causal-NLP. Cladder. Hugging Face dataset, 2023. URL [https://huggingface.co/
568 datasets/causal-nlp/CLadder](https://huggingface.co/datasets/causal-nlp/CLadder).
- 569
570 Victor Chernozhukov, Denis Chetverikov, Mert Demirer, Esther Duflo, Christian Hansen, Whit-
571 ney Newey, and James Robins. Double/debiased machine learning for treatment and structural
572 parameters, 2018.
- 573
574 Haoang Chi, He Li, Wenjing Yang, Feng Liu, Long Lan, Xiaoguang Ren, Tongliang Liu, and Bo Han.
575 Unveiling causal reasoning in large language models: Reality or mirage? In *Advances in Neural
576 Information Processing Systems (NeurIPS)*, 2024.
- 577
578 Zhenyu Chi et al. Unveiling causal reasoning in large language models via g^2 -reasoner. *arXiv
579 preprint arXiv:2503.13846*, 2025.
- 580
581 Paul F. Christiano, Jan Leike, Tom B. Brown, Miljan Martic, Shane Legg, and Dario Amodei. Deep
582 reinforcement learning from human preferences. In *Advances in Neural Information Processing
583 Systems (NeurIPS)*, 2017.
- 584
585 Xinrun Du, Yifan Yao, Kaijing Ma, Bingli Wang, Tianyu Zheng, King Zhu, Minghao Liu, Yiming
586 Liang, Xiaolong Jin, Zhenlin Wei, et al. Supergpqa: Scaling llm evaluation across 285 graduate
587 disciplines. *arXiv preprint arXiv:2502.14739*, 2025.
- 588
589 Tom Everitt, Victoria Krakovna, Laurent Orseau, Marcus Hutter, and Shane Legg. Reinforcement
590 learning with a corrupted reward channel. *arXiv preprint arXiv:1705.08417*, 2017.
- 591
592 Leo Gao, John Schulman, and Jacob Hilton. Scaling laws for reward model overoptimization. In
593 *Proceedings of the 40th International Conference on Machine Learning, ICML*. PMLR, 2023.
URL <https://proceedings.mlr.press/v202/gao23h/gao23h.pdf>.
- Robert Geirhos, Jörn-Henrik Jacobsen, Claudio Michaelis, Richard Zemel, Wieland Brendel, Matthias
Bethge, and Felix A Wichmann. Shortcut learning in deep neural networks. *Nature Machine
Intelligence*, 2(11):665–673, 2020.
- Gene H. Golub and Charles F. Van Loan. *Matrix Computations*. Johns Hopkins University Press, 4 edi-
tion, 2013. URL [https://math.ecnu.edu.cn/~jypan/Teaching/books/2013%
20Matrix%20Computations%204th.pdf](https://math.ecnu.edu.cn/~jypan/Teaching/books/2013%20Matrix%20Computations%204th.pdf).

- 594 Google DeepMind. Big-bench extra hard (bbeh). GitHub repository, 2025. URL [https://](https://github.com/google-deepmind/bbeh)
595 github.com/google-deepmind/bbeh.
596
- 597 Neel Guha, Julian Nyarko, Daniel E. Ho, Christopher Ré, Adam Chilton, Aditya Narayana, Alex
598 Chohlas-Wood, Austin Peters, Brandon Waldon, Daniel N. Rockmore, Diego Zambrano, Dmitry
599 Talisman, Enam Hoque, Faiz Surani, Frank Fagan, Galit Sarfaty, Gregory M. Dickinson, Haggai
600 Porat, Jason Hegland, Jessica Wu, Joe Nudell, Joel Niklaus, John Nay, Jonathan H. Choi, Kevin
601 Tobia, Margaret Hagan, Megan Ma, Michael Livermore, Nikon Rasumov-Rahe, Nils Holzenberger,
602 Noam Kolt, Peter Henderson, Sean Rehaag, Sharad Goel, Shang Gao, Spencer Williams, Sunny
603 Gandhi, Tom Zur, Varun Iyer, and Zehua Li. Legalbench: A collaboratively built benchmark for
604 measuring legal reasoning in large language models, 2023.
- 605 Sören Heimersheim et al. How to use and interpret activation patching. *arXiv preprint*
606 *arXiv:2405.05964*, 2024.
- 607 Dan Hendrycks, Collin Burns, Saurav Kadavath, Akul Arora, Steven Basart, Eric Tang, Dawn Song,
608 and Jacob Steinhardt. Measuring mathematical problem solving with the math dataset. *arXiv*
609 *preprint arXiv:2103.03874*, 2021.
- 610 Zhijing Jin, Jiarui Liu, Zhiheng Lyu, Spencer Poff, Mrinmaya Sachan, Rada Mihalcea, Mona T.
611 Diab, and Bernhard Schölkopf. Can large language models infer causation from correlation? In
612 *International Conference on Learning Representations (ICLR)*, 2024.
- 613 Adam Tauman Kalai, Ofir Nachum, Santosh S. Vempala, and Edwin Zhang. Why language models
614 hallucinate, 2025. URL <https://arxiv.org/abs/2509.04664>.
- 615 Prithish Kamath, Akilesh Tangella, Danica Sutherland, and Nathan Srebro. Does invariant risk
616 minimization capture invariance? In Arindam Banerjee and Kenji Fukumizu, editors, *Proceedings*
617 *of The 24th International Conference on Artificial Intelligence and Statistics*, volume 130 of
618 *Proceedings of Machine Learning Research*, pages 4069–4077. PMLR, 13–15 Apr 2021. URL
619 <https://proceedings.mlr.press/v130/kamath21a.html>.
- 620 Divyansh Kaushik, Eduard Hovy, and Zachary C. Lipton. Explaining the efficacy of counterfactually
621 augmented data. *arXiv preprint arXiv:2010.02114*, 2020. URL [https://arxiv.org/abs/](https://arxiv.org/abs/2010.02114)
622 [2010.02114](https://arxiv.org/abs/2010.02114).
- 623 Pang Wei Koh and Percy Liang. Understanding black-box predictions via influence functions. In
624 *Proceedings of ICML*, 2017. URL [https://proceedings.mlr.press/v70/koh17a/](https://proceedings.mlr.press/v70/koh17a/koh17a.pdf)
625 [koh17a.pdf](https://proceedings.mlr.press/v70/koh17a/koh17a.pdf).
- 626 Zhe Li, Wei Zhao, Yige Li, and Jun Sun. Do influence functions work on large language models?
627 *arXiv preprint arXiv:2409.19998*, 2024.
- 628 Hunter Lightman, Vineet Kosaraju, Yuri Burda, Harrison Edwards, Bowen Baker, Teddy Lee, Jan
629 Leike, John Schulman, Ilya Sutskever, and Karl Cobbe. Let’s verify step by step. In *The Twelfth*
630 *International Conference on Learning Representations*, 2023.
- 631 Stephanie Lin, Jacob Hilton, and Owain Evans. Truthfulqa: Measuring how models mimic human
632 falsehoods. *arXiv preprint arXiv:2109.07958*, 2021.
- 633 Haotian Liu et al. Lms as narcissistic evaluators: When ego inflates evaluation scores. *arXiv preprint*
634 *arXiv:2311.01964*, 2023.
- 635 lucasmccabe. Logiqa. Hugging Face dataset, 2021. URL [https://huggingface.co/](https://huggingface.co/datasets/lucasmccabe/logiqa)
636 [datasets/lucasmccabe/logiqa](https://huggingface.co/datasets/lucasmccabe/logiqa).
- 637 Dung Nguyen Manh, Thang Phan Chau, Nam Le Hai, Thong T Doan, Nam V Nguyen, Quang
638 Pham, and Nghi DQ Bui. Codemmlu: A multi-task benchmark for assessing code understanding
639 capabilities of codellms. *CoRR*, 2024.
- 640 James Martens. Deep learning via hessian-free optimization. In *Proceedings of ICML*, pages
641 735–742, 2010. URL [https://www.cs.toronto.edu/~jmartens/docs/Deep_](https://www.cs.toronto.edu/~jmartens/docs/Deep_HessianFree.pdf)
642 [HessianFree.pdf](https://www.cs.toronto.edu/~jmartens/docs/Deep_HessianFree.pdf).

- 648 Meta AI. Llama-3.2-3b-instruct — model card. [https://huggingface.co/meta-llama/](https://huggingface.co/meta-llama/Llama-3.2-3B-Instruct)
649 [Llama-3.2-3B-Instruct](https://huggingface.co/meta-llama/Llama-3.2-3B-Instruct), 2024a. Accessed 2025-09-09.
- 650
651 Meta AI. Meta-llama-3-8b-instruct — model card. [https://huggingface.co/](https://huggingface.co/meta-llama/Meta-Llama-3-8B-Instruct)
652 [meta-llama/Meta-Llama-3-8B-Instruct](https://huggingface.co/meta-llama/Meta-Llama-3-8B-Instruct), 2024b. Accessed 2025-09-09.
- 653 Paulius Micikevicius, Sharan Narang, Jonah Alben, Gregory Diamos, Erich Elsen, David Garcia,
654 Boris Ginsburg, Michael Houston, Oleksii Kuchaiev, Ganesh Venkatesh, and Hao Wu. Mixed
655 precision training. In *International Conference on Learning Representations (ICLR)*, 2018. URL
656 <https://arxiv.org/abs/1710.03740>.
- 657
658 Microsoft. Phi-3.5-mini-instruct — model card. [https://huggingface.co/microsoft/](https://huggingface.co/microsoft/Phi-3.5-mini-instruct)
659 [Phi-3.5-mini-instruct](https://huggingface.co/microsoft/Phi-3.5-mini-instruct), 2025. Accessed 2025-09-09.
- 660
661 NVIDIA Corporation. Nvidia a100 tensor core gpu architecture. Whitepaper, 2020. URL
662 [https://images.nvidia.com/aem-dam/en-zz/Solutions/data-center/](https://images.nvidia.com/aem-dam/en-zz/Solutions/data-center/nvidia-ampere-architecture-whitepaper.pdf)
663 [nvidia-ampere-architecture-whitepaper.pdf](https://images.nvidia.com/aem-dam/en-zz/Solutions/data-center/nvidia-ampere-architecture-whitepaper.pdf).
- 664 Miruna Oprescu, Vasilis Syrgkanis, and Zhiwei Steven Wu. Orthogonal random forest for causal
665 inference. In *International Conference on Machine Learning*, pages 4932–4941. PMLR, 2019.
- 666
667 Long Ouyang, Jeff Wu, Xu Jiang, Diogo Almeida, Carroll L. Wainwright, Pamela Mishkin, Chong
668 Zhang, and et al. Training language models to follow instructions with human feedback. In
669 *Advances in Neural Information Processing Systems (NeurIPS)*, 2022.
- 670 Long Pan et al. Spontaneous reward hacking in iterative self-refinement. *arXiv preprint*
671 *arXiv:2307.04549*, 2023.
- 672
673 Long Pan et al. Feedback loops with language models drive in-context reward hacking. *arXiv preprint*
674 *arXiv:2402.06627*, 2024.
- 675
676 Judea Pearl. Direct and indirect effects. In *Proceedings of UAI*, 2001. URL [https://ftp.cs.](https://ftp.cs.ucla.edu/pub/stat_ser/r273-jsm05.pdf)
677 [ucla.edu/pub/stat_ser/r273-jsm05.pdf](https://ftp.cs.ucla.edu/pub/stat_ser/r273-jsm05.pdf).
- 678
679 Judea Pearl. *Causality: Models, Reasoning, and Inference*. Cambridge University Press, 2 edition,
680 2009.
- 681
682 Barak A. Pearlmutter. Fast exact multiplication by the hessian. *Neural Computation*, 6(1):
683 147–160, 1994. URL [https://direct.mit.edu/neco/article/6/1/147/5766/](https://direct.mit.edu/neco/article/6/1/147/5766/Fast-Exact-Multiplication-by-the-Hessian)
[Fast-Exact-Multiplication-by-the-Hessian](https://direct.mit.edu/neco/article/6/1/147/5766/Fast-Exact-Multiplication-by-the-Hessian).
- 684
685 Qwen Team. Qwen3-1.7b — model card. <https://huggingface.co/Qwen/Qwen3-1.7B>,
686 2025a. Accessed 2025-09-09.
- 687
688 Qwen Team. Qwen3-4b-thinking-2507 — model card. [https://huggingface.co/Qwen/](https://huggingface.co/Qwen/Qwen3-4B-Thinking-2507)
[Qwen3-4B-Thinking-2507](https://huggingface.co/Qwen/Qwen3-4B-Thinking-2507), 2025b. Accessed 2025-09-09.
- 689
690 Rafael Rafailov, Yaswanth Chittepudi, Ryan Park, Harshit Sushil Sikchi, Joey Hejna, Brad Knox,
691 Chelsea Finn, and Scott Niekum. Scaling laws for reward model overoptimization in direct
692 alignment algorithms. *Advances in Neural Information Processing Systems*, 37:126207–126242,
693 2024.
- 694
695 Shauli Ravfogel, Yanai Elazar, Hila Gonen, Michael Twiton, and Yoav Goldberg. Null it out:
696 Guarding protected attributes by iterative nullspace projection. *arXiv preprint arXiv:2004.07667*,
697 2020.
- 698
699 Diederik M Roijers, Peter Vamplew, Shimon Whiteson, and Richard Dazeley. A survey of multi-
700 objective sequential decision-making. *Journal of Artificial Intelligence Research*, 48:67–113,
701 2013.
- 702
703 Elan Rosenfeld, Pradeep Ravikumar, and Andrej Risteski. The risks of invariant risk minimization.
arXiv preprint arXiv:2010.05761, 2020.

- 702 Andrew Slavin Ross, Michael C. Hughes, and Finale Doshi-Velez. Right for the right reasons:
703 Training differentiable models by constraining their explanations. In *International Joint Conference*
704 *on Artificial Intelligence (IJCAI)*, 2017.
- 705 John Schulman, Sergey Levine, Pieter Abbeel, Michael Jordan, and Philipp Moritz. Trust region
706 policy optimization. In *International conference on machine learning*, pages 1889–1897. PMLR,
707 2015.
- 708 John Schulman, Filip Wolski, Prafulla Dhariwal, Alec Radford, and Oleg Klimov. Proximal policy
709 optimization algorithms. *arXiv preprint arXiv:1707.06347*, 2017.
- 710 Zhen Shao et al. Deepseekmath: Pushing the limits of mathematical reasoning in open language
711 models. *arXiv preprint arXiv:2402.03300*, 2024. Introduces GRPO.
- 712 Guangming Sheng, Chi Zhang, Zilingfeng Ye, Xibin Wu, Wang Zhang, Ru Zhang, Yanghua Peng,
713 Haibin Lin, and Chuan Wu. Hybridflow: A flexible and efficient rlhf framework. *arXiv preprint*
714 *arXiv: 2409.19256*, 2024.
- 715 Ivaxi Sheth, Bahare Fatemi, and Mario Fritz. Causalgraph2llm: Evaluating llms for causal queries.
716 In *Findings of ACL: NAACL 2025*, 2025.
- 717 Rahul Shrama et al. Towards understanding sycophancy in language models. *arXiv preprint*
718 *arXiv:2310.13548*, 2023.
- 719 Shang Hong Sim, Tej Deep Pala, Vernon Toh, Hai Leong Chieu, Amir Zadeh, Chuan Li, Navonil
720 Majumder, and Soujanya Poria. Lessons from training grounded llms with verifiable rewards.
721 *arXiv preprint arXiv:2506.15522*, 2025.
- 722 Nisan Stiennon, Long Ouyang, Jeffrey Wu, Daniel M. Ziegler, Ryan Lowe, and et al. Learning
723 to summarize with human feedback. In *Advances in Neural Information Processing Systems*
724 *(NeurIPS)*, 2020.
- 725 Harsha Tanneru et al. On the hardness of verifying chain-of-thought faithfulness. *arXiv preprint*
726 *arXiv:2411.04047*, 2024.
- 727 Tyler J. VanderWeele. Controlled direct and mediated effects: Definition, identification and bounds.
728 *Scandinavian Journal of Statistics*, 38(3):551–563, 2011. doi: 10.1111/j.1467-9469.2010.00722.x.
729 URL <https://pmc.ncbi.nlm.nih.gov/articles/PMC4193506/>.
- 730 Roman Vershynin. *High-dimensional probability: An introduction with applications in data science*,
731 volume 47. Cambridge university press, 2018.
- 732 Boxin Wang et al. Large language models are not fair evaluators. *arXiv preprint arXiv:2310.05492*,
733 2023.
- 734 Chaoqi Wang, Zhuokai Zhao, Yibo Jiang, Zhaorun Chen, Chen Zhu, Yuxin Chen, Jiayi Liu, Lizhu
735 Zhang, Xiangjun Fan, Hao Ma, et al. Beyond reward hacking: Causal rewards for large language
736 model alignment. *arXiv preprint arXiv:2501.09620*, 2025.
- 737 Jason Wei, Xuezhi Wang, Dale Schuurmans, Maarten Bosma, Brian Ichter, Fei Xia, Ed H. Chi, Quoc
738 Le, and Denny Zhou. Chain-of-thought prompting elicits reasoning in large language models. In
739 *Advances in Neural Information Processing Systems (NeurIPS)*, 2022.
- 740 Yilun Wen et al. Language models learn to mislead humans via rlhf. *arXiv preprint arXiv:2409.12822*,
741 2024.
- 742 Joong-Ho Won, Jason Xu, and Kenneth Lange. Projection onto minkowski sums with application to
743 constrained learning. In *International Conference on Machine Learning*, pages 3642–3651. PMLR,
744 2019.
- 745 Liang Yu et al. Towards better causal reasoning in language models. In *NAACL Long Papers*, 2025.
- 746 Wenting Yu, Meng Jiang, Peter Clark, and Ashish Sabharwal. Ifqa: Counterfactual open-domain qa.
747 GitHub repository, 2023. URL <https://github.com/wyu97/IfQA>.

756 Weizhe Yuan, Richard Yuanzhe Pang, Kyunghyun Cho, Xian Li, Sainbayar Sukhbaatar, Jing Xu, and
757 Jason Weston. Self-rewarding language models. *arXiv preprint arXiv:2401.10020*, 2025.
758

759 Tianyi Zhang, Varsha Kishore, Felix Wu, Kilian Q Weinberger, and Yoav Artzi. Bertscore: Evaluating
760 text generation with bert. *arXiv preprint arXiv:1904.09675*, 2019.

761 Lianmin Zheng, Wei-Lin Chiang, Ying Sheng, Siyuan Zhuang, Zhanghao Wu, Yonghao Zhuang,
762 Zi Lin, Zhuohan Li, Dacheng Li, Eric Xing, et al. Judging llm-as-a-judge with mt-bench and
763 chatbot arena. *Advances in neural information processing systems*, 36:46595–46623, 2023.
764

765 Yujun Zhou, Zhenwen Liang, Haolin Liu, Wenhao Yu, Kishan Panaganti, Linfeng Song, Dian Yu,
766 Xiangliang Zhang, Haitao Mi, and Dong Yu. Evolving language models without labels: Majority
767 drives selection, novelty promotes variation. *arXiv preprint arXiv:2509.15194*, 2025.
768
769
770
771
772
773
774
775
776
777
778
779
780
781
782
783
784
785
786
787
788
789
790
791
792
793
794
795
796
797
798
799
800
801
802
803
804
805
806
807
808
809

810 A RELATED WORKS

811 A.1 REWARD HACKING IN LLM

812 Reward hacking occurs when models exploit imperfections in reward functions rather than achieving
 813 the intended objective (Amodei et al., 2016; Everitt et al., 2017; Zhou et al., 2025). In large language
 814 models, this problem emerges prominently under alignment: proxy reward models can be gamed,
 815 leading to behaviors such as generating superficially persuasive but incorrect answers Wen et al.
 816 (2024), sycophancy Shrama et al. (2023), and exploiting evaluator biases in LLM-as-judge settings Liu
 817 et al. (2023); Wang et al. (2023). More recently, iterative refinement has revealed *in-context reward*
 818 *hacking*, where models adapt outputs to maximize evaluation scores rather than genuinely improving
 819 reasoning quality Pan et al. (2023; 2024). Mitigation strategies typically involve improved supervision
 820 (e.g., rationale-level feedback Ross et al. (2017)), evaluator debiasing, or structural interventions such
 821 as counterfactual augmentation Kaushik et al. (2020) and invariance penalties Arjovsky et al. (2020).
 822 However, these often rely on external annotations nor emphasized robustness–accuracy trade-offs.
 823 Our method introduces differentiable, counterfactual-enhanced causal signals as internal rewards,
 824 jointly optimizing accuracy and causal coherence to directly curb shortcut exploitation.
 825

826 A.2 CAUSAL INFERENCE AND REASONING IN LLMs

827 There have been several studies centering on causal reasoning capability in LLMs. Corr2Cause finds
 828 poor correlation-to-causation transfer by pure LLM reasoning (Jin et al., 2024), CoT often fails to
 829 be causally responsible for answers (Bao et al., 2024a), and faithful CoT remains challenging even
 830 with fine-tuning or activation edits (Tanneru et al., 2024). Methods that *impose* causal structure are
 831 emerging (e.g., G²-Reasoner combines knowledge and goal prompts to improve causal reasoning)
 832 (Chi et al., 2025). Mechanistic approaches intervene directly in internal representations—editing
 833 factual circuits and testing hypotheses via causal/activation patching—to probe and modify pathways
 834 (Heimersheim et al., 2024). There also several studies utilizing SCM (Yu et al., 2025) and causal
 835 graph (Sheth et al., 2025) for better reasoning. Yet most studies are *post-hoc*. We build on this
 836 direction by using interventional, pathway-aware signals during post-training, and by explicitly
 837 aggregating them with task rewards to control the accuracy–causal-coherence trade-off.
 838

839 B THEORETICAL GUARANTEES

840 **Setup.** Let π, π_{old} be stationary policies, $\gamma \in (0, 1)$ the discount factor, and d^π the normalized
 841 discounted state distribution. Fix the baseline policy π_{old} and let $A(s, a) = A_{\pi_{\text{old}}}(s, a)$ be its
 842 advantage under the CE-PO reward r_{CE} (the Minkowski combiner of all channels). For each
 843 state s , denote by $\mathcal{N}_s \subset \mathbb{R}^{|\mathcal{A}|}$ the subspace spanned by *nuisance directions* (e.g., verbosity/length)
 844 estimated via counterfactual perturbations, and let P_s be the Euclidean orthogonal projector onto
 845 \mathcal{N}_s . Write the action simplex $\Delta^{|\mathcal{A}|} = \{p \in \mathbb{R}^{|\mathcal{A}|} : p \geq 0, \mathbf{1}^\top p = 1\}$ with tangent space
 846 $\mathcal{T}_s = \{v \in \mathbb{R}^{|\mathcal{A}|} : \mathbf{1}^\top v = 0\}$. Define the *orthogonalized advantage*
 847

$$848 A_\perp(s, \cdot) := (I - P_s) A(s, \cdot), \quad \varepsilon_\perp := \max_s \|A_\perp(s, \cdot)\|_\infty.$$

849 and let

$$850 \alpha := \max_s D_{\text{TV}}(\pi(\cdot|s), \pi_{\text{old}}(\cdot|s)) = \frac{1}{2} \max_s \|\pi(\cdot|s) - \pi_{\text{old}}(\cdot|s)\|_1.$$

851 The TRPO-style surrogate is

$$852 L_{\pi_{\text{old}}}^{\text{CE}}(\pi) = \eta_{\text{CE}}(\pi_{\text{old}}) + \frac{1}{1 - \gamma} \mathbb{E}_{s \sim d^{\pi_{\text{old}}}} [\langle A(s, \cdot), \pi(\cdot|s) - \pi_{\text{old}}(\cdot|s) \rangle],$$

853 noting that $\langle A(s, \cdot), \pi_{\text{old}}(\cdot|s) \rangle = 0$ for all s .

860 Assumptions and Propositions

861 **P0 (Fixed reward at π_{old})** All CE components used to construct r_{CE} (Jacobians, counterfactual
 862 subspaces \mathcal{N}_s , projectors P_s , and channel scores) are computed at π_{old} and held fixed while
 863 evaluating $\eta_{\text{CE}}(\pi)$ over one update step.

A0 (Residualization on the action-simplex tangent) By construction of r_{CE} via counterfactual projection, the CE advantage is locally invariant along nuisance directions at π_{old} :

$$\langle A(s, \cdot), v \rangle = 0 \quad \text{for all } v \in \mathcal{N}_s \cap \mathcal{T}_s.$$

Equivalently, $P_s A(s, \cdot) = 0$ and hence $A(s, \cdot) = A_{\perp}(s, \cdot) \in \mathcal{N}_s^{\perp}$.

A1 (Well-posedness & bounded rewards) Base rewards and channel scores are in $[0, 1]$ and CE weights sum to one, so $r_{\text{CE}} \in [0, 1]$. The MDP is standard (finite or σ -finite), ensuring existence of $V^{\pi}, Q^{\pi}, A_{\pi}$.

Theorem 1 (Orthogonalized TRPO Lower Bound for CE-PO). *For any π , with $\alpha = \max_s D_{\text{TV}}(\pi(\cdot|s), \pi_{\text{old}}(\cdot|s))$, the CE-PO performance satisfies*

$$\eta_{\text{CE}}(\pi) \geq L_{\pi_{\text{old}}}^{\text{CE}}(\pi) - \frac{4\gamma}{(1-\gamma)^2} \varepsilon_{\perp} \alpha^2. \quad (8)$$

Proof. Step 1: Performance-difference decomposition. The performance-difference lemma gives

$$\eta_{\text{CE}}(\pi) - \eta_{\text{CE}}(\pi_{\text{old}}) = \frac{1}{1-\gamma} \mathbb{E}_{s \sim d^{\pi}} [\langle A(s, \cdot), \pi(\cdot|s) \rangle].$$

Add and subtract the surrogate:

$$\eta_{\text{CE}}(\pi) = L_{\pi_{\text{old}}}^{\text{CE}}(\pi) + \frac{1}{1-\gamma} \mathbb{E}_s \left[(d^{\pi} - d^{\pi_{\text{old}}})(s) \langle A(s, \cdot), \pi(\cdot|s) \rangle \right] =: L_{\pi_{\text{old}}}^{\text{CE}}(\pi) - \frac{1}{1-\gamma} T,$$

so it suffices to bound $|T|$.

Step 2: Bounding the occupancy shift. By a standard coupling argument, $\|d^{\pi} - d^{\pi_{\text{old}}}\|_1 \leq \frac{2\gamma}{1-\gamma} \alpha$. Thus

$$|T| \leq \left(\max_s |\langle A(s, \cdot), \pi(\cdot|s) \rangle| \right) \|d^{\pi} - d^{\pi_{\text{old}}}\|_1 \leq \frac{2\gamma}{1-\gamma} \alpha \max_s |\langle A(s, \cdot), \pi(\cdot|s) \rangle|.$$

Step 3: Using CE-PO orthogonalization (A1). For each s , by A1 we have $A(s, \cdot) = A_{\perp}(s, \cdot) \in \mathcal{N}_s^{\perp}$, and $\langle A(s, \cdot), \pi_{\text{old}}(\cdot|s) \rangle = 0$ (definition of advantage). Hence

$$\langle A(s, \cdot), \pi(\cdot|s) \rangle = \langle A_{\perp}(s, \cdot), \pi(\cdot|s) - \pi_{\text{old}}(\cdot|s) \rangle \leq \|A_{\perp}(s, \cdot)\|_{\infty} \|\pi(\cdot|s) - \pi_{\text{old}}(\cdot|s)\|_1 \leq 2\varepsilon_{\perp} \alpha.$$

Taking the supremum over s yields $\max_s |\langle A(s, \cdot), \pi(\cdot|s) \rangle| \leq 2\varepsilon_{\perp} \alpha$, and therefore

$$|T| \leq \frac{2\gamma}{1-\gamma} \alpha (2\varepsilon_{\perp} \alpha) = \frac{4\gamma}{1-\gamma} \varepsilon_{\perp} \alpha^2.$$

Step 4: Final bound. Plugging into Step 1,

$$\eta_{\text{CE}}(\pi) \geq L_{\pi_{\text{old}}}^{\text{CE}}(\pi) - \frac{4\gamma}{(1-\gamma)^2} \varepsilon_{\perp} \alpha^2,$$

which is inequality (8). \square

Properties. (*Tightening via nuisance removal, ℓ_2 statement*). Let $\varepsilon_{\text{base}}^{(2)} := \max_s \|A_{\text{base}}(s, \cdot)\|_2$ be the corresponding constant under a non-residualized baseline reward, and define $\varepsilon_{\perp}^{(2)} := \max_s \|A_{\perp}(s, \cdot)\|_2$. If there exists $\rho \in (0, 1]$ such that $\|P_s A_{\text{base}}(s, \cdot)\|_2 \geq \rho \|A_{\text{base}}(s, \cdot)\|_2$ for all s (*nontrivial nuisance energy*), then by Pythagoras for orthogonal projectors,

$$\varepsilon_{\perp}^{(2)} \leq \sqrt{1-\rho^2} \varepsilon_{\text{base}}^{(2)}.$$

Since $\|\cdot\|_{\infty} \leq \|\cdot\|_2$, we also have $\varepsilon_{\perp} \leq \varepsilon_{\perp}^{(2)}$, yielding a strictly smaller penalty constant in the ℓ_2 analogue of equation 8 (and a looser but still improved constant in ℓ_{∞} up to norm equivalence).

(*Null penalty for hacking-only updates*). If for some s the per-state policy change lies entirely in \mathcal{N}_s , i.e., $\pi(\cdot|s) - \pi_{\text{old}}(\cdot|s) \in \mathcal{N}_s$, then $(I - P_s)(\pi - \pi_{\text{old}}) = 0$ and hence $\langle A_{\perp}(s, \cdot), \pi(\cdot|s) - \pi_{\text{old}}(\cdot|s) \rangle = 0$, so that state contributes nothing to the penalty term in equation 8.

KL-based corollaries. By Pinsker’s inequality $D_{\text{TV}}(p, q) \leq \sqrt{\frac{1}{2}\text{KL}(p\|q)}$, if $\delta := \max_s \text{KL}(\pi_{\text{old}}(\cdot|s)\|\pi(\cdot|s))$, then

$$\eta_{\text{CE}}(\pi) \geq L_{\pi_{\text{old}}}^{\text{CE}}(\pi) - \frac{2\gamma}{(1-\gamma)^2} \varepsilon_{\perp} \delta. \tag{9}$$

If only an *expected* per-state KL is enforced, $\mathbb{E}_{s \sim d^{\pi_{\text{old}}}} [\text{KL}(\pi_{\text{old}}\|\pi)] \leq \bar{\delta}$, then equation 9 holds with $\delta = C \bar{\delta}$ for a method-dependent constant $C \geq 1$ (e.g., via ratio clipping or per-state caps used in TRPO/PPO).

From counterfactual Jacobians to action-space nuisance. Let \mathcal{U}_s be a finite set of representation-space directions obtained from counterfactual Jacobians (e.g., top singular vectors of counterfactual J_{AB} blocks), and let g_s be a smooth channel-score functional (e.g., spectral energy share). Write $G_s := \nabla_{\pi(\cdot|s)} g_s$ for its Jacobian with respect to action probabilities. Define

$$\mathcal{N}_s := \text{span}\{G_s^{\top} u : u \in \mathcal{U}_s\} \cap \mathcal{T}_s.$$

If each channel score is *residualized* by projecting J_{AB} onto \mathcal{U}_s^{\perp} , then, at π_{old} , the resulting CE reward satisfies $\langle \nabla_{\pi} r_{\text{CE}}(s, \cdot), v \rangle = 0$ for all $v \in \mathcal{N}_s \cap \mathcal{T}_s$, implying $P_s A(s, \cdot) = 0$ (Assumption **A0**).

C ADDITIONAL TECHNIQUES

Jacobian causal score is time and space-efficient in LLMs. We design the Jacobian-based causal reward to be *practical at scale*: it reuses forward activations, avoids any explicit Jacobian/Hessian materialization, and replaces costly subspace operations with a constant-time projection. Concretely, we rely on three tricks. (i) *Local Jacobian via vJP (vector–Jacobian Projection)*. Sensitivities are computed with one forward and two backward passes using vector–Jacobian products (reverse-mode autodiff), i.e., for a test direction U we use $\text{vJP}(s_B, E(A); U) = \langle J_{AB}, U \rangle_F = J_{AB}^{\top} \text{vec}(U)$ and the scalar sensitivity $\|J_{AB}^{\top} \text{vec}(U)\|_2$; so no dense $J \in \mathbb{R}^{Td \times Td}$ (or Hessian) is ever formed; working memory stays at $\mathcal{O}(Td)$ for token embeddings (Baydin et al., 2018; Pearlmutter, 1994). (ii) *Residual projection*. The counterfactual direction is removed with a single inner product: if g is the local Jacobian direction and \hat{c} the (reshuffled) counterfactual direction, we use $r = g - (g^{\top} \hat{c}) \hat{c}$, avoiding SVD/subspace construction and extra buffers (Golub and Van Loan, 2013). (iii) *Mixed precision*. Forward/backward in FP16/bfloat16 with dynamic loss scaling provides wall-clock gains while roughly halving activation memory, without changing the algorithmic interface (Micikevicius et al., 2018; NVIDIA Corporation, 2020).

Throughput and footprint. On a single modern GPU (e.g., A100 80 GB) and Llama-3-8B-Instruct, our implementation returns a reward in ≈ 2 s, with negligible persistent memory beyond base activations. Simple Jacobian takes $\mathcal{O}(Td)$ space complexity and would take $\mathcal{O}(C_{\text{fwd}})$ time complexity, where C_{fwd} means the time for one single forward in LLM reasoning. Yet for Hessian score, there needs to be continuous iteration of simple Jacobian to converge for a stable score and space complexity would also be $\mathcal{O}((Td)^2)$, which is not practical. To provide a more quantitative analysis, we benchmarked our Jacobian computation against a BERTScore baseline on Qwen-3-14B (on a single node of 4 A100 GPU) using 100 randomly generated samples. We found the baseline took 1.20s and peaked at 28.00 GB, while our Jacobian step took 6.98s and peaked at 30.18 GB. This represents a 5.82x wall-clock time overhead and, more importantly, only a 1.078x GPU memory usage multiplier.

D ADDITIONAL PLOTS

Question. Context: For husbands that don’t set the alarm, the probability of ringing alarm is 90%. For husbands that set the alarm, the probability of ringing alarm is 25%. Question: For husbands that set the alarm, would it be less likely to see ringing alarm if the husband had not set the alarm?

Ground Truth: No

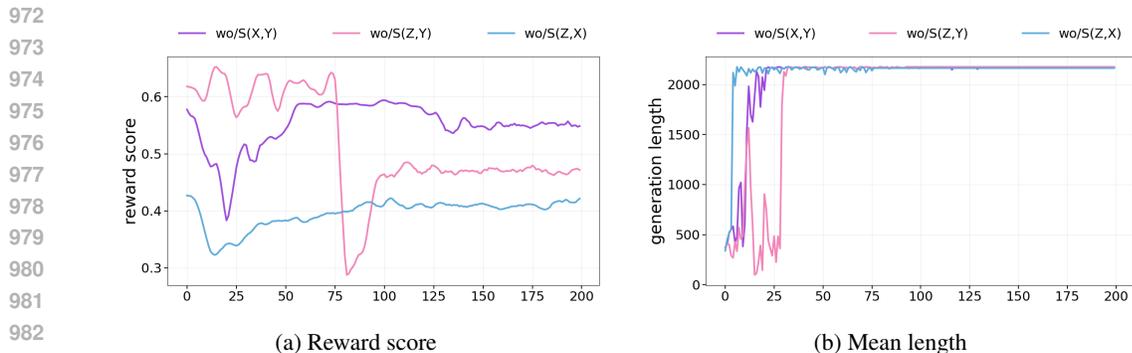


Figure 11: This figure supplements the ablation analysis and demonstrates that omitting a single Jacobian term induces reward degradation accompanied by length inflation.

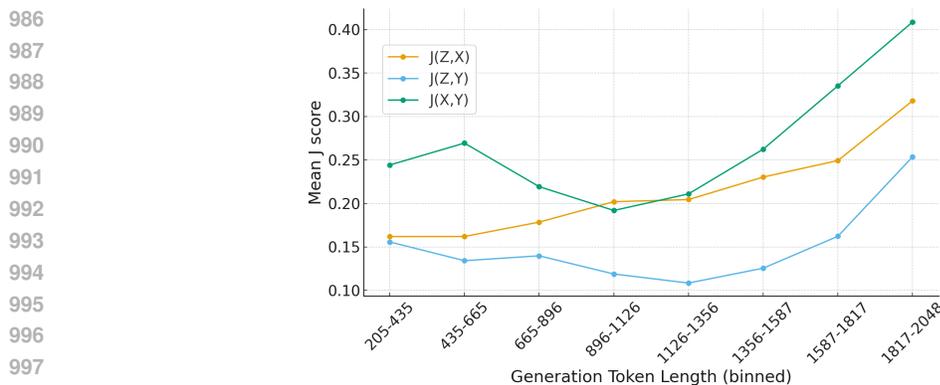


Figure 13: This figure illustrates the binned raw Jacobian scores distribution among different generation token length, echoing the reward hacking by length statement in Figure 2.

Method	Prediction (Reasoning)
BERTScore	Predicts “Yes”; relies on surface overlap (...compares 25% vs 90% directly but ignores that the subject of “less likely” is “not set”).
Jacobian Only	Predicts “Yes”; observes 25% < 90% but misplaces the subject, treating “set” as the focus of “less likely” instead of “not set”.
CE-GRPO	Predicts “No”; correctly checks whether 90% < 25% (false) and concludes “not set” actually makes ringing more likely, matching ground truth.

Takeaways.

- **Jacobian Only** tends to *template-follow*: latches onto “less likely \Rightarrow smaller %” and ignores subject/polarity flips, causing causal misread.
- **BERTScore** measures overlap, not causal/conditional structure; it cannot reason about which condition the “less likely” modifies.

Figure 12: Comparison of methods on the alarm-setting counterfactual example on Qwen-3-4B-Thinking.

D.1 RESPONSE TEXT OF FIGURE 12

CE-GRPO: <think> Okay, let’s try to figure out this probability question. First, I need to understand the context. The problem says that for husbands who don’t set the alarm, there’s a 90% chance the alarm rings. For those who do set the alarm, the probability is 25%. The question is asking whether it’s less likely to see the alarm ring if the husband had not set the alarm. Wait, that seems a bit confusing. Let me break it down.

1026 The question is about the probability of the alarm ringing given whether the husband set the alarm or not. The
 1027 context gives two probabilities: 90% for not setting the alarm and 25% for setting it. But the question is phrased
 1028 a bit oddly. It says, "would it be less likely to see ringing alarm if the husband had not set the alarm?" Hmm.
 1029 Maybe I need to rephrase that. Wait, maybe the question is asking if a husband didn't set the alarm, is it less
 1030 likely that the alarm rings compared to if he did set it? But that's exactly what the context says. Because 90% is
 1031 higher than 25%, so if the husband didn't set the alarm (90% chance), it's more likely the alarm rings. But the
 1032 question is phrased as "less likely to see ringing alarm if the husband had not set the alarm," which would be the
 1033 opposite. So maybe the answer is no, it's more likely. But maybe I'm misunderstanding the question. Let me
 1034 check again. The context is about the probability of the alarm ringing given whether the husband set it. The
 1035 question is asking if the husband didn't set the alarm, is it less likely to see the alarm ring? But according to the
 1036 context, it's 90%, which is higher than 25%, so it's more likely. Therefore, the answer is no, it's not less likely;
 1037 it's more likely. But maybe the question is trying to trick me into thinking about the reverse. Like, if the alarm
 1038 doesn't ring, is it less likely the husband set it? But that's a different question. The original question is about the
 1039 probability of the alarm ringing given whether the husband set it. So the answer is that it's more likely when the
 1040 husband doesn't set the alarm. Therefore, the answer is no, it's not less likely. </think> The context states that
 1041 husbands who **do not set the alarm** have a **90% probability** of the alarm ringing, while those who **do**
 1042 **set the alarm** have a **25% probability**. The question asks whether it is **less likely** to see the alarm ring
 1043 if the husband **did not set the alarm**. This is directly contradicted by the given probabilities: - **90% (high**
 1044 **probability)** of the alarm ringing when the husband **does not set** the alarm. - **25% (low probability)**
 1045 of the alarm ringing when the husband **does set** the alarm. Thus, it is **not less likely** to see the alarm
 1046 ring if the husband did not set the alarm. Instead, it is **more likely** (90% vs. 25%). **Answer:** No, it is
 1047 more likely to see the alarm ring if the husband did not set the alarm (90% vs. 25%)."

1045 **BERTScore:** "<think> Okay, let's try to figure out this probability question. The user is asking whether it's
 1046 less likely for a husband to ring the alarm if he had not set the alarm, compared to when he did set it. Wait, the
 1047 question is a bit confusing. Let me parse it again. The context says: For husbands that don't set the alarm, the
 1048 probability of ringing the alarm is 90%. For those that do set the alarm, it's 25%. The question is: For husbands
 1049 that set the alarm, would it be less likely to see ringing alarm if the husband had not set the alarm? Hmm, that
 1050 seems contradictory. Wait, maybe there's a typo or misunderstanding. Let me check again. Wait, the question is
 1051 phrased as: "For husbands that set the alarm, would it be less likely to see ringing alarm if the husband had not
 1052 set the alarm?" That doesn't make sense. If a husband has set the alarm, then he's in the "set the alarm" group.
 1053 The question is asking if, for those who set the alarm, would it be less likely to ring the alarm if they hadn't set
 1054 it? That seems like a paradox. Maybe the user meant to ask: For husbands that set the alarm, is it less likely to
 1055 ring the alarm compared to those who didn't set it? Because the numbers are 25% vs 90%. So, if a husband
 1056 sets the alarm, the probability of ringing is 25%, which is lower than 90% for those who don't set it. So, yes,
 1057 it's less likely. But the way the question is phrased is confusing. Let me make sure I'm interpreting it correctly.
 1058 Alternatively, maybe the question is trying to ask: If a husband didn't set the alarm, is it less likely to ring the
 1059 alarm compared to if he had set it? But that's not what's written. The original question says: "For husbands
 1060 that set the alarm, would it be less likely to see ringing alarm if the husband had not set the alarm?" So, it's a
 1061 conditional. If the husband had not set the alarm, would it be less likely? But that's not a standard way to phrase
 1062 it. Maybe the user is confused. Let me think again. Wait, maybe the question is trying to ask: For husbands who
 1063 set the alarm, is it less likely that the alarm rings compared to those who don't set it? Because the numbers are
 1064 25% vs 90%. So, yes, setting the alarm reduces the probability of the alarm ringing. But the way the question is
 1065 phrased is confusing. The user might be mixing up the conditions. Let me try to rephrase the question. The user
 1066 might be asking: If a husband sets the alarm, is it less likely that the alarm rings compared to if he didn't set it?
 1067 In that case, yes, because 25% is less than 90%. So the answer would be yes. But the original question is phrased
 1068 as "For husbands that set the alarm, would it be less likely to see ringing alarm if the husband had not set the
 1069 alarm?" Which is a bit of a paradox. Maybe it's a translation error or a misphrasing. But based on the numbers
 1070 given, the answer is yes. So, the answer is yes, it's less likely. But I need to make sure I'm not misinterpreting.
 1071 Let me check again. The question is: For husbands that set the alarm, would it be less likely to see ringing alarm
 1072 if the husband had not set the alarm? So, if the husband had not set the alarm, would it be less likely? But that's
 1073 not a standard question. Maybe the user meant: For husbands that set the alarm, is it less likely that the alarm
 1074 rings compared to those who didn't set it? In that case, yes. So, the answer is yes. But the way the question is
 1075 phrased is confusing. However, given the numbers, the answer is yes. </think> The question seems to have a
 1076 slight ambiguity, but based on the context provided, here's the clarification: - **If a husband sets the alarm**,
 1077 the probability of the alarm ringing is **25%**. - **If a husband does not set the alarm**, the probability of
 1078 the alarm ringing is **90%**. Thus, **it is less likely** for the alarm to ring if the husband **set the alarm**
 1079 (25%) compared to if he **did not set it** (90%). **Answer:** Yes, it is less likely for the alarm to ring if the
 1080 husband set the alarm (25%) compared to if he did not (90%). The probability is lower when the alarm is set."

1076 **Jacobian-Only:** "<think> Okay, let's try to figure out this probability question. So the context is about husbands
 1077 and alarms. The problem says that for husbands who don't set the alarm, the probability of the alarm ringing
 1078 is 90%. And for those who do set the alarm, the probability is 25%. The question is asking: For husbands
 1079 that set the alarm, would it be less likely to see the alarm ring if the husband had not set the alarm? Wait, that
 1080 seems a bit confusing. Let me parse that again. Wait, the question is phrased as: "For husbands that set the
 1081 alarm, would it be less likely to see ringing alarm if the husband had not set the alarm?" Hmm. Maybe there's

1080 a typo or confusion in the wording. Let me think. Maybe it's asking whether, given that a husband set the
 1081 alarm, is it less likely that the alarm rings compared to if he hadn't set it? But that doesn't make sense because
 1082 the probability of the alarm ringing is given based on whether the husband set it or not. So if a husband set
 1083 the alarm, the probability is 25%, and if he didn't, it's 90%. So, if a husband set the alarm, the probability of
 1084 the alarm ringing is 25%, which is lower than 90% if he hadn't set it. So, in that case, yes, it would be less
 1085 likely to see the alarm ring if the husband had set it (25%) compared to if he hadn't (90%). But the question
 1086 is phrased as "For husbands that set the alarm, would it be less likely to see ringing alarm if the husband had
 1087 not set the alarm?" That seems contradictory. Wait, maybe it's a trick question. Let me try to rephrase. The
 1088 question is asking: If a husband sets the alarm, is it less likely that the alarm rings compared to if he had not
 1089 set the alarm? But that's not the case. Because if he set the alarm, the probability is 25%, and if he didn't, it's
 1090 90%. So, actually, it's more likely that the alarm rings if he didn't set it. So the answer would be no, it's not less
 1091 likely. Wait, but the question is phrased as "would it be less likely to see ringing alarm if the husband had not
 1092 set the alarm?" So maybe the question is trying to ask if, given that a husband set the alarm, is it less likely to
 1093 see the alarm ring compared to if he had not set it. But that's not the case. The probability is lower when he
 1094 sets the alarm. So the answer would be yes, it's less likely. But the wording is confusing. Let me check again.
 1095 Original question: "For husbands that set the alarm, would it be less likely to see ringing alarm if the husband
 1096 had not set the alarm?" So, the phrasing is a bit odd. Maybe it's trying to ask: If a husband set the alarm, is it
 1097 less likely that the alarm rings compared to if he had not set the alarm? In that case, the answer is yes, because
 1098 25% is less than 90%. But the way the question is phrased is confusing. Maybe it's a trick question where they
 1099 are trying to test understanding of conditional probability. Alternatively, maybe it's a case of confusing the two
 1100 probabilities. Let me think again. The problem gives two probabilities: $P(\text{ring} \mid \text{not set}) = 90\%$, and $P(\text{ring} \mid \text{set}) = 25\%$.
 1101 The question is asking: For husbands that set the alarm, would it be less likely to see ringing alarm
 1102 if the husband had not set the alarm? Wait, that seems to be mixing up the conditions. If a husband set the alarm,
 1103 then the probability of the alarm ringing is 25%. If he had not set the alarm, the probability is 90%. So, if we are
 1104 considering a husband who set the alarm, the probability of the alarm ringing is 25%, which is less than 90% if
 1105 he hadn't set it. Therefore, the answer is yes, it would be less likely. But the question is phrased in a way that
 1106 seems to be asking if, given that the husband set the alarm, is it less likely to see the alarm ring compared to if
 1107 he had not set it. So the answer is yes. But maybe the question is trying to trick you into thinking about it in a
 1108 different way. Alternatively, maybe it's a case of confusing the two probabilities. Let me check if there's any
 1109 other interpretation. Alternatively, maybe the question is asking: If a husband set the alarm, is it less likely that
 1110 the alarm rings compared to if he had not set it? Which is exactly what the numbers say. So the answer is yes.
 1111 But the way the question is phrased is confusing. Maybe the user made a typo. Alternatively, maybe it's a case
 1112 of confusing the conditional probabilities. For example, if you have a husband who set the alarm, the probability
 1113 of the alarm ringing is 25%, which is lower than 90% if he hadn't set it. Therefore, the answer is yes, it would
 1114 be less likely. So the answer should be yes. But I need to make sure I'm not missing anything here. Let me think
 1115 again. The question is phrased as: For husbands that set the alarm, would it be less likely to see ringing alarm if
 1116 the husband had not set the alarm? So, the phrasing is a bit confusing. Maybe it's trying to ask: If a husband set
 1117 the alarm, is it less likely that the alarm rings compared to if he had not set it? Which is the case. So the answer
 1118 is yes. Therefore, the answer is yes, it would be less likely. But maybe the user is confused because they are
 1119 mixing up the conditions. So the answer is yes, because the probability of the alarm ringing is 25% when the
 1120 husband set it, which is less than 90% when he didn't. Therefore, for husbands that set the alarm, it is less likely
 1121 to see the alarm ring compared to if they had not set it. So the answer is yes. </think> The question is asking
 1122 whether, for husbands who set the alarm, it would be less likely to see the alarm ring if the husband had not set
 1123 the alarm. **Key Information:** - **If a husband does not set the alarm**, the probability of the alarm ringing
 1124 is **90%**. - **If a husband sets the alarm**, the probability of the alarm ringing is **25%**. **Analysis:** -
 1125 The probability of the alarm ringing is **lower (25%)** when the husband sets the alarm compared to when he
 1126 does not (90%). - Therefore, for husbands who set the alarm, it is **less likely** to see the alarm ring than if the
 1127 husband had not set it. **Answer:** Yes, it would be less likely to see the alarm ring if the husband had not set
 1128 the alarm. The probability of the alarm ringing is 25% when the husband sets it, which is significantly lower
 1129 than the 90% probability when he does not set it."

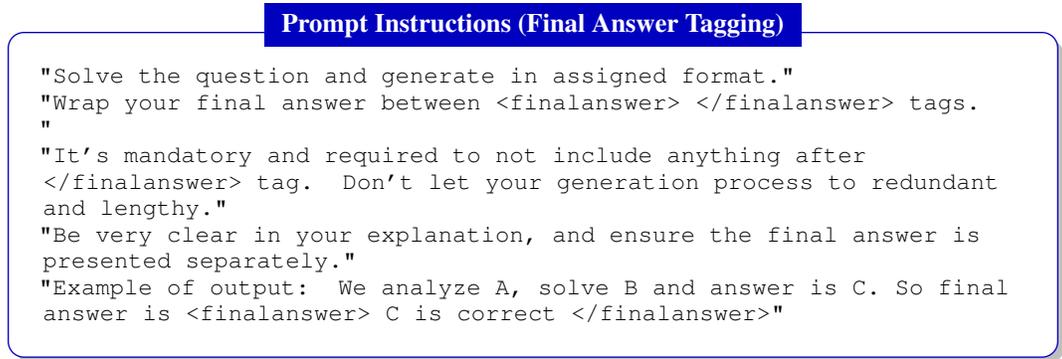
1123 E BROADER IMPACT

1125 The proposed CE-PO framework has the potential to substantially improve the reliability and trustworthiness
 1126 of large language models. By explicitly rewarding causal coherence in addition to task accuracy, the method
 1127 reduces the tendency of models to exploit superficial cues or engage in reward hacking, which can lead to
 1128 misleading or unfaithful reasoning. In practice, this enables LLMs to generate explanations and rationales that
 1129 are not only correct in outcome but also grounded in meaningful reasoning processes—an essential property for
 1130 high-stakes domains such as law, science, medicine, and education.
 1131 Beyond improving alignment, CE-PO can serve as a general methodology for building AI systems that are
 1132 "right for the right reasons," supporting transparency, auditability, and robustness under distribution shifts. This
 1133 contributes to safer deployment of generative models in real-world decision-making pipelines. At the same time,
 emphasizing causal regularization may help mitigate harmful biases that arise when models rely on spurious
 correlations, thereby supporting more equitable AI systems.

F PROMPT TEMPLATES

We presented the prompt template for clear Z , X , Y separation in Figure 14 and detailed instruction for LLM-as-Judge in Figure 15.

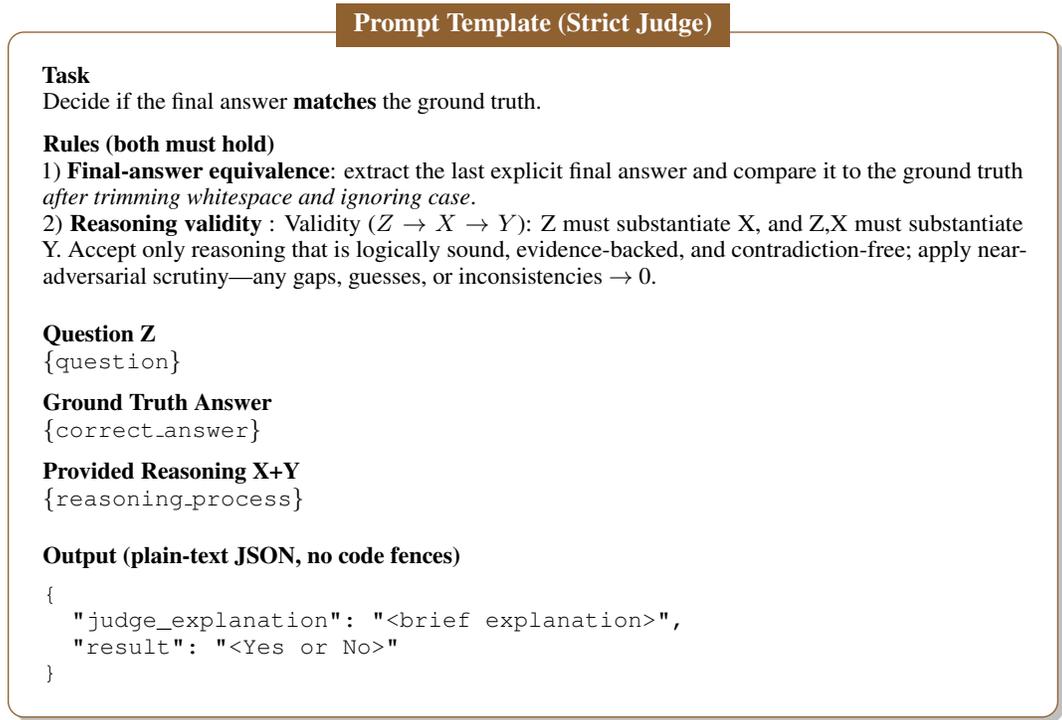
1134
1135
1136
1137
1138
1139
1140
1141
1142
1143
1144
1145
1146
1147
1148
1149



1150
1151

Figure 14: Verbatim instructions for final-answer tagging.

1152
1153



1176
1177
1178

Figure 15: Prompt template for strict binary judgment requiring answer equivalence and valid reasoning.

1179
1180

Table 5: Welch's t -tests comparing $S(A, B)$ distributions between Incorrect and Correct groups.

1182
1183
1184
1185
1186
1187

Metric	t	p -value
$S(Z, X)$	3.55	0.0005
$S(X, Y)$	3.23	0.0014
$S(Z, Y)$	3.15	0.0018

G DATASET FEATURES

Table 6: Training Dataset

Dataset	Question	Answer
BBEHCausal	Question: Alice (AF) and Bob (BF) each fire a bullet at a window, simultaneously striking the window. The window only shatters (WS) if it is hit by two bullets. Is Alice firing the bullet a sufficient cause for the window shattering? Reply based on the answer a logician would give.	No
CasHOLD	Context: limitations was equitably tolled fails because he did not raise this argument before the district court. See <i>Hinton v. Pac. Enters.</i> , 5 F.3d 391, 395 (9th Cir.1993) (HOLDING). Grace’s remaining contentions lack merit. Select the correct holding (1 to 4): “holding that the burden to allege facts sufficient to establish jurisdiction resides with plaintiff”, “holding that the plaintiff bears the burden when relying on the discovery rule”, “recognizing the validity of the doctrine but holding no equitable tolling on the facts presented”, “holding that plaintiff bears the burden to timely allege facts supporting equitable tolling”, “holding that the burden is on the plaintiff to allege facts sufficient to establish jurisdiction”	3: “holding that plaintiff bears the burden to timely allege facts supporting equitable tolling”
MATHHARD	There are thirty-five red, yellow, orange, and white marbles in a bag. If half the number of red marbles equals two less than the number of yellow marbles, equals a third the number of orange marbles, and equals a third of three more than the number of white marbles, how many red marbles are there?	8
IfQA	If Caroline Flack’s mother had felt the name suited her, what would have been Caroline’s name?	Caroline

H THE USE OF LARGE LANGUAGE MODELS (LLMs)

LLMs are the primary subject of our reinforcement learning (RL) experiments: we fine-tune pretrained models under CE-PO and baseline objectives to study stability, reward hacking, and causal alignment. Separately, we use LLMs as utilities for controlled text transformations—rephrasing, grammar correction, and during writing phase.

1242
1243
1244
1245
1246
1247
1248
1249
1250
1251
1252
1253
1254
1255
1256
1257
1258
1259
1260
1261
1262
1263
1264
1265
1266
1267
1268
1269
1270
1271
1272
1273
1274
1275
1276
1277
1278
1279
1280
1281
1282
1283
1284
1285
1286
1287
1288
1289
1290
1291
1292
1293
1294
1295

Table 7: Testing Dataset

Dataset	Question	Answer
BBEHMATH	Consider the operations $a \square b =$	-10038
	$\begin{cases} (a-b), & a+b > 0 \\ -a-b, & \text{otherwise} \end{cases}, \quad a \ ; \ b =$	
	$\begin{cases} (b-a) * b, & a-b < 2 \\ (a-b) * a, & \text{otherwise} \end{cases}, \quad a \ ! \ b =$	
	$\begin{cases} (2 \square b) * a, & a > b \\ b, & \text{otherwise} \end{cases}, \quad \text{and } a \ @ \ b =$	
	$\begin{cases} (a!b), & a-b > 0 \\ a * b, & \text{otherwise} \end{cases}; \quad \text{define } A = (((-5 ;$	
	$-7) + (2 @ -4) \square ((-9 \square -2) ! (6 \cdot 1)) \square$	
	$((7 - 6) \square (6 + -6)) \cdot ((-2 \square 9) ; (-1 - 5)),$	
	$B = (((9 ; 7) \square (6 \cdot 7)) \square ((2 ! -2) + (3 \square -6))) ;$	
	$((1 - 8) + (-5 @ -5)); ((-4 - 10); (-3 @ 6)), \text{ and}$	
	$C = (((-10! -1) ! (-5 + -6)) @ ((-8 + 5) \cdot (-10 +$	
	$10)); (((-3 - 3); (2 \square -5)) - ((2 \cdot 9) - (-8 \square 2)));$	
	compute $A + B - C$ (answer as a number).	
CLadder	Context: For husbands that don't set the alarm, the probability of ringing alarm is 42%. For husbands that set the alarm, the probability of ringing alarm is 51%. Question: For husbands that set the alarm, would it be less likely to see ringing alarm if the husband had not set the alarm?	yes
LegalBench	Context: Overall, the percentages and correlation coefficients in nine of the ten endogenous are sufficient for purposes of the second Gingles precondition. See 478 U.S. at 56. Plaintiffs have shown minority political cohesiveness in these [**105] appellate-judgeship elections. Nonetheless, the stark statistic remains that, in Alabama's history, only two African Americans have been elected to statewide office for a total of [**186] three general elections and three primary elections. At the statewide level, this factor weighs in favor of Plaintiffs. Question: Based on the excerpt above, did the judge's legal reasoning rely on statistical evidence (e.g., regression analysis) when determining whether there was a causal link? Answer Yes or No	Yes
LogiQA	Some purple clay pots are alive. Therefore, some living things have good or bad quality.Question: Which of the following judgments, if true, would best strengthen the argument?Options: ['The quality of purple clay teapots is different from good to bad.', 'Some purple clay pots are inanimate.', 'There is no difference in the quality of purple clay teapots.', 'Some living things are not purple clay pots.']	0, The quality of purple clay teapots is different from good to bad.