
Rethinking KL Regularization in RLHF: From Wrong Value Estimation to Correct Gradient Optimization

Anonymous Authors¹

Abstract

Reinforcement Learning from Human Feedback (RLHF) leverages a Kullback-Leibler (KL) divergence loss to stabilize training and prevent overfitting. However, in methods such as GRPO, its implementation may be guided by principles from numerical value estimation — a practice that overlooks the term’s functional role as an optimization loss. To analyze this issue, we establish a unified framework that connects two seemingly distinct implementation styles: using the mathematical term k_n as a detached coefficient for the policy’s score function (k_n in reward) or as a direct loss function through which gradients are propagated (k_n as loss). We show that the latter can always be analyzed via an equivalent gradient coefficient in the former, unifying the two perspectives. Through this framework, we first prove that conclusions from value estimation fail to guide proper KL loss design, using the k_1 as loss as a counterexample. We then prove the conventional k_1 in reward (like PPO) is the principled loss for Reverse KL (RKL) regularization. We further establish a key finding: under on-policy conditions, the k_2 as loss formulation is, in fact, gradient-equivalent to k_1 in reward. This equivalence, first proven in our work, identifies both as the theoretically sound implementations of the RKL objective. In contrast, we show that the recently adopted k_3 as loss (like GRPO) is merely a first-order, biased approximation of the principled loss. Furthermore, we argue that common off-policy implementations of k_n as loss methods are biased due to neglected importance sampling, and we propose a principled correction. Our findings provide a comprehensive, gradient-based rationale for choosing and correctly implementing KL regularization, paving the way for more robust and effective RLHF systems.

1. Introduction

Training state-of-the-art Large Language Models (LLMs) is a multi-stage process. Following large-scale pretraining and Supervised Fine-Tuning (SFT), a final post-training stage, Reinforcement Learning from Human Feedback (RLHF), is often employed to align models with preferences and to improve performance on demanding tasks such as mathematical reasoning (Ouyang et al., 2022; Shao et al., 2024). A core component of RLHF is **KL regularization**: a penalty derived from the Kullback–Leibler divergence (Kullback & Leibler, 1951) that constrains the policy to stay close to a reference model (typically the SFT checkpoint), thereby improving training stability and reducing reward overfitting (Ouyang et al., 2022; Stiennon et al., 2020).

Despite the critical role of the KL loss, its theoretical foundations in optimization remain underexplored. Across algorithms and codebases, the KL penalty takes several subtly different forms, and a common justification is how well a single-sample term estimates the KL *value* (e.g., bias/variance). However, in RLHF, the KL term serves as an *optimization loss*, not merely a diagnostic: because the sampling distribution depends on the policy parameters, good value estimation does not guarantee that the induced gradient matches the intended regularizer. As a result, implementations that look nearly identical in code—often differing only in whether a term is detached or backpropagated—can correspond to qualitatively different update directions and stability behavior.

This paper argues that a gradient-centric perspective is essential for designing robust and effective RLHF algorithms. We establish a unified framework that connects two common implementation styles: using the term k_n as a detached coefficient (k_n in reward) versus differentiating a surrogate penalty directly (k_n as loss). This framework allows us to compare implementations through the scalar coefficient they induce on the score function $\nabla_{\theta} \log \pi_{\theta}$.

Our main contributions are threefold:

1. **A gradient-correctness criterion for KL loss design.** We show that conclusions from KL value estimation do not necessarily translate into correct optimization be-

¹. Correspondence to: Anonymous Author.

havior. In particular, k_1 as loss provides a clean counterexample: despite k_1 being an unbiased KL value estimator, it has zero expected gradient under on-policy sampling and thus provides no KL constraint.

2. **Identifying principled RKL implementations.** We derive the exact on-policy Reverse KL (RKL) gradient and prove that k_1 in reward (PPO-style) and k_2 as loss are gradient-equivalent, establishing both as theoretically sound choices for RLHF regularization.
3. **Diagnosing GRPO-style surrogates and off-policy pitfalls.** We analyze k_3 as loss (used in GRPO) as a biased first-order surrogate with mismatched tail behavior. We further identify missing importance-sampling and clipping corrections as a key source of bias when “as loss” heads are reused off-policy, and provide a principled correction.

2. Related Work

Our work intersects three lines of research: KL value estimation, RLHF systems, and KL loss design. We briefly review each and clarify how our gradient-centric perspective complements prior contributions.

KL Value Estimation. The KL divergence is often estimated via Monte Carlo sampling, motivating a family of single-sample terms k_1, k_2, k_3 (John, 2020). Prior discussions evaluate them as *value estimators*, characterizing k_1 as unbiased but high-variance, k_2 as biased but lower-variance, and k_3 as an “optimal” low-variance unbiased estimator under regularity assumptions. We complement this view in two ways: we show that the advertised superiority of k_3 can fail under distribution shift and heavy-tailed importance ratios (Appendix I); and more importantly, we argue that value-estimation criteria are insufficient when these terms are used as *optimization losses*, where the induced gradient is what determines the regularization effect.

RLHF Systems and Methods. Open-source RLHF frameworks (e.g., OpenRLHF (Hu et al., 2024), Verl (Sheng et al., 2024), slime (Zhu et al., 2025), and ROLL (Wang et al., 2025)) support PPO-style updates (Ouyang et al., 2022) and more recent critic-free variants. They emphasize systems-level challenges such as efficient rollouts (e.g., OpenRLHF uses vLLM (Kwon et al., 2023)) and actor-critic stability. Several recent works focus on the critic bottleneck: VAPO (Yue et al., 2025) proposes critic pretraining, while GRPO and Reinforce++ remove the critic altogether to enable larger actor scaling. Most methods retain KL regularization, though some rule-based reward algorithms (e.g., DAPO) propose removing it for performance; ProRL (Liu et al., 2025b) mitigates this trade-off by periodically resetting the reference model, highlighting that KL remains

important for long-horizon stability. Our contribution is complementary: we analyze the gradient correctness of the KL regularizer itself, which impacts stability regardless of the surrounding loss.

KL Loss in RLHF. The practice of adding a KL penalty *in the reward* traces back to InstructGPT (Ouyang et al., 2022), where the log-ratio acts as a detached coefficient multiplying the policy score function. Jaques et al. (2019) noted a potential connection between adding KL to the reward versus to the loss, but without a formal gradient-level analysis. More recently, GRPO (Shao et al., 2024) (adopted in DeepSeek-R1 (Guo et al., 2025)) popularized directly differentiating through a k_3 -style surrogate, often justified by its value-estimation properties (John, 2020). Our work provides a unified framework to compare these choices at the gradient level, identifies when they match the intended Reverse KL objective, and discusses connections to Forward KL in Appendix K.

3. Preliminaries

Before analyzing KL implementations, we introduce the necessary background: the single-sample terms commonly used to approximate KL divergences, and the standard RLHF objective they are meant to regularize.

3.1. Single-Sample Terms for KL Estimation

The Kullback–Leibler (KL) divergence from a distribution $q(x)$ to a reference $p(x)$ is defined as

$$D_{\text{KL}}(q \parallel p) = \mathbb{E}_{x \sim q} \left[\log \frac{q(x)}{p(x)} \right]. \quad (1)$$

When this expectation is intractable, practitioners estimate it from Monte Carlo samples. Defining the likelihood ratio $\delta(x) = p(x)/q(x)$, common single-sample terms include (John, 2020):

$$\begin{aligned} k_1(x) &= -\log \delta(x), \\ k_2(x) &= \frac{1}{2} (\log \delta(x))^2, \\ k_3(x) &= \delta(x) - 1 - \log \delta(x). \end{aligned}$$

In RLHF, we instantiate these terms with $q(\cdot) = \pi_\theta(\cdot|x)$ (the current policy) and $p(\cdot) = \pi_{\text{ref}}(\cdot|x)$ (the reference). The term k_3 is designed to be nonnegative and to reduce variance when p and q are close; however, the “strictly better estimator” claim (John, 2020) relies on regularity conditions that can fail under distribution shift (see Appendix I for counterexamples). More fundamentally, when these terms are used for *regularization*, their suitability depends on the induced gradient, not on value-estimation properties—a distinction we formalize in the next section. In particular, when the support or tails of p and q differ significantly,

k_3 can exhibit heavy-tailed behavior: its sample mean can be dominated by rare events and its variance may become extremely large (or even infinite) without additional assumptions. Therefore, using k_3 as a “low-variance unbiased” value estimator requires verifying conditions that are often violated in practice; Appendix I provides counterexamples and discussion.

Gap in prior work. The blog (John, 2020) cited by GRPO to justify k_3 discusses only *value estimation*—not *loss design*. Even as a value estimator, k_3 can fail under distribution shift (Appendix I).

3.2. The RLHF Objective

RLHF fine-tunes a policy π_θ to produce responses y to prompts x that maximize a learned reward $r(x, y)$. Pure reward maximization can lead to reward hacking and drift from a trusted SFT policy. A standard remedy is to regularize toward a fixed reference π_{ref} via the Reverse KL (RKL):

$$\begin{aligned}
 \mathcal{J}_{\text{RLHF}}(\theta) &= \mathbb{E}_{x \sim \mathcal{D}} \mathbb{E}_{y \sim \pi_\theta(\cdot|x)} [r(x, y)] \\
 &\quad - \beta D_{\text{KL}}(\pi_\theta(\cdot|x) \parallel \pi_{\text{ref}}(\cdot|x)) \quad (2) \\
 &\triangleq \mathcal{J}_{\text{Reward}}(\theta) - \beta \mathcal{J}_{\text{RKL}}(\theta).
 \end{aligned}$$

Here \mathcal{D} is the prompt distribution, $\beta > 0$ controls the strength of regularization, and in practice y is sampled from a detached snapshot π_θ that is numerically equal to π_θ at sampling time.

The high-level objective is clear, but it leaves open an important implementation choice: *how should the KL term be integrated into the policy-gradient update?* Different choices lead to different induced gradients. We address this question next.

4. A Unified Framework for KL Regularization

Although most RLHF algorithms optimize the same high-level objective (Equation (2)), their concrete implementations of the KL term differ. This section introduces a unified framework that makes these differences explicit and comparable.

The key observation is that any practical KL implementation, regardless of its form, ultimately induces a scalar coefficient $c(x, y)$ multiplying the policy score function $\nabla_\theta \log \pi_\theta(y|x)$. This coefficient determines both the expected update direction and the variance of the stochastic gradient. By making $c(x, y)$ explicit, we can directly compare implementations: a KL loss is *gradient-correct* if and only if its induced coefficient matches that of the target KL objective.

Notation.

Key convention (red vs. black). We use θ (red) to denote trainable parameters that carry gradients, while θ (black) denotes a detached snapshot (no gradients).

Samples y are drawn from the detached snapshot $\pi_\theta(\cdot|x)$, which is numerically equal to $\pi_\theta(\cdot|x)$ at sampling time. Scalar coefficients multiplying the score function are always treated as detached.¹

4.1. Core Components of the RLHF Objective

The practical RLHF objective consists of two terms—reward maximization and KL regularization—both estimated via Monte Carlo sampling.

Reward Maximization. The primary goal is to maximize expected reward: $\mathcal{J}_{\text{Reward}}(\theta) = \mathbb{E}_{x \sim \mathcal{D}, y \sim \pi_\theta(\cdot|x)} [r(x, y)]$. Since sampling from π_θ is non-differentiable, the gradient is estimated via the log-derivative trick (REINFORCE (Williams, 1992); see Appendix B):

$$\nabla_\theta \mathcal{J}_{\text{Reward}}(\theta) = \mathbb{E}_{x \sim \mathcal{D}, y \sim \pi_\theta(\cdot|x)} [r(x, y) \cdot \nabla_\theta \log \pi_\theta(y|x)]. \quad (3)$$

Here $r(x, y)$ is typically a shaped advantage signal rather than a raw reward (details in Appendix A).

KL Regularization. The RKL term constrains the policy toward the reference:

$$\begin{aligned}
 \mathcal{J}_{\text{RKL}}(\theta) &= \mathbb{E}_{x \sim \mathcal{D}} [D_{\text{KL}}(\pi_\theta(\cdot|x) \parallel \pi_{\text{ref}}(\cdot|x))] \\
 &= \mathbb{E}_{x \sim \mathcal{D}, y \sim \pi_\theta(\cdot|x)} [\log \pi_\theta(y|x) - \log \pi_{\text{ref}}(y|x)]. \quad (4)
 \end{aligned}$$

Like the reward term, practical implementations estimate this expectation with samples $y \sim \pi_\theta(\cdot|x)$ and differentiate only through π_θ (see Appendix C). However, *how* the KL-derived term interacts with the gradient varies across algorithms. We distinguish two common styles:

- k_n in reward (detached coefficient):** The KL-derived scalar k_n multiplies the score function as a *detached* coefficient—gradients do not flow through k_n itself. This is the style used in PPO, with $k_1 = \log \pi_\theta - \log \pi_{\text{ref}}$:

$$\mathcal{J}_{k_n \text{ in reward}}(\theta) := \mathbb{E}_{x \sim \mathcal{D}, y \sim \pi_\theta(\cdot|x)} [k_n(\pi_\theta(y|x), \pi_{\text{ref}}(y|x)) \log \pi_\theta(y|x)]. \quad (5)$$

- k_n as loss (direct differentiation):** The KL-derived term is treated as a standalone loss, with gradients

¹We treat the entire response y as a single action, working with the joint probability $\pi(y|x)$ rather than token-level factors. This bandit-style simplification suffices for analyzing the core gradient properties.

propagated *through* it:

$$k_3 = \frac{\pi_{\text{ref}}}{\pi_{\theta}} - \log \frac{\pi_{\text{ref}}}{\pi_{\theta}} - 1 \quad (\text{GRPO}), \quad k_2 = \frac{1}{2} \left(\log \frac{\pi_{\theta}}{\pi_{\text{ref}}} \right)^2 \quad (\text{Ours}).$$

The corresponding objective is:

$$\mathcal{J}_{k_n \text{ as loss}}(\theta) := \mathbb{E}_{x \sim \mathcal{D}, y \sim \pi_{\theta}(\cdot|x)} [k_n(\pi_{\theta}(y|x), \pi_{\text{ref}}(y|x))]. \quad (6)$$

Both forms are *optimization surrogates*: their scalar values need not equal a KL divergence; what matters is the gradient they induce.

Connecting the two styles. Although k_n as loss differentiates through k_n directly, its gradient can always be written in the $k_{n'}$ in reward form for some equivalent coefficient $k_{n'}$. This coefficient is obtained by differentiating k_n with respect to the log-policy:

Equivalent coefficient formula:

$$k_{n'} := \left. \frac{\partial k_n}{\partial \log \pi_{\theta}} \right|_{\pi_{\theta} = \pi_{\theta}}. \quad (7)$$

This formula allows direct comparison of the two styles. A key example: $k_{2'} = k_1$. In Section 5.2, we prove that k_1 in reward and k_2 as loss are therefore gradient-equivalent.

4.2. Integration Forms and Algorithm Mapping

The KL formulation also determines how it is integrated with the reward objective in practice.

Combined vs. Decoupled Forms. Since k_n in reward shares the same score function as the reward term, its coefficient can be merged to form a single “advantage”:

$$\begin{aligned} \mathcal{L}_{\text{Combined}}(\theta) &= -\mathbb{E}_{x,y} [A_{\text{comb}}(x,y) \log \pi_{\theta}(y|x)], \\ A_{\text{comb}}(x,y) &:= r(x,y) - \beta k_n. \end{aligned} \quad (8)$$

In contrast, k_n as loss naturally yields a **decoupled** objective with a separate KL head:

$$\begin{aligned} \mathcal{L}_{\text{Decoupled}}(\theta) &= -\mathbb{E}_{x \sim \mathcal{D}, y \sim \pi_{\theta}} [r(x,y) \cdot \log \pi_{\theta}(y|x)] \\ &\quad + \beta \mathbb{E}_{x \sim \mathcal{D}, y \sim \pi_{\theta}} [k_n(\pi_{\theta}(y|x), \pi_{\text{ref}}(y|x))]. \end{aligned} \quad (9)$$

Off-policy considerations. In off-policy updates (e.g., PPO with multiple gradient steps per rollout), the combined form inherits importance-sampling (IS) and clipping corrections via the standard PPO surrogate $\pi_{\theta}/\pi_{\theta_k}$. In contrast, k_n as loss implementations must apply IS and clipping to the KL head itself—a step commonly omitted in practice, leading to biased gradients (see Appendix G).

Positioning PPO and GRPO. Table 1 summarizes where common algorithms fall in this taxonomy. PPO uses k_1 in reward in a combined form; GRPO uses k_3 as loss in a decoupled form.

With the framework in place, we now analyze which KL losses are gradient-correct for the RKL objective and diagnose failure modes when they are not.

5. Gradient-Based Analysis of KL Implementations

We now apply the coefficient framework to analyze common KL implementations. The analysis proceeds in three steps: (i) a counterexample showing that k_1 as loss provides no regularization signal; (ii) derivation of the principled RKL gradient and its equivalent surrogates; and (iii) diagnosis of k_3 as loss as a biased first-order approximation.

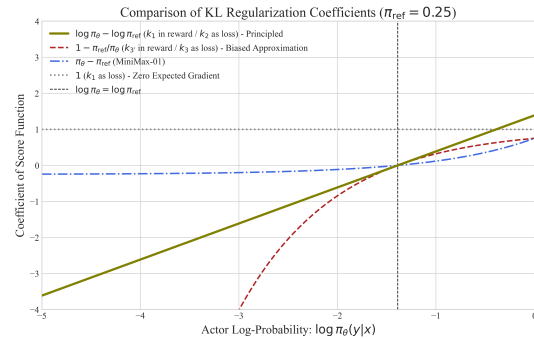


Figure 1. Induced score-function coefficients as a function of $\log \pi_{\theta}(y|x)$ (with $\pi_{\text{ref}} = 0.25$). The principled RKL implementations (k_1 in reward and k_2 as loss) yield coefficient k_1 . The GRPO-style k_3 as loss uses coefficient $1 - \delta$ (where $\delta = \pi_{\text{ref}}/\pi_{\theta}$), a first-order approximation with mismatched tails. The direct k_1 as loss yields a constant 1, producing zero-mean noise.

5.1. Counterexample: Why k_1 as loss Fails

We begin with a counterexample that illustrates the gap between value estimation and optimization. Although $k_1 = \log(\pi_{\theta}/\pi_{\text{ref}})$ is an unbiased estimator of the KL value, using it directly as a loss provides no regularization.

Consider the direct-loss formulation with on-policy sampling:

$$\mathcal{J}_{k_1 \text{ as loss}}(\theta) = \mathbb{E}_{x \sim \mathcal{D}, y \sim \pi_{\theta}(\cdot|x)} [\log \pi_{\theta}(y|x) - \log \pi_{\text{ref}}(y|x)]. \quad (10)$$

Since π_{ref} does not depend on θ , differentiating yields

$$\nabla_{\theta} \mathcal{J}_{k_1 \text{ as loss}}(\theta) = \mathbb{E}_{x \sim \mathcal{D}, y \sim \pi_{\theta}(\cdot|x)} [\nabla_{\theta} \log \pi_{\theta}(y|x)]. \quad (11)$$

The gradient is independent of π_{ref} —it carries *no* directional signal toward the reference. Moreover, by the zero-mean score identity (Lemma C.2), this expectation equals zero. In finite-sample (minibatch) practice, the gradient is

Table 1. Taxonomy of RLHF KL implementations. k_n in reward can also be used in a decoupled form.

Algorithm	Typical k_n	Style	Form	Off-Policy Notes
PPO / REINFORCE	k_1	in reward	Combined	IS/clipping inherited from PPO surrogate.
GRPO	k_3	as loss	Decoupled	Requires explicit IS/clipping; often omitted.

nonzero but has zero mean; k_1 as loss thus injects reference-independent noise that increases variance and can destabilize training, without providing any meaningful constraint.

This counterexample underscores a key point: good value estimation does not imply good optimization. A KL regularizer must be evaluated by its induced gradient, not its scalar estimate.

5.2. Principled RKL Surrogates: k_1 in reward $\Leftrightarrow k_2$ as loss

Having seen a failure mode, we now derive the gradient that a principled KL regularizer *should* produce. Applying the product rule and log-derivative trick to the RKL objective (Equation (4); see Appendix C) gives:

$$\nabla_{\theta} \mathcal{J}_{\text{RKL}}(\theta) = \mathbb{E}_{x \sim \mathcal{D}} \left[\sum_y \nabla_{\theta} \pi_{\theta}(y|x) \left(\log \frac{\pi_{\theta}(y|x)}{\pi_{\text{ref}}(y|x)} + 1 \right) \right]. \quad (12)$$

By the zero-mean score identity in Lemma C.2, the term '+1' vanishes in expectation, resulting in the practical form of the policy gradient:

Target RKL Gradient (the ‘gold standard’):

$$\nabla_{\theta} \mathcal{J}_{\text{RKL}}(\theta) = \mathbb{E}_{x \sim \mathcal{D}, y \sim \pi_{\theta}} \left[\underbrace{(\log \pi_{\theta} - \log \pi_{\text{ref}})}_{\text{coefficient } k_1} \underbrace{\nabla_{\theta} \log \pi_{\theta}}_{\text{score function}} \right] \quad (13)$$

Any surrogate that correctly implements RKL must reproduce this target gradient. The following theorem shows that two seemingly different designs achieve this.

Theorem 5.1 (Gradient equivalence of principled RKL surrogates). *Under on-policy sampling $y \sim \pi_{\theta}(\cdot|x)$, the following two objectives are **gradient-equivalent**—both produce Equation (13):*

$$\mathcal{J}_{k_1 \text{ in reward}} = \mathbb{E} \left[\underbrace{(\log \pi_{\theta} - \log \pi_{\text{ref}})}_{\text{coeff. } k_1 \text{ (no grad)}} \cdot \log \pi_{\theta} \right] \quad (14)$$

$$\mathcal{J}_{k_2 \text{ as loss}} = \mathbb{E} \left[\underbrace{\frac{1}{2} (\log \pi_{\theta} - \log \pi_{\text{ref}})^2}_{\text{gradients flow through}} \right] \quad (15)$$

Key insight. In Equation (14), k_1 is a fixed scalar—gradients flow only through $\log \pi_{\theta}$. In Equation (15), gradients flow through the squared term. Yet differentiating the

square produces exactly k_1 as the coefficient:

$$\nabla_{\theta} \frac{1}{2} (\log \frac{\pi_{\theta}}{\pi_{\text{ref}}})^2 = \underbrace{(\log \frac{\pi_{\theta}}{\pi_{\text{ref}}})}_{\text{same } k_1!} \cdot \nabla_{\theta} \log \pi_{\theta}. \quad (16)$$

Thus, both forms induce the same update: $k_1 \times$ score function. See Appendix C for the full proof.

Why prefer k_2 as loss? While gradient-equivalent to k_1 in reward, the squared form $k_2 = \frac{1}{2} (\log \pi_{\theta} / \pi_{\text{ref}})^2$ offers practical advantages: (i) it is *always non-negative*, avoiding sign cancellation that increases variance in k_1 ; (ii) its gradient coefficient k_1 emerges naturally from differentiation, making the implementation straightforward; and (iii) unlike k_3 , its variance does not depend on the chi-squared divergence and remains well-behaved under moderate distribution shift. These properties make k_2 as loss a robust default for on-policy RKL regularization.

Off-policy updates require additional IS and clipping corrections (Appendix G).

5.3. The GRPO Surrogate: k_3 as loss as a First-Order Approximation

We now turn to k_3 as loss, popularized by GRPO:

$$\mathcal{J}_{k_3 \text{ as loss}}(\theta) = \mathbb{E}_{x \sim \mathcal{D}, y \sim \pi_{\theta}(\cdot|x)} \left[\frac{\pi_{\text{ref}}}{\pi_{\theta}} - \log \frac{\pi_{\text{ref}}}{\pi_{\theta}} - 1 \right]. \quad (17)$$

Differentiating yields the induced coefficient $k_{3'}$:

$$\begin{aligned} \nabla_{\theta} \mathcal{J}_{k_3 \text{ as loss}}(\theta) &= \mathbb{E}_{x \sim \mathcal{D}, y \sim \pi_{\theta}(\cdot|x)} \left[\underbrace{\left(1 - \frac{\pi_{\text{ref}}(y|x)}{\pi_{\theta}(y|x)} \right)}_{k_{3'} \text{ coefficient}} \nabla_{\theta} \log \pi_{\theta}(y|x) \right] \\ &= \nabla_{\theta} \mathcal{J}_{k_{3'} \text{ in reward}}(\theta). \end{aligned} \quad (18)$$

Let $\delta = \pi_{\text{ref}} / \pi_{\theta}$. The principled RKL coefficient and the k_3 -induced coefficient are:

$$c^* = -\log \delta, \quad c_{3'} = 1 - \delta. \quad (19)$$

Around $\delta = 1$, Taylor expansion gives $-\log \delta \approx 1 - \delta + \mathcal{O}((\delta - 1)^2)$, so $c_{3'}$ is a first-order approximation of c^* . The mismatch beyond first order leads to three issues (see Appendix E for formal statements and Figure 1 for visualization):

1. **Bias.** For $\delta \neq 1$, the update direction differs from the true RKL gradient.
2. **Mismatched tails.** When $\pi_\theta > \pi_{\text{ref}}$ ($\delta \rightarrow 0$), the principled coefficient $-\log \delta \rightarrow +\infty$ provides a strong restoring force, whereas $1 - \delta \rightarrow 1$ saturates—weakening the constraint precisely when drift is large. Conversely, when $\pi_\theta < \pi_{\text{ref}}$ ($\delta \rightarrow \infty$), $1 - \delta \rightarrow -\infty$ diverges faster than $-\log \delta$, inducing potentially unbounded updates.
3. **Variance sensitivity.** $\text{Var}[1 - \delta] = \chi^2(\pi_{\text{ref}} \parallel \pi_\theta)$, the chi-square divergence, which can be unstable under distribution shift.

Appendix K provides an alternative view: $c_{3'} = 1 - \delta$ equals the Forward KL coefficient $-\delta$ plus an implicit baseline, making k_3 as loss a variance-reduced FKL estimator rather than an RKL surrogate.

5.4. Summary of Recommendations

Table 2 summarizes the key practical consequences of our analysis. The ‘‘Coefficient’’ column shows the scalar c multiplying the score function $\nabla_{\theta} \log \pi_{\theta}$; a method is RKL-correct if $c = k_1$. Based on our analysis, we offer the following practical guidance (see Appendix G and Appendix F for details on the last two points):

- **Avoid k_1 as loss.** Its expected gradient is zero and independent of the reference—it provides no regularization, only noise.
- **Use k_1 in reward or k_2 as loss for principled RKL regularization.** These are gradient-equivalent under on-policy sampling and correctly implement the RKL objective.
- **Recognize k_3 as loss as a biased surrogate.** It can impose weaker constraints when drift is large and may induce unbounded updates in the opposite tail.
- **Apply IS/clipping in off-policy settings.** When using k_n as loss with PPO-style updates, explicit importance sampling and clipping are required for the KL head; omitting them introduces systematic bias.
- **Consider bounded alternatives if stability is paramount.** The MSE-based penalty (e.g., MiniMax-01 loss) induces a coefficient bounded in $[-1, 1]$; see Appendix F.

We empirically validate these predictions in Section 6.

6. Experimental Validation

We validate our gradient analysis with controlled RLHF experiments on a mathematical reasoning task. The experiments address two questions derived from Section 5: (i)

Does k_1 as loss provide any regularization beyond noise? (ii) Does the principled k_2 as loss enforce a stronger constraint than its first-order surrogate k_3 as loss? We report 7B-scale results here; 1.5B-scale dynamics and downstream benchmarks appear in Appendix M and N.

6.1. Setup

Dataset. We use a curated subset of OpenR1-Math-220k,² comprising NuminaMath 1.5 prompts with reference solutions verified by Math-Verify.³ After removing sequences exceeding 2048 tokens, 7,300 prompts remain.

Training Configuration. To isolate gradient effects, we use fully on-policy training: rollout batch size 32, 8 responses per prompt, update batch size 256, and sampling temperature 1.0. The actor is Qwen2.5-Math-7B (Yang et al., 2024); rewards combine regex-based format checking and Math-Verify accuracy. We disable entropy regularization. Unless otherwise stated, we set the KL weight to $\beta = 0.5$; to better reflect practical settings with weaker constraints, we additionally report diagnostics at $\beta = 0.001$ (which permits larger drift and makes instabilities easier to observe). A numerical stability issue with group normalization and our mitigation are discussed in Appendix J.

6.2. Results

We track reward, accuracy, KL and log-probability gaps to the reference, group-level reward standard deviation, and response length.

Experiment 1: k_1 as loss vs. no KL. Figure 2 compares k_1 as loss to a no-KL baseline. As derived in Section 5, the gradient of k_1 as loss is independent of π_{ref} and has zero expectation; it should behave similarly to no KL, contributing only variance. Indeed, the two runs are comparable in reward/accuracy, but k_1 as loss exhibits larger KL and log-probability gaps at later stages, consistent with added stochasticity rather than regularization. Notably, both the no-KL and k_1 as loss settings reach higher reward than runs with effective KL constraints (Figure 3), reinforcing that k_1 as loss does not meaningfully limit drift.

Experiment 2: k_2 as loss vs. k_3 as loss. Figure 3 compares the principled k_2 as loss to the GRPO-style k_3 as loss. As a first-order approximation, k_3 as loss is asymmetric: it penalizes more strongly than k_2 as loss when $\pi_\theta < \pi_{\text{ref}}$ but saturates when $\pi_\theta > \pi_{\text{ref}}$. Since RL training typically drives the policy toward higher-reward regions

²<https://huggingface.co/datasets/open-r1/OpenR1-Math-220k>

³<https://github.com/huggingface/Math-Verify>

Table 2. KL regularization selection guide. Target coefficient for RKL: $c^* = k_1 = \log(\pi_\theta / \pi_{\text{ref}})$. Let $\delta = \pi_{\text{ref}} / \pi_\theta$.

Implementation	Coeff.	RKL?	Off-policy	Recommendation
In-reward form: $\mathcal{L} = \mathbb{E}[k_n \cdot \log \pi_\theta]$ (coefficient k_n does not carry gradients)				
k_1 in reward	k_1	✓	Inherits PPO clip	Recommended. Principled RKL.
k_2 in reward	k_2	✗	–	Avoid. Wrong coefficient.
k_3 in reward	k_3	✗	–	Avoid. Mixed RKL–FKL.
MiniMax-01	$\pi_\theta - \pi_{\text{ref}}$	–	Inherits PPO clip	Bounded $[-1, 1]$; MSE-based.
As-loss form: $\mathcal{L} = \mathbb{E}[k_n(\pi_\theta, \pi_{\text{ref}})]$ (gradients flow through k_n)				
k_1 as loss	1	✗	–	Avoid. Zero expected gradient.
k_2 as loss	k_1	✓	Needs IS/clip [†]	Recommended. Low-variance, stable.
k_3 as loss	$1 - \delta$	≈	Needs IS/clip [†]	Biased 1st-order; mismatched tails.
DeepSeek-V3.2	$\rho \cdot k_1$	✓*	Built-in IS	Equivalent to k_1 in reward; unclipped unstable.

✓ = RKL-correct, ✗ = incorrect, ≈ = 1st-order approx. *Full grad required; detaching $\rho \rightarrow 1 - \delta$. [†]App. G.

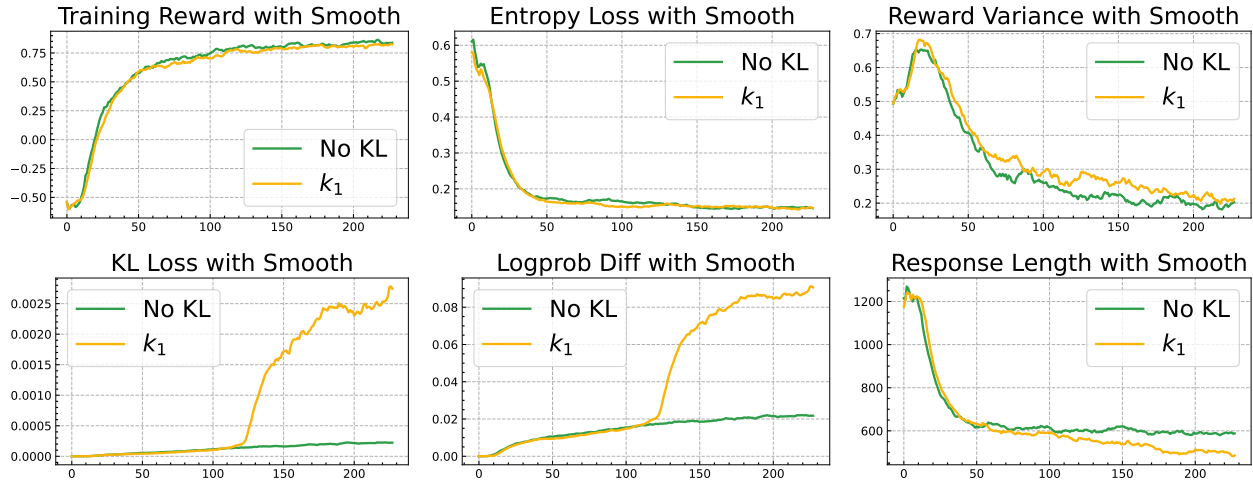


Figure 2. **Experiment 1:** 7B-scale comparison of k_1 as loss versus no KL. As predicted, k_1 as loss provides no effective regularization; it can increase stochasticity and lead to sharper drift (larger KL and log-prob gaps).

where $\pi_\theta > \pi_{\text{ref}}$, the constraint from k_3 as loss weakens in later stages. Indeed, both methods achieve comparable reward, but k_2 as loss maintains smaller log-probability gaps, indicating tighter coupling to the reference. Under a weaker and more realistic constraint ($\beta = 0.001$), k_3 as loss can show an unstable spike in the KL-loss curve (Figure 4), while k_2 as loss does not in the same setting.

Downstream Performance. Table 3 reports benchmark accuracy after training. Methods without effective KL constraints (no KL, k_1 as loss) achieve highest scores because they impose no limit on policy drift, while principled regularizers (k_2 as loss, k_3 as loss) trade some accuracy for stability—consistent with the intended role of KL regularization. Full results including 1.5B-scale and general reasoning benchmarks appear in Appendix N.

Summary. These results corroborate the theoretical analysis: k_1 as loss provides no regularization signal, while k_2 as loss enforces a more stable constraint than k_3 as loss.

Table 3. Downstream math accuracy (7B, $\beta = 0.5$). Methods without effective KL regularization allow unconstrained drift, yielding higher scores but less stable training.

Method	AIME	AMC	MATH	Minerva	Oly.	Avg.
Qwen2.5-Math-7B	8.2	31.3	43.6	7.4	15.6	19.0
+ RL w/o KL	17.5	55.6	78.6	36.8	42.4	41.4
+ k_1 as loss	15.4	56.0	80.6	40.8	43.0	41.8
+ k_2 as loss	11.5	48.5	64.2	16.9	24.9	29.6
+ k_3 as loss	13.2	48.9	65.4	18.8	29.0	31.4

Case Study: DeepSeek-V3.2. DeepSeek-V3.2 (Liu et al., 2025a) appears to recognize the bias in vanilla k_3 as loss and proposes a corrected IS-weighted variant:

$$\mathcal{L}_{\text{KL,DS}}(\theta) = \frac{\pi_\theta}{\pi_{\text{old}}} \cdot k_3(\theta). \quad (20)$$

Notably, the ratio $\rho(\theta) = \pi_\theta / \pi_{\text{old}}$ is *gradient-carrying*—unlike standard IS where the ratio is detached. Through our framework, we show this design choice is deliberate: the product-rule expansion $\nabla(\rho \cdot k_3)$ generates two terms

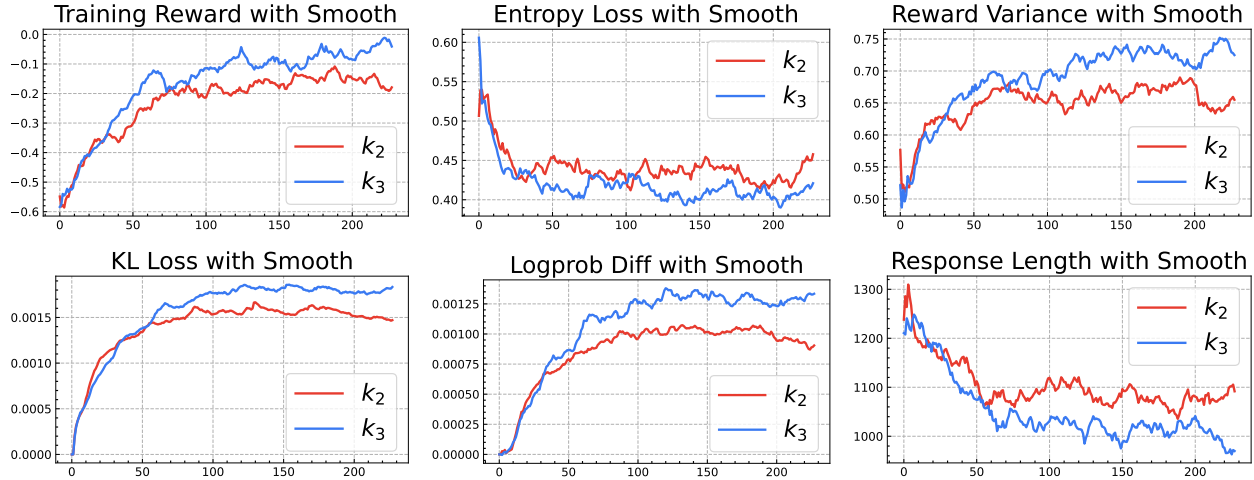


Figure 3. **Experiment 2:** 7B-scale comparison of k_2 as loss (principled RKL) and k_3 as loss (first-order surrogate). Both constrain the policy, but k_2 as loss maintains tighter coupling to the reference and yields more stable training dynamics.

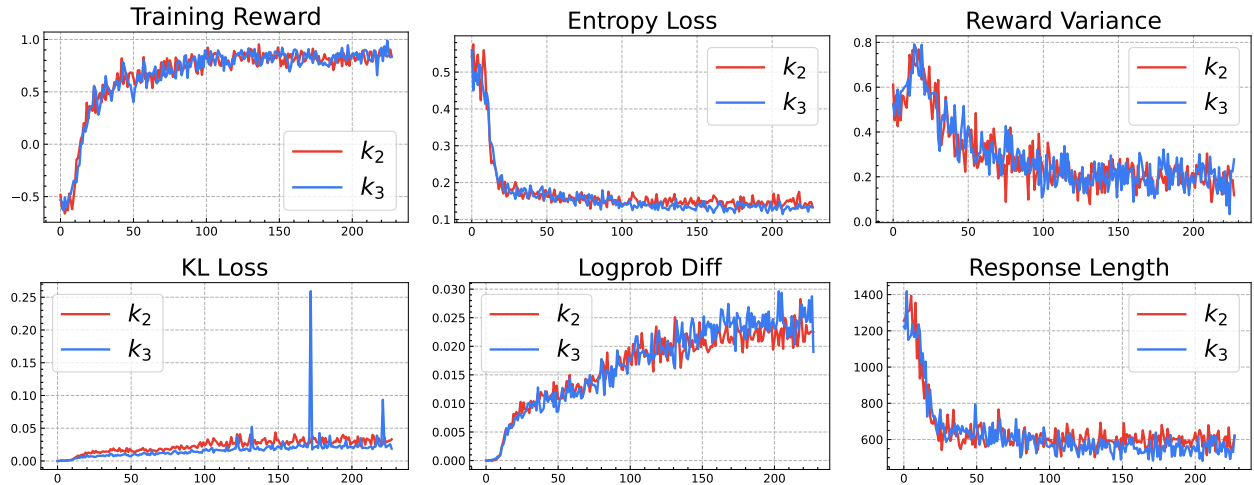


Figure 4. 7B-scale training diagnostics for k_2 as loss and k_3 as loss at a weaker KL weight ($\beta = 0.001$). While most metrics track closely, the **KL Loss** panel (bottom-left) reveals a visible spike for k_3 as loss (reaching ~ 0.25) that is absent for k_2 as loss, demonstrating the k_3 surrogate’s instability when drift is larger. See Figure 7 for 1.5B-scale results showing the same pattern.

whose biases cancel, recovering the principled k_1 coefficient. However, this creates a practical tension: preserving the full gradient can yield large updates when ρ is large, while detaching ρ loses the cancellation and reduces to the biased k_3 direction. Appendix L provides the full derivation.

7. Conclusion

This paper develops a gradient-centric framework for KL regularization in RLHF, expressing each implementation as a scalar coefficient on the policy score function. This perspective yields a simple correctness criterion and three main findings: (i) k_1 as loss has zero expected gradient despite being an unbiased value estimator; (ii) k_1 in reward and k_2 as loss are gradient-equivalent and correctly

implement Reverse KL; (iii) k_3 as loss is a biased first-order surrogate with mismatched tail behavior. We also identify missing importance-sampling corrections in off-policy “as loss” implementations.

Experiments at 7B scale validate these predictions, with k_2 as loss showing greater stability than k_3 as loss under weak regularization. We recommend k_1 in reward or k_2 as loss for principled RLHF regularization.

Impact Statement

This work is methodological and aims to improve the stability and reproducibility of RLHF training by clarifying the effects of KL regularization on gradients. Better-grounded KL implementations can reduce training instability and wasted compute, and may lower the risk of unintended behaviors caused by uncontrolled policy drift during alignment. At the same time, improvements to RLHF optimization can be used in both beneficial and harmful applications; we encourage practitioners to evaluate downstream impacts and follow responsible deployment practices.

References

- Clark, P., Cowhey, I., Etzioni, O., Khot, T., Sabharwal, A., Schoenick, C., and Tafjord, O. Think you have solved question answering? try arc, the ai2 reasoning challenge. *arXiv preprint arXiv:1803.05457*, 2018.
- Guo, D., Yang, D., Zhang, H., Song, J., Zhang, R., Xu, R., Zhu, Q., Ma, S., Wang, P., Bi, X., et al. Deepseek-r1: Incentivizing reasoning capability in llms via reinforcement learning. *arXiv preprint arXiv:2501.12948*, 2025.
- Hendrycks, D., Burns, C., Kadavath, S., Arora, A., Basart, S., Tang, E., Song, D., and Steinhardt, J. Measuring mathematical problem solving with the math dataset. *arXiv preprint arXiv:2103.03874*, 2021.
- Hu, J., Wu, X., Zhu, Z., Wang, W., Zhang, D., Cao, Y., et al. Openrlhf: An easy-to-use, scalable and high-performance rlhf framework. *arXiv preprint arXiv:2405.11143*, 2024.
- Jaques, N., Ghandeharioun, A., Shen, J. H., Ferguson, C., Lapedriza, A., Jones, N., Gu, S., and Picard, R. Way off-policy batch deep reinforcement learning of implicit human preferences in dialog. *arXiv preprint arXiv:1907.00456*, 2019.
- John, S. Approximating kl divergence, 2020. URL <http://joschu.net/blog/kl-approx.html>.
- Kullback, S. and Leibler, R. A. On information and sufficiency. *The annals of mathematical statistics*, 22(1): 79–86, 1951.
- Kwon, W., Li, Z., Zhuang, S., Sheng, Y., Zheng, L., Yu, C. H., Gonzalez, J. E., Zhang, H., and Stoica, I. Efficient memory management for large language model serving with pagedattention. In *Proceedings of the ACM SIGOPS 29th Symposium on Operating Systems Principles*, 2023.
- Li, A., Gong, B., Yang, B., Shan, B., Liu, C., Zhu, C., Zhang, C., Guo, C., Chen, D., Li, D., et al. Minimax-01: Scaling foundation models with lightning attention. *arXiv preprint arXiv:2501.08313*, 2025.
- Li, J., Beeching, E., Tunstall, L., Lipkin, B., Soletskyi, R., Huang, S., Rasul, K., Yu, L., Jiang, A. Q., Shen, Z., et al. Numnamath: The largest public dataset in ai4maths with 860k pairs of competition math problems and solutions. *Hugging Face repository*, 13(9):9, 2024.
- Liu, A., Mei, A., Lin, B., Xue, B., Wang, B., Xu, B., Wu, B., Zhang, B., Lin, C., Dong, C., et al. Deepseek-v3. 2: Pushing the frontier of open large language models. *arXiv preprint arXiv:2512.02556*, 2025a.
- Liu, M., Diao, S., Lu, X., Hu, J., Dong, X., Choi, Y., Kautz, J., and Dong, Y. Prorl: Prolonged reinforcement learning expands reasoning boundaries in large language models. *arXiv preprint arXiv:2505.24864*, 2025b.
- Liu, Z., Chen, C., Li, W., Qi, P., Pang, T., Du, C., Lee, W. S., and Lin, M. Understanding r1-zero-like training: A critical perspective. *arXiv preprint arXiv:2503.20783*, 2025c.
- Ouyang, L., Wu, J., Jiang, X., Almeida, D., Wainwright, C., Mishkin, P., Zhang, C., Agarwal, S., Slama, K., Ray, A., et al. Training language models to follow instructions with human feedback. *Advances in neural information processing systems*, 35:27730–27744, 2022.
- Rein, D., Hou, B. L., Stickland, A. C., Petty, J., Pang, R. Y., Dirani, J., Michael, J., and Bowman, S. R. Gpqa: A graduate-level google-proof q&a benchmark. In *First Conference on Language Modeling*, 2024.
- Shao, Z., Wang, P., Zhu, Q., Xu, R., Song, J., Bi, X., Zhang, H., Zhang, M., Li, Y., Wu, Y., et al. Deepseekmath: Pushing the limits of mathematical reasoning in open language models. *arXiv preprint arXiv:2402.03300*, 2024.
- Sheng, G., Zhang, C., Ye, Z., Wu, X., Zhang, W., Zhang, R., Peng, Y., Lin, H., and Wu, C. Hybridflow: A flexible and efficient rlhf framework. *arXiv preprint arXiv:2409.19256*, 2024.
- Stiennon, N., Ouyang, L., Wu, J., Ziegler, D., Lowe, R., Voss, C., Radford, A., Amodei, D., and Christiano, P. F. Learning to summarize with human feedback. *Advances in neural information processing systems*, 33:3008–3021, 2020.
- Tang, Y. and Munos, R. On a few pitfalls in kl divergence gradient estimation for rl. *arXiv preprint arXiv:2506.09477*, 2025.
- Wang, W., Xiong, S., Chen, G., Gao, W., Guo, S., He, Y., Huang, J., Liu, J., Li, Z., Li, X., et al. Reinforcement learning optimization for large-scale learning: An efficient and user-friendly scaling library. *arXiv preprint arXiv:2506.06122*, 2025.

- 495 Wang, Y., Ma, X., Zhang, G., Ni, Y., Chandra, A., Guo, S.,
496 Ren, W., Arulraj, A., He, X., Jiang, Z., et al. Mmlu-pro:
497 A more robust and challenging multi-task language un-
498 derstanding benchmark. *Advances in Neural Information*
499 *Processing Systems*, 37:95266–95290, 2024.
- 500 Williams, R. J. Simple statistical gradient-following algo-
501 rithms for connectionist reinforcement learning. *Machine*
502 *learning*, 8(3):229–256, 1992.
- 503 Yang, A., Zhang, B., Hui, B., Gao, B., Yu, B., Li, C., Liu,
504 D., Tu, J., Zhou, J., Lin, J., et al. Qwen2. 5-math techni-
505 cal report: Toward mathematical expert model via self-
506 improvement. *arXiv preprint arXiv:2409.12122*, 2024.
- 507 Yue, Y., Yuan, Y., Yu, Q., Zuo, X., Zhu, R., Xu, W., Chen,
508 J., Wang, C., Fan, T., Du, Z., et al. Vapo: Efficient and
509 reliable reinforcement learning for advanced reasoning
510 tasks. *arXiv preprint arXiv:2504.05118*, 2025.
- 511 Zhu, Z., Xie, C., Lv, X., and slime Contributors. slime: An
512 llm post-training framework for rl scaling. [https://](https://github.com/THUDM/slime)
513 github.com/THUDM/slime, 2025. GitHub reposi-
514 tory. Corresponding author: Xin Lv.
- 515
- 516
- 517
- 518
- 519
- 520
- 521
- 522
- 523
- 524
- 525
- 526
- 527
- 528
- 529
- 530
- 531
- 532
- 533
- 534
- 535
- 536
- 537
- 538
- 539
- 540
- 541
- 542
- 543
- 544
- 545
- 546
- 547
- 548
- 549

Appendix Overview

This appendix provides comprehensive supporting material for the theoretical framework and empirical analysis presented in the main text. To help readers navigate efficiently, we organize the content into four thematic groups.

I. Foundations and Notation (Section A–Section B). These sections establish the practical context and mathematical foundations that underpin our analysis.

- **Section A:** Implementation details and reward shaping in RLHF—how minibatches are structured, how rewards are normalized, and how the KL term is integrated.
- **Section B:** Policy-gradient derivation for reward maximization—clarifying the distinction between objectives, gradients, and surrogate losses.

II. Core Theoretical Results (Section C–Section E). These sections contain the formal proofs supporting the main claims.

- **Section C:** Full proof that k_1 in reward and k_2 as loss are gradient-equivalent under on-policy sampling—the “gold standard” for RKL regularization.
- **Section D:** Surrogate objectives and Monte Carlo estimation—connecting population-level theory to practical minibatch implementations.
- **Section E:** Formal analysis of the k_3 as loss surrogate—proving its first-order bias, tail asymmetry, and variance sensitivity.

III. Extensions and Alternative Methods (Section F–Section L). These sections extend the gradient-centric framework to related settings and alternative designs.

- **Section F:** Derivation of the bounded MiniMax-01 regularizer as an MSE-based alternative with bounded coefficients.
- **Section G:** Off-policy corrections for KL regularization—principled importance sampling and clipping when reusing samples.
- **Section H:** Python code for visualizing the coefficient functions in Figure 1.
- **Section I:** Statistical instability of the k_3 value estimator—conditions under which its variance becomes infinite.
- **Section J:** Group normalization stability issues and a proposed fix.
- **Section K:** Forward KL vs. Reverse KL—reinterpreting k_3 as loss as a variance-reduced FKL estimator.
- **Section L:** Analysis of the DeepSeek-V3.2 importance-weighted KL estimator and its trade-offs.

IV. Empirical Supplements (Section M–Section N). These sections provide additional experimental evidence.

- **Section M:** 1.5B-scale experiments complementing the 7B-scale results in the main text.
- **Section N:** Downstream benchmark performance on math and general reasoning tasks.

A. Detailed Implementation of RLHF Methods

The main text introduces a unified framework that analyzes KL implementations through their induced gradient coefficients. Before diving into the theoretical analysis, it is helpful to establish the practical context: how are RLHF algorithms actually implemented, and where do the different KL formulations fit in?

This section provides that context. We describe the standard minibatch structure, the reward shaping techniques used in practice, and the two canonical ways to integrate KL regularization—the “combined form” (k_n in reward) and the “decoupled form” (k_n as loss). These details will be referenced throughout the theoretical sections that follow.

Minibatch structure. In a typical training step, we draw N prompts $\{x^{(i)}\}_{i=1}^N$ from the dataset \mathcal{D} . For each prompt $x^{(i)}$, we sample G responses $y^{(i,1)}, \dots, y^{(i,G)} \sim \pi_\theta(\cdot | x^{(i)})$ from the detached policy snapshot. The G responses for each prompt form a *group*; the full minibatch contains $N \times G$ prompt–response pairs. We use π_θ for the trainable policy (gradients flow through it) and π_θ for its detached snapshot.

Baseline subtraction and normalization. Having established the minibatch structure, we now describe how the raw reward signal is processed before being used in the policy gradient. A key technique is baseline subtraction: subtracting an action-independent baseline does not change the expected policy gradient (see Equation (36)) and typically reduces variance (Williams, 1992). We consider the following operators applied to a scalar signal r :

$$f_{\text{group-bl}}(r) = r - \text{mean}_{\text{group}}(r), \quad (21)$$

$$f_{\text{batch-bl}}(r) = r - \text{mean}_{\text{batch}}(r). \quad (22)$$

In practice, normalization is also used. Normalization is *not* unbiased but can improve numerical stability by controlling the scale of the reward signal:

$$f_{\text{BN}}(r) = \frac{r - \text{mean}_{\text{batch}}(r)}{\text{std}_{\text{batch}}(r)}, \quad (23)$$

$$f_{\text{GN}}(r) = \frac{r - \text{mean}_{\text{group}}(r)}{\text{std}_{\text{group}}(r)}. \quad (24)$$

How common algorithms shape the reward signal. With these operators defined, we can now specify how different RLHF algorithms process the raw reward signal. Let $r_{\text{raw}}(x, y)$ denote the raw score from the reward model. Different algorithms post-process this score as follows:

$$\text{REINFORCE: } r(x, y) = f_{\text{BN}}(r_{\text{raw}}(x, y)), \quad (25)$$

$$\text{PPO/REINFORCE++: } r(x, y) = f_{\text{BN}}(r_{\text{raw}}(x, y)), \quad (26)$$

$$\text{GRPO: } r(x, y) = f_{\text{GN}}(r_{\text{raw}}(x, y)), \quad (27)$$

$$\text{Dr-GRPO (Liu et al., 2025c): } r(x, y) = f_{\text{group-bl}}(r_{\text{raw}}(x, y)), \quad (28)$$

$$\text{REINFORCE++baseline: } r(x, y) = f_{\text{BN}}(f_{\text{group-bl}}(r_{\text{raw}}(x, y))). \quad (29)$$

Integration of the KL regularization loss. The transforms above shape only the reward signal. A separate and critical design choice is how to integrate the KL regularizer. As discussed in the main text (Section 4), there are two canonical approaches:

- (i) **Combined form (k_n in reward):** Used by REINFORCE/PPO methods. A combined reward signal is formed first,

$$A_{\text{combined}}(x, y) = r_{\text{raw}}(x, y) - \beta k_1(\pi_\theta(y | x), \pi_{\text{ref}}(y | x)),$$

and then baseline/normalization is applied to this combined signal A_{combined} before it multiplies the score function.

- (ii) **Decoupled form (k_n as loss):** Used by GRPO methods. The KL penalty $k_n(\pi_\theta(\cdot | x), \pi_{\text{ref}}(\cdot | x))$ is optimized as a separate, unnormalized loss term, added to the policy-gradient loss driven by the shaped reward $r(x, y)$.

Complete on-policy objectives. Combining the reward shaping and KL integration described above, we can now write down the complete on-policy objectives for representative algorithms. These expressions make explicit how the two KL styles manifest in practice.

REINFORCE/PPO (Monte Carlo minibatch):

$$\mathcal{L}_{\text{REINFORCE/PPO,MC}}(\theta) = -\frac{1}{NG} \sum_{i=1}^N \sum_{j=1}^G \left\{ A(x^{(i)}, y^{(i,j)}) \log \pi_\theta(y^{(i,j)} | x^{(i)}) \right\}, \quad (30)$$

$$\text{where } A(x, y) = f_{\text{BN}}(r_{\text{raw}}(x, y) - \beta k_1(\pi_\theta(y | x), \pi_{\text{ref}}(y | x))). \quad (31)$$

Here, the k_1 term is evaluated at the current policy snapshot (no gradients flow through it), consistent with the policy-gradient framework where it acts as a coefficient for the score function.

GRPO (Monte Carlo minibatch):

$$\begin{aligned} \mathcal{L}_{\text{GRPO,MC}}(\theta) = & -\frac{1}{NG} \sum_{i=1}^N \sum_{j=1}^G \left\{ r(x^{(i)}, y^{(i,j)}) \log \pi_{\theta}(y^{(i,j)} | x^{(i)}) \right\} \\ & + \frac{\beta}{NG} \sum_{i=1}^N \sum_{j=1}^G k_3 \left(\pi_{\theta}(y^{(i,j)} | x^{(i)}), \pi_{\text{ref}}(y^{(i,j)} | x^{(i)}) \right). \end{aligned} \quad (32)$$

In this decoupled form, the shaped reward $r(x, y)$ drives the policy-gradient term, while the KL penalty is a separate loss where gradients flow directly through π_{θ} inside $k_3(\cdot)$.

Remark. While baseline subtraction is an unbiased variance-reduction technique, normalization is a biased but often crucial heuristic for practical stability. Both are important engineering details, even if omitted from simplified theoretical analyses.

B. Policy Gradient Derivation for Reward Maximization

The gradient-centric framework developed in the main text relies on expressing any KL implementation as a coefficient multiplying the policy score function $\nabla_{\theta} \log \pi_{\theta}(y|x)$. Before analyzing KL regularization specifically, it is instructive to review the standard policy gradient derivation for reward maximization, as the same techniques—the log-derivative trick and the distinction between trainable and detached parameters—apply throughout.

This section derives the policy gradient for the reward-maximization objective, clarifying the distinction between the true objective, its gradient, and surrogate loss functions. We adopt the same notation convention as the main text: θ denotes parameters that carry gradients, while θ denotes detached parameters (e.g., in the sampling distribution).

Objective. The goal is to maximize expected reward:

$$\mathcal{J}_{\text{reward}}(\theta) = \mathbb{E}_{x \sim \mathcal{D}, y \sim \pi_{\theta}(\cdot|x)} [r(x, y)] = \mathbb{E}_{x \sim \mathcal{D}} \sum_y [r(x, y) \cdot \pi_{\theta}(y|x)]. \quad (33)$$

We assume standard regularity conditions that permit the interchange of differentiation and expectation operators.

Policy gradient derivation. We compute the gradient of the objective function $\mathcal{J}_{\text{reward}}(\theta)$ using the log-derivative trick. The distinction between θ (trainable) and θ (detached) is crucial: the former appears wherever gradients must flow, while the latter marks quantities that are held fixed during differentiation.

$$\begin{aligned} \nabla_{\theta} \mathcal{J}_{\text{reward}}(\theta) &= \nabla_{\theta} \mathbb{E}_{x \sim \mathcal{D}} \sum_y [r(x, y) \cdot \pi_{\theta}(y|x)] \\ &= \mathbb{E}_{x \sim \mathcal{D}} \sum_y [r(x, y) \cdot \nabla_{\theta} \pi_{\theta}(y|x)] \\ &= \mathbb{E}_{x \sim \mathcal{D}} \sum_y \left[r(x, y) \cdot \pi_{\theta}(y|x) \cdot \frac{\nabla_{\theta} \pi_{\theta}(y|x)}{\pi_{\theta}(y|x)} \right] \\ &= \mathbb{E}_{x \sim \mathcal{D}} \sum_y \pi_{\theta}(y|x) [r(x, y) \cdot \nabla_{\theta} \log \pi_{\theta}(y|x)] \\ &= \mathbb{E}_{x \sim \mathcal{D}, y \sim \pi_{\theta}(\cdot|x)} [r(x, y) \cdot \nabla_{\theta} \log \pi_{\theta}(y|x)]. \end{aligned} \quad (34)$$

The key step occurs in the third line, where we multiply and divide by $\pi_{\theta}(y|x)$. This serves two purposes: (i) it converts the sum over y into an expectation under π_{θ} , enabling Monte Carlo estimation; and (ii) it introduces the log-derivative identity $\nabla_{\theta} \log \pi_{\theta} = \nabla_{\theta} \pi_{\theta} / \pi_{\theta}$. In the final line, the gradient is expressed as an expectation over samples from π_{θ} , where the sampling process itself has no gradient path—only the score function $\nabla_{\theta} \log \pi_{\theta}$ carries gradients. This structure, where a scalar coefficient (without gradients) multiplies the score function, is the foundation for the unified framework in Section 4.

C. Formal Proof of the Principled KL Regularization

This section provides the complete proof of Theorem 5.1 from the main text—the central theoretical result establishing that k_1 in reward and k_2 as loss are gradient-equivalent implementations of the Reverse KL (RKL) regularizer.

The proof proceeds in four stages. First, we state the regularity assumptions required for the analysis (Section C.1). Second, we establish two fundamental identities—the log-derivative identity and the zero-mean score property—that underpin all policy-gradient derivations (Section C.2). Third, we derive the exact gradient of the true RKL objective (Section C.3). Finally, we show that both surrogates reproduce this target gradient (Section C.4).

Throughout, we use $\pi_{\theta}(\cdot|x)$ for the gradient-carrying policy and $\pi_{\theta}(\cdot|x)$ for its detached snapshot (numerically identical at the current iterate). Sampling measures and scalar coefficients multiplying the score function are always treated as detached.

C.1. Assumptions and Notation

We begin by stating the regularity assumptions that ensure all subsequent derivations are well-defined. Let \mathcal{D} be a data distribution over prompts x , $\pi_{\text{ref}}(\cdot|x)$ a fixed reference policy, and $\pi_{\theta}(\cdot|x)$ a differentiable policy parameterized by θ . All logarithms are natural.

(A1) (Valid policy) For each x , the function $y \mapsto \pi_{\theta}(y|x)$ is a valid probability mass/density: $\pi_{\theta}(y|x) > 0$ on its support, it is differentiable in θ , and normalizes to one, i.e., $\sum_y \pi_{\theta}(y|x) = 1$ (or $\int \pi_{\theta}(y|x) dy = 1$).

(A2) (Interchange) The interchange of expectation/summation and differentiation is valid (standard regularity conditions).

(A3) (Independence) The data distribution \mathcal{D} and reference policy π_{ref} do not depend on the trainable parameters θ .

(A4) (Support) The KL divergence is well-defined: for all x and all y in the support of $\pi_{\theta}(\cdot|x)$, we have $\pi_{\text{ref}}(y|x) > 0$.

Notation convention: Throughout this section, π_{θ} (black θ) denotes the *detached* copy of π_{θ} (red θ), numerically equal at the current iterate but with no gradient path. Unless stated otherwise, expectations over y are taken with respect to the detached sampling distribution $y \sim \pi_{\theta}(\cdot|x)$.

C.2. Fundamental Identities

Before deriving the RKL gradient, we establish two lemmas and two corollaries that will be used repeatedly. These identities formalize the key properties of the score function that make policy-gradient methods tractable.

Lemma C.1 (Log-derivative identity with detached denominator). *For any fixed x and any y with $\pi_{\theta}(y|x) > 0$,*

$$\nabla_{\theta} \log \pi_{\theta}(y|x) = \frac{\nabla_{\theta} \pi_{\theta}(y|x)}{\pi_{\theta}(y|x)}. \quad (35)$$

Proof. By the chain rule, $\nabla_{\theta} \log \pi_{\theta} = (\nabla_{\theta} \pi_{\theta}) / \pi_{\theta}$. Replacing the denominator with its detached, numerically identical copy π_{θ} preserves the numerical value while making the no-gradient path explicit. \square \square

Lemma C.2 (Zero-mean score). *For any fixed x ,*

$$\mathbb{E}_{y \sim \pi_{\theta}(\cdot|x)} [\nabla_{\theta} \log \pi_{\theta}(y|x)] = 0. \quad (36)$$

Proof. Using Lemma C.1,

$$\sum_y \pi_{\theta}(y|x) \frac{\nabla_{\theta} \pi_{\theta}(y|x)}{\pi_{\theta}(y|x)} = \sum_y \nabla_{\theta} \pi_{\theta}(y|x) = \nabla_{\theta} \sum_y \pi_{\theta}(y|x) = \nabla_{\theta} (1) = 0. \quad (37)$$

Corollary C.0.1 (Score-function reweighting). *For any function $z(y, x)$ that does not depend on θ ,*

$$\sum_y \nabla_{\theta} \pi_{\theta}(y|x) z(y, x) = \mathbb{E}_{y \sim \pi_{\theta}(\cdot|x)} [z(y, x) \nabla_{\theta} \log \pi_{\theta}(y|x)]. \quad (38)$$

770 *Proof.* Multiply and divide by $\pi_\theta(y|x)$, then apply Lemma C.1:

$$771 \sum_y \nabla_\theta \pi_\theta(y|x) z(y, x) = \sum_y \pi_\theta(y|x) \frac{\nabla_\theta \pi_\theta(y|x)}{\pi_\theta(y|x)} z(y, x) = \sum_y \pi_\theta(y|x) z(y, x) \nabla_\theta \log \pi_\theta(y|x).$$

772 The right-hand side is precisely the expectation $\mathbb{E}_{y \sim \pi_\theta(\cdot|x)}[z(y, x) \nabla_\theta \log \pi_\theta(y|x)]$. \square

773 **Corollary C.0.2** (Baseline invariance). *For any function $b(x)$ that does not depend on θ ,*

$$774 \mathbb{E}_{y \sim \pi_\theta(\cdot|x)}[b(x) \nabla_\theta \log \pi_\theta(y|x)] = b(x) \cdot \mathbb{E}_{y \sim \pi_\theta(\cdot|x)}[\nabla_\theta \log \pi_\theta(y|x)] = 0. \quad (39)$$

775 Thus, adding an action-independent baseline $b(x)$ (which does not carry gradients) to any coefficient does not change the expected gradient.

776 C.3. Derivation of the True RKL Gradient

777 With the fundamental identities in hand, we now derive the gradient of the true RKL objective. This gradient serves as the “target” against which we will compare the two surrogate formulations. The RKL divergence objective is:

$$778 \mathcal{J}_{\text{RKL}}(\theta) = \mathbb{E}_{x \sim \mathcal{D}} \left[\sum_y \pi_\theta(y|x) \log \frac{\pi_\theta(y|x)}{\pi_{\text{ref}}(y|x)} \right]. \quad (40)$$

779 **Step 1 (Differentiate under the expectation).** By (A2), the gradient operator is moved inside the expectation and sum.

780 **Step 2 (Apply product rule).** For each y , we differentiate the term $\pi_\theta(y|x) \log \pi_\theta(y|x)$ using the product rule. Let $f = \pi_\theta$ and $g = \log \pi_\theta$. Then:

$$781 \begin{aligned} \nabla_\theta [\pi_\theta \log \pi_\theta] &= (\nabla_\theta \pi_\theta) \log \pi_\theta + \pi_\theta \cdot \nabla_\theta \log \pi_\theta \\ &= (\nabla_\theta \pi_\theta) \log \pi_\theta + \pi_\theta \cdot \frac{\nabla_\theta \pi_\theta}{\pi_\theta} \quad (\text{by Lemma C.1}) \\ &= (\nabla_\theta \pi_\theta) (\log \pi_\theta + 1). \end{aligned} \quad (41)$$

782 Here, the factors $\log \pi_\theta$ and π_θ that are not being differentiated are written as their detached copies to emphasize the no-gradient path. For the reference term, by (A3) π_{ref} does not depend on θ , so:

$$783 \nabla_\theta [-\pi_\theta(y|x) \log \pi_{\text{ref}}(y|x)] = -(\nabla_\theta \pi_\theta(y|x)) \log \pi_{\text{ref}}(y|x). \quad (42)$$

784 **Step 3 (Collect terms).** Combining the two terms from Step 2, we obtain:

$$785 \nabla_\theta \mathcal{J}_{\text{RKL}}(\theta) = \mathbb{E}_{x \sim \mathcal{D}} \left[\sum_y \nabla_\theta \pi_\theta(y|x) \left(\log \frac{\pi_\theta(y|x)}{\pi_{\text{ref}}(y|x)} + 1 \right) \right]. \quad (43)$$

786 At this stage, the gradient is expressed as a sum over y , weighted by $\nabla_\theta \pi_\theta(y|x)$. To convert this into a form amenable to Monte Carlo estimation, we apply Corollary C.0.1.

787 **Step 4 (Apply score-function reweighting).** Setting the scalar coefficient $z(y, x) := \log \frac{\pi_\theta(y|x)}{\pi_{\text{ref}}(y|x)} + 1$ (which does not carry gradients) and applying Corollary C.0.1:

$$788 \nabla_\theta \mathcal{J}_{\text{RKL}}(\theta) = \mathbb{E}_{x \sim \mathcal{D}, y \sim \pi_\theta(\cdot|x)} \left[\left(\log \frac{\pi_\theta(y|x)}{\pi_{\text{ref}}(y|x)} + 1 \right) \nabla_\theta \log \pi_\theta(y|x) \right]. \quad (44)$$

789 **Step 5 (Simplify using the zero-mean score property).** The coefficient in Equation (44) contains a constant $+1$ term. By Lemma C.2, this term contributes zero to the expected gradient:

$$790 \mathbb{E}_{y \sim \pi_\theta(\cdot|x)} [1 \cdot \nabla_\theta \log \pi_\theta(y|x)] = 0.$$

791 Dropping this zero-contribution term yields the final, simplified gradient:

$$792 \nabla_\theta \mathcal{J}_{\text{RKL}}(\theta) = \mathbb{E}_{x \sim \mathcal{D}, y \sim \pi_\theta(\cdot|x)} \left[\log \frac{\pi_\theta(y|x)}{\pi_{\text{ref}}(y|x)} \nabla_\theta \log \pi_\theta(y|x) \right]. \quad (45)$$

C.4. The Gold Standard: k_1 in reward $\Leftrightarrow k_2$ as loss

Having derived the true RKL gradient in Equation (45), we now show that two structurally different surrogates— k_1 in reward and k_2 as loss—both reproduce this exact gradient under on-policy sampling. This equivalence is the core theoretical result that validates both as principled implementations of RKL regularization.

Surrogate 1: k_1 in reward. The first surrogate treats the log-ratio as a scalar coefficient multiplying the score function:

$$\mathcal{J}_{k_1 \text{ in reward}}(\theta) = \mathbb{E}_{x \sim \mathcal{D}, y \sim \pi_\theta(\cdot|x)} \left[\underbrace{\left(\log \frac{\pi_\theta(y|x)}{\pi_{\text{ref}}(y|x)} \right)}_{\text{coefficient } k_1} \log \pi_\theta(y|x) \right]. \quad (46)$$

Since the coefficient k_1 does not carry gradients (it is evaluated at the current snapshot π_θ), differentiating yields:

$$\nabla_\theta \mathcal{J}_{k_1 \text{ in reward}}(\theta) = \mathbb{E}_{x \sim \mathcal{D}, y \sim \pi_\theta(\cdot|x)} \left[\left(\log \frac{\pi_\theta(y|x)}{\pi_{\text{ref}}(y|x)} \right) \nabla_\theta \log \pi_\theta(y|x) \right], \quad (47)$$

which is identical to Equation (45). This is the style used in PPO.

Surrogate 2: k_2 as loss. The second surrogate differentiates *through* the log-ratio via a squared penalty:

$$\mathcal{J}_{k_2 \text{ as loss}}(\theta) = \mathbb{E}_{x \sim \mathcal{D}, y \sim \pi_\theta(\cdot|x)} \left[\frac{1}{2} \left(\log \pi_\theta(y|x) - \log \pi_{\text{ref}}(y|x) \right)^2 \right]. \quad (48)$$

Applying the chain rule to differentiate the squared term:

$$\nabla_\theta \left[\frac{1}{2} \left(\log \pi_\theta - \log \pi_{\text{ref}} \right)^2 \right] = \left(\log \pi_\theta - \log \pi_{\text{ref}} \right) \cdot \nabla_\theta \log \pi_\theta. \quad (49)$$

Evaluating the scalar multiplier at the current policy snapshot π_θ and taking the expectation:

$$\nabla_\theta \mathcal{J}_{k_2 \text{ as loss}}(\theta) = \mathbb{E}_{x \sim \mathcal{D}, y \sim \pi_\theta(\cdot|x)} \left[\left(\log \pi_\theta(y|x) - \log \pi_{\text{ref}}(y|x) \right) \nabla_\theta \log \pi_\theta(y|x) \right], \quad (50)$$

which is also identical to Equation (45). Thus, despite their different functional forms, both surrogates induce the same expected gradient—the principled RKL gradient.

C.5. Conclusion and Implementation Guidance

We have now established the central result: under Assumptions (A1)–(A4), the true RKL objective and both surrogates— k_1 in reward and k_2 as loss—share the same expected gradient, given by Equation (45). This equivalence rests on the following key conventions, which practitioners should adhere to for correct implementation.

Sampling Measure: Samples are drawn from the current policy snapshot $y \sim \pi_\theta(\cdot|x)$ (e.g., via an external rollout engine).

Scalar Coefficients: The scale coefficients that multiply the score function do not carry gradients—they are evaluated at the current snapshot π_θ . Applying Corollary C.0.2, any action-independent baseline $b(x)$ can be added to reduce variance.

Gradient Path: Gradients propagate only through terms explicitly parameterized by θ .

Implementation Notes. We now translate these theoretical results into concrete implementation patterns. For a single on-policy sample $y \sim \pi_\theta(\cdot|x)$:

- **k_1 in reward:** The surrogate loss is

$$\mathcal{L}_{k_1 \text{ in reward}} = \left(\log \pi_\theta(y|x) - \log \pi_{\text{ref}}(y|x) \right) \cdot \log \pi_\theta(y|x),$$

where the coefficient $k_1 = \log \pi_\theta - \log \pi_{\text{ref}}$ is evaluated at the current snapshot (stop gradient) and only $\log \pi_\theta$ carries gradients. Gradient descent on $\mathcal{L}_{k_1 \text{ in reward}}$ performs descent on \mathcal{J}_{RKL} . In typical RLHF, where we maximize $\mathcal{J}_{\text{Reward}} - \beta \mathcal{J}_{\text{RKL}}$, the combined loss to minimize is $-(r(x, y) - \beta k_1) \cdot \log \pi_\theta(y|x)$.

- k_2 as loss: The surrogate loss is

$$\mathcal{L}_{k_2 \text{ as loss}} = \frac{1}{2} \left(\log \pi_{\theta}(y|x) - \log \pi_{\text{ref}}(y|x) \right)^2,$$

where the entire log-ratio carries gradients through π_{θ} . Gradient descent on $\mathcal{L}_{k_2 \text{ as loss}}$ also performs descent on \mathcal{J}_{RKL} , since differentiating yields the same coefficient $k_1 = \log \pi_{\theta} - \log \pi_{\text{ref}}$ multiplying $\nabla_{\theta} \log \pi_{\theta}$.

Remarks. (i) For continuous spaces, replace sums by integrals; the proof is unchanged provided densities are positive on their support. (ii) The equivalence requires on-policy sampling. If samples are drawn from a stale policy π_{old} , exact correction uses importance weights $\rho(x, y) = \pi_{\theta}(y|x)/\pi_{\text{old}}(y|x)$ inside the expectations.

D. Surrogate Objective: Full-Vocabulary vs. Monte Carlo

The preceding sections derive policy gradients as expectations over the full action space. In practice, however, we cannot enumerate all possible responses y for a given prompt x —the vocabulary is astronomically large for language models. This section bridges the gap between theory and practice by formalizing the connection between population-level objectives and their minibatch Monte Carlo estimators.

We detail two conceptually distinct implementations: a **full-vocabulary loss**, which computes the exact inner expectation over all actions (theoretically exact but computationally infeasible), and a **Monte Carlo (MC) loss**, which replaces this sum with i.i.d. samples (unbiased and practical). Understanding this connection is essential for verifying that practical implementations are gradient-correct.

Conventions and gradient paths. To ensure consistency with the main text, we restate the key conventions governing gradient flow:

1. The trainable policy π_{θ} carries gradients; its numerically identical snapshot at the current iterate (which does not carry gradients) is denoted π_{θ} .
2. All scalars that multiply the score function do not carry gradients: the reward $r(x, y)$, any KL-derived term $k_n(\cdot)$, and their combination $r(x, y) - \beta k_n(\cdot)$.
3. Gradients flow only through $\log \pi_{\theta}(y|x)$; everything inside the coefficient $c(x, y)$ does not carry gradients.
4. We adopt the naming from the main text: “reward coefficient,” “ k_n coefficient,” and “combined form coefficient.”

We express the objective using a generic scalar coefficient $c(x, y)$, which can take several forms:

$$c(x, y) \in \left\{ \begin{array}{ll} r(x, y) & \text{reward coefficient} \\ k_n(\pi_{\theta}(y|x), \pi_{\text{ref}}(y|x)) & k_n \text{ coefficient} \\ r(x, y) - \beta k_n(\pi_{\theta}(y|x), \pi_{\text{ref}}(y|x)) & \text{combined form coefficient} \end{array} \right\}. \quad (51)$$

A typical KL choice is

$$k_1(y|x) = \log \pi_{\theta}(y|x) - \log \pi_{\text{ref}}(y|x). \quad (52)$$

Policy gradient in expectation form (with baseline). For the population objective $\mathcal{J}_{\text{true}}(\theta) = \mathbb{E}_{x \sim \mathcal{D}, y \sim \pi_{\theta}(\cdot|x)}[c(x, y)]$, using an action-independent baseline $b(x)$ and the log-derivative identity (where the denominator is evaluated at the current snapshot),

$$\nabla_{\theta} \log \pi_{\theta}(y|x) = \frac{\nabla_{\theta} \pi_{\theta}(y|x)}{\pi_{\theta}(y|x)}, \quad (53)$$

the unbiased policy gradient is

$$\nabla_{\theta} \mathcal{J}_{\text{true}}(\theta) = \mathbb{E}_{x \sim \mathcal{D}, y \sim \pi_{\theta}(\cdot|x)} \left[(c(x, y) - b(x)) \nabla_{\theta} \log \pi_{\theta}(y|x) \right]. \quad (54)$$

This relies on the zero-mean score property under $y \sim \pi_{\theta}(\cdot|x)$ as proved in Equation (36), ensuring $b(x)$ does not change the expected gradient.

Population surrogate loss. A surrogate loss whose negative gradient recovers Equation (54) is

$$\mathcal{L}_{\text{sur}}(\theta) = -\mathbb{E}_{x \sim \mathcal{D}, y \sim \pi_{\theta}(\cdot|x)} \left[(c(x, y) - b(x)) \log \pi_{\theta}(y|x) \right]. \quad (55)$$

Two interchangeable minibatch implementations. We now provide two equivalent minibatch estimators of Equation (55). The first computes the exact inner expectation over the discrete action space \mathcal{V} by summing all actions with a sampling weight—this is the “full-vocabulary” version. The second replaces this inner sum with i.i.d. on-policy samples, yielding an unbiased estimate conditional on the minibatch prompts—this is the “Monte Carlo” version used in practice.

$$\mathcal{L}_{\text{sur,Full}}(\theta) = -\frac{1}{N} \sum_{i=1}^N \sum_{y^{(i)} \in \mathcal{V}} \underbrace{\pi_{\theta}(y^{(i)} | x^{(i)})}_{\text{sampling weight}} \left(c(x^{(i)}, y^{(i)}) - b(x^{(i)}) \right) \log \pi_{\theta}(y^{(i)} | x^{(i)}). \quad (56)$$

Here, the sampling weight π_{θ} does not carry gradients; the gradient path is solely via $\log \pi_{\theta}$.

$$\mathcal{L}_{\text{sur,MC}}(\theta) = -\frac{1}{N} \sum_{i=1}^N \frac{1}{G} \sum_{j=1}^G \left(c(x^{(i)}, y^{(i,j)}) - b(x^{(i)}) \right) \log \pi_{\theta}(y^{(i,j)} | x^{(i)}), \quad y^{(i,j)} \sim \pi_{\theta}(\cdot|x^{(i)}). \quad (57)$$

In Equation (57), $\{y^{(i,j)}\}_{j=1}^G$ are i.i.d. samples from the current policy snapshot $\pi_{\theta}(\cdot|x^{(i)})$; increasing G reduces variance while preserving unbiasedness.

Unbiasedness and practical considerations. For any fixed $x^{(i)}$ and function f ,

$$\mathbb{E}_{\{y^{(i,j)}\}_{j=1}^G \text{ i.i.d.} \sim \pi_{\theta}(\cdot|x^{(i)})} \left[\frac{1}{G} \sum_{j=1}^G f(y^{(i,j)}) \right] = \sum_{y^{(i)} \in \mathcal{V}} \pi_{\theta}(y^{(i)} | x^{(i)}) f(y^{(i)}), \quad (58)$$

hence $\mathbb{E}[\mathcal{L}_{\text{sur,MC}} | \{x^{(i)}\}] = \mathcal{L}_{\text{sur,Full}}$. In practice, computing $r(x, y)$ or $k_n(\cdot)$ over the full vocabulary is infeasible for LLMs due to GPU memory constraints; MC estimation is therefore standard.

Alternative decoupled formulation: k_n as loss. In addition to incorporating KL via the coefficient $c(x, y)$, one may add a separate penalty-only loss that differentiates directly through the log-ratio. Let $\psi_n : \mathbb{R} \rightarrow \mathbb{R}$ be differentiable (e.g., $\psi_1(t) = t, \psi_2(t) = \frac{1}{2}t^2$). Define

$$\mathcal{L}_{k_n \text{ as loss,MC}}(\theta) = -\frac{1}{NG} \sum_{i=1}^N \sum_{j=1}^G \psi_n \left(\underbrace{\log \pi_{\theta}(y^{(i,j)} | x^{(i)})}_{\text{with grad}} - \underbrace{\log \pi_{\text{ref}}(y^{(i,j)} | x^{(i)})}_{\text{fixed}} \right). \quad (59)$$

Its gradient takes the score-like form

$$\nabla_{\theta} \mathcal{L}_{k_n \text{ as loss,MC}}(\theta) = -\frac{1}{NG} \sum_{i=1}^N \sum_{j=1}^G \psi'_n(\cdot) \nabla_{\theta} \log \pi_{\theta}(y^{(i,j)} | x^{(i)}), \quad (60)$$

where (\cdot) denotes the log-ratio in Equation (59). Here the prime denotes differentiation with respect to the scalar argument:

$$\psi'_n(t) \triangleq \frac{d}{dt} \psi_n(t). \quad (61)$$

Under on-policy sampling, evaluating the scalar coefficient $\psi'_n(\cdot)$ at the current iterate yields the gradient-equivalence established in the main text (see Theorem in Section 5.2 and Section C). A common choice $\psi_2(t) = \frac{1}{2}t^2$ recovers the squared log-ratio penalty. The total objective is then $\mathcal{L}_{\text{reward,MC}} - \beta \cdot \mathcal{L}_{k_n \text{ as loss,MC}}$, with π_{ref} fixed (no gradients flow through it).

E. Formal Analysis of the k_3 as loss Gradient Surrogate

The preceding section established that k_1 in reward and k_2 as loss are gradient-correct implementations of the RKL regularizer. A natural question is whether other formulations—particularly k_3 as loss, which is widely used in GRPO—also implement the same gradient. This section provides the formal analysis showing that they do not.

Specifically, we prove that k_3 as loss is a first-order, biased surrogate for the principled RKL gradient, and we rigorously characterize its three core deficiencies:

1. **Local bias:** The induced coefficient differs from the true RKL coefficient except at $\pi_\theta = \pi_{\text{ref}}$.
2. **Tail asymmetry:** The surrogate saturates when the policy over-covers the reference and diverges when it under-covers.
3. **Statistical instability:** The variance of the induced coefficient equals the chi-squared divergence, which can be unbounded.

Setup and notation. Throughout this section, we fix a prompt x and consider on-policy samples $y \sim \pi_\theta(\cdot|x)$. To simplify notation, we define the probability ratio:

$$\delta(y) := \frac{\pi_{\text{ref}}(y|x)}{\pi_\theta(y|x)}. \quad (62)$$

The ratio δ captures how much the reference distribution favors action y relative to the current policy: $\delta > 1$ means the reference assigns more probability to y than the current policy (under-coverage), while $\delta < 1$ means the opposite (over-coverage). Our analysis compares the coefficient induced by k_3 as loss against the principled RKL gradient coefficient, $c^*(y) = -\log \delta(y)$. All scalar coefficients multiplying the score function $\nabla_\theta \log \pi_\theta(y|x)$ are treated as detached.

Gradient-Equivalent Coefficient of k_3 as loss. The k_3 as loss objective is given by:

$$\mathcal{J}_{k_3 \text{ as loss}}(\theta) = \mathbb{E}_{x \sim \mathcal{D}, y \sim \pi_\theta(\cdot|x)} \left[\frac{\pi_{\text{ref}}(y|x)}{\pi_\theta(y|x)} - 1 - \log \frac{\pi_{\text{ref}}(y|x)}{\pi_\theta(y|x)} \right]. \quad (63)$$

Differentiating and evaluating the resulting scalar multiplier at the detached snapshot yields:

$$\nabla_\theta \mathcal{J}_{k_3 \text{ as loss}}(\theta) = \mathbb{E}_{x \sim \mathcal{D}, y \sim \pi_\theta(\cdot|x)} \left[\left(1 - \frac{\pi_{\text{ref}}(y|x)}{\pi_\theta(y|x)} \right) \nabla_\theta \log \pi_\theta(y|x) \right]. \quad (64)$$

This confirms that k_3 as loss is gradient-equivalent (under on-policy sampling) to an ‘in-reward’ update with the following scalar coefficient:

$$c_{3'}(y) := 1 - \delta(y). \quad (65)$$

$$\mathcal{J}_{k_{3'} \text{ in reward}}(\theta) = \mathbb{E}_{x \sim \mathcal{D}, y \sim \pi_\theta(\cdot|x)} \left[\left(1 - \frac{\pi_{\text{ref}}(y|x)}{\pi_\theta(y|x)} \right) \log \pi_\theta(y|x) \right]. \quad (66)$$

The remainder of this section formally analyzes the deficiencies of this proxy, $c_{3'}$, when compared to the principled target, $c^* = -\log \delta$.

Lemma E.1 (First-order agreement and second-order bias). *The proxy $1 - \delta$ is the first-order Taylor approximation of the principled coefficient $-\log \delta$ around $\delta = 1$. The approximation error (bias) is of second order:*

$$\text{Bias}(\delta) = (-\log \delta) - (1 - \delta) = \frac{1}{2}(\delta - 1)^2 - \frac{1}{3}(\delta - 1)^3 + O((\delta - 1)^4). \quad (67)$$

Proof sketch. The result is obtained by expanding $-\log \delta$ in a Taylor series at $\delta = 1$ and subtracting the term $(1 - \delta)$. \square

Lemma E.2 (One-sided domination and asymmetric tails). *For all $\delta > 0$, the proxy is a strict lower bound, $1 - \delta \leq -\log \delta$, with equality holding only at $\delta = 1$. Their tail behaviors are pathologically asymmetric:*

- **Over-coverage** ($\delta \rightarrow 0^+$): *The proxy provides a weak, saturating restoring force ($\lim_{\delta \rightarrow 0^+} (1 - \delta) = 1$), whereas the principled coefficient provides an unbounded penalty ($\lim_{\delta \rightarrow 0^+} (-\log \delta) = +\infty$).*

- **Under-coverage** ($\delta \rightarrow \infty$): The proxy induces an unbounded linear penalty ($\lim_{\delta \rightarrow \infty} (1 - \delta) = -\infty$), while the principled coefficient grows only logarithmically.

Proof sketch. The inequality follows from the fundamental property $\log \delta \leq \delta - 1$. The limits are elementary. \square

Theorem E.1 (Variance equals chi-squared divergence). Assuming $\text{supp}(\pi_{\text{ref}}(\cdot|x)) \subseteq \text{supp}(\pi_{\theta}(\cdot|x))$, the proxy coefficient $c_{3'}$ has zero mean under the on-policy sampling distribution, and its variance is exactly the chi-squared divergence:

$$\begin{aligned} \mathbb{E}_{y \sim \pi_{\theta}} [1 - \delta(y)] &= 0, \\ \text{Var}_{y \sim \pi_{\theta}} [1 - \delta(y)] &= \chi^2(\pi_{\text{ref}}(\cdot|x) \parallel \pi_{\theta}(\cdot|x)). \end{aligned} \quad (68)$$

If the support condition is violated, the variance is infinite.

Proof sketch. $\mathbb{E}[\delta] = \sum_y \pi_{\theta}(y|x) \frac{\pi_{\text{ref}}(y|x)}{\pi_{\theta}(y|x)} = 1$, thus $\mathbb{E}[1 - \delta] = 0$. The variance identity then follows directly from the definition of $\chi^2(p \parallel q)$. \square

Corollary E.1.1 (Implication for stochastic gradient variance). The variance of the stochastic gradient term induced by k_3 as loss is directly governed by the chi-squared divergence, an often unstable metric:

$$\mathbb{E} \left[\|(1 - \delta(y)) \nabla_{\theta} \log \pi_{\theta}(y|x)\|^2 \right] = \mathbb{E} \left[(1 - \delta(y))^2 \|\nabla_{\theta} \log \pi_{\theta}(y|x)\|^2 \right]. \quad (69)$$

Summary. These results provide a rigorous, gradient-centric justification for the claims in the main text:

- The k_3 as loss formulation does *not* implement the true RKL gradient.
- It deploys a first-order proxy ($c_{3'} = 1 - \delta$) that is accurate only when the policy is very close to the reference ($\delta \approx 1$).
- Its pathological tail behavior—saturating at +1 in one direction while diverging to $-\infty$ in the other—introduces optimization challenges not present in the principled k_1 in reward or k_2 as loss formulations.
- Its variance equals the chi-squared divergence, which can be unbounded under distribution shift.

This analysis underscores the importance of selecting regularization losses based on their gradient properties, not merely their characteristics as value estimators.

F. Derivation of an Alternative Regularizer: The MiniMax-01 Loss

The preceding analysis revealed that KL-based coefficients can be unbounded: the principled RKL coefficient $\log(\pi_{\theta}/\pi_{\text{ref}})$ grows without bound as the policy diverges from the reference. While this provides strong regularization, it can also induce large gradient updates that destabilize training in some settings.

This section derives the MiniMax-01 loss (Li et al., 2025) as an alternative regularizer with a *bounded* coefficient. We show that it originates from a mean squared error (MSE) objective in probability space and fits naturally within our gradient-centric framework. The key advantage is that the induced coefficient $\pi_{\theta} - \pi_{\text{ref}}$ is always bounded in $[-1, 1]$, providing a more conservative regularization force.

Objective. We minimize the full-vocabulary MSE between the policy and the reference:

$$\mathcal{J}_{\text{MSE}}(\theta) = \mathbb{E}_{x \sim \mathcal{D}} \left[\frac{1}{2} \sum_y \left(\pi_{\theta}(y|x) - \pi_{\text{ref}}(y|x) \right)^2 \right]. \quad (70)$$

On-policy gradient. Differentiating Equation (70) with respect to θ , evaluating the scalar multiplier at the detached snapshot π_θ (our standard on-policy convention), and converting to the score-function form yields:

$$\begin{aligned}
 \nabla_{\theta} \mathcal{J}_{\text{MSE}}(\theta) &= \mathbb{E}_{x \sim \mathcal{D}} \sum_y \left[\frac{1}{2} \nabla_{\theta} (\pi_{\theta}(y|x) - \pi_{\text{ref}}(y|x))^2 \right] \\
 &= \mathbb{E}_{x \sim \mathcal{D}} \sum_y \underbrace{(\pi_{\theta}(y|x) - \pi_{\text{ref}}(y|x))}_{\text{scalar coefficient}} \nabla_{\theta} \pi_{\theta}(y|x) \\
 &= \mathbb{E}_{x \sim \mathcal{D}} \sum_y \pi_{\theta}(y|x) (\pi_{\theta}(y|x) - \pi_{\text{ref}}(y|x)) \nabla_{\theta} \log \pi_{\theta}(y|x) \\
 &= \mathbb{E}_{x \sim \mathcal{D}, y \sim \pi_{\theta}(\cdot|x)} \left[(\pi_{\theta}(y|x) - \pi_{\text{ref}}(y|x)) \nabla_{\theta} \log \pi_{\theta}(y|x) \right].
 \end{aligned} \tag{71}$$

The last line reveals that MSE regularization induces a score-function update whose scalar coefficient is the probability difference $\pi_\theta - \pi_{\text{ref}}$.

MiniMax-01 surrogate loss (Monte Carlo). Using a on-policy sampler that draws G responses $y^{(i,j)} \sim \pi_{\theta}(\cdot | x^{(i)})$ per prompt, the unbiased minibatch surrogate whose negative gradient recovers Equation (71) is

$$\mathcal{L}_{\text{MSE,MC(MiniMax-01)}}(\theta) = -\frac{1}{NG} \sum_{i=1}^N \sum_{j=1}^G \left(\pi_{\theta}(y^{(i,j)} | x^{(i)}) - \pi_{\text{ref}}(y^{(i,j)} | x^{(i)}) \right) \log \pi_{\theta}(y^{(i,j)} | x^{(i)}). \tag{72}$$

This head shares the same in-reward score-function structure as our principled KL implementations: the coefficient does not carry gradients, and only $\log \pi_{\theta}$ contributes to the gradient.

Key properties and implications.

1. **Bounded gradient coefficient.** Since $0 \leq \pi_{\theta}(y|x), \pi_{\text{ref}}(y|x) \leq 1$, the coefficient satisfies $-1 \leq \pi_{\theta}(y|x) - \pi_{\text{ref}}(y|x) \leq 1$. This boundedness enhances stability against large or pathological updates, in contrast to the unbounded log-ratio used by KL (see Figure 1). This supports our recommendation in Section 5 to consider bounded alternatives when stability is paramount.
2. **Symmetry in probability space.** The MSE penalty is symmetric with respect to probability differences, providing more conservative corrections when policies diverge, compared to the logarithmic penalty of Reverse KL.
3. **Off-policy compatibility.** Owing to its in-reward form (where the coefficient does not carry gradients), this head is fully compatible with importance sampling and clipping, following the same correction rules as in Section G.

Remark. Consistent with our KL analysis, Equation (71) is obtained by evaluating scalar multipliers at the detached snapshot π_{θ} (on-policy). This expresses the regularizer in the same coefficient-times-score-function form, enabling direct comparison of the induced update dynamics.

G. Off-Policy Correction for KL Regularization

All derivations so far have assumed *on-policy* sampling: the samples y are drawn from the current policy π_{θ} at the moment of gradient computation. In practice, however, many RLHF implementations—particularly PPO with multiple gradient steps per rollout—operate in an *off-policy* setting where samples are drawn from a stale behavior policy π_{θ_k} .

Off-policy updates require importance sampling (IS) and clipping for *both* the reward head and the KL head. When KL is implemented in the combined “in reward” form, the IS correction is inherited automatically from the PPO surrogate. However, when KL is implemented “as loss” in a decoupled form, the required IS/clipping on the KL term is easy to overlook, leading to biased gradients.

This section derives the principled correction and provides concrete guidance for both combined and decoupled objective structures.

G.1. From On-Policy to Off-Policy Gradients

In the policy-gradient view, updates are driven by a scalar coefficient $c(x, y)$ (which does not carry gradients) multiplying the score function. The on-policy gradient estimator is:

$$\nabla_{\theta} \mathcal{J}_c(\theta) = \mathbb{E}_{x \sim \mathcal{D}, y \sim \pi_{\theta}(\cdot | x)} [c(x, y) \nabla_{\theta} \log \pi_{\theta}(y | x)], \quad (73)$$

where π_{θ} is a detached snapshot numerically equal to π_{θ} at the time of gradient evaluation. For samples drawn from a behavior policy $y \sim \pi_{\theta_k}(\cdot | x)$, an unbiased off-policy estimator requires IS, assuming the behavior policy has support over the sampled data ($\pi_{\theta_k}(y | x) > 0$):

$$\nabla_{\theta} \mathcal{J}_c(\theta) = \mathbb{E}_{x \sim \mathcal{D}, y \sim \pi_{\theta_k}(\cdot | x)} \left[\underbrace{\frac{\pi_{\theta}(y | x)}{\pi_{\theta_k}(y | x)}}_{\text{IS weight (no gradient)}} c(x, y) \nabla_{\theta} \log \pi_{\theta}(y | x) \right]. \quad (74)$$

In practice, PPO replaces this detached IS weight with the gradient-carrying ratio $\rho_k(\theta) = \frac{\pi_{\theta}(y|x)}{\pi_{\theta_k}(y|x)}$ (where gradients flow only through the numerator) and employs a clipped surrogate objective to reduce variance. For any scalar coefficient $c(x, y)$ (which does not carry gradients), the clipped objective to be maximized is:

$$\mathcal{J}_{c, \text{clipped}}(\theta) = \mathbb{E}_{x \sim \mathcal{D}, y \sim \pi_{\theta_k}(\cdot | x)} \left[\min \left(\rho_k(\theta) c(x, y), \text{clip}(\rho_k(\theta), 1 - \epsilon, 1 + \epsilon) c(x, y) \right) \right]. \quad (75)$$

G.2. Correcting k_n as loss by Converting to an In-Reward Head

A k_n as loss head is gradient-correct only on-policy. To adapt it for off-policy use, it must first be converted to its gradient-equivalent in-reward form. This is achieved by defining a detached (stop-gradient) coefficient $k_{n'}(x, y)$ that reproduces the on-policy gradient of the original loss. For a differentiable penalty $k_n(\pi_{\theta}(y|x), \pi_{\text{ref}}(y|x))$, this coefficient is its derivative with respect to the policy’s log-probability, evaluated at the *current* detached snapshot:

$$k_{n'}(\pi_{\theta}(y|x), \pi_{\text{ref}}(y|x)) := \frac{\partial}{\partial \log \pi_{\theta}} k_n(\pi_{\theta}(y|x), \pi_{\text{ref}}(y|x)) \Big|_{\pi_{\theta} = \pi_{\theta}}. \quad (76)$$

This conversion precisely aligns with the theoretical equivalences established in the main text:

- **Principled k_2 as loss:** $k_2 = \frac{1}{2}(\log \pi_{\theta} - \log \pi_{\text{ref}})^2 \Rightarrow k_{2'} = \log \pi_{\theta} - \log \pi_{\text{ref}}$ (the k_1 in reward coefficient).
- **Proxy k_3 as loss:** $k_3 = \frac{\pi_{\text{ref}}}{\pi_{\theta}} - 1 - \log \frac{\pi_{\text{ref}}}{\pi_{\theta}} \Rightarrow k_{3'} = 1 - \frac{\pi_{\text{ref}}}{\pi_{\theta}}$ (the $k_{3'}$ in reward coefficient).

Once expressed as a scalar coefficient $k_{n'}(x, y)$ (which does not carry gradients), the KL head is handled off-policy exactly like any other score-function head via Equation (75). In PPO with multiple epochs per batch, $k_{n'}$ should be computed once using the rollout policy π_{θ_k} and held fixed across all epochs, consistent with standard PPO treatment of advantage estimates. The importance sampling ratio $\rho_k(\theta) = \pi_{\theta}/\pi_{\theta_k}$ still updates with the policy across epochs, providing the necessary off-policy correction.

G.3. Two Principled Off-Policy Integration Strategies

With the correctly derived coefficient $k_{n'}$ in hand, there are two principled ways to integrate it into the PPO objective, mirroring the on-policy discussion in Section 4.

1. Combined Form (Single Clipped Head). Merge the reward advantage and the KL coefficient *before* applying the PPO machinery:

$$A_{\text{combined}}(x, y) := r(x, y) - \beta k_{n'}(\pi_{\theta}(y|x), \pi_{\text{ref}}(y|x)). \quad (77)$$

The clipped surrogate is then applied to this combined head:

$$\mathcal{J}_{\text{RLHF}}(\theta) = \mathbb{E}_{y \sim \pi_{\theta_k}} \left[\min \left(\rho_k(\theta) A_{\text{combined}}, \text{clip}(\rho_k(\theta), 1 - \epsilon, 1 + \epsilon) A_{\text{combined}} \right) \right]. \quad (78)$$

This is a robust and straightforward approach, as IS and clipping are consistently applied to both components. For correct PPO semantics, form A_{combined} prior to any baseline subtraction or normalization. This preserves the trade-off set by β , which would be distorted by shifting or rescaling the KL component.

1210 **2. Decoupled Form (Two Clipped Heads).** Maintain separate reward and KL objectives, each with its own IS correction
 1211 and clipping scheme:

$$1212 \mathcal{J}_{\text{reward}}(\theta) = \mathbb{E}_{y \sim \pi_{\theta, k}} \left[\min \left(\rho_k(\theta) r(x, y), \text{clip}(\rho_k(\theta), 1 - \epsilon_1, 1 + \epsilon_2) r(x, y) \right) \right], \quad (79)$$

$$1213 \mathcal{J}_{\text{KL}}(\theta) = \mathbb{E}_{y \sim \pi_{\theta, k}} \left[\min \left(\rho_k(\theta) k_{n'}(\pi_{\theta}(y|x), \pi_{\text{ref}}(y|x)), \right. \right. \\ 1214 \left. \left. \text{clip}(\rho_k(\theta), 1 - \epsilon, 1 + \epsilon) k_{n'}(\pi_{\theta}(y|x), \pi_{\text{ref}}(y|x)) \right) \right]. \quad (80)$$

1218 The final objective to maximize is:

$$1219 \mathcal{J}_{\text{RLHF}}(\theta) := \mathcal{J}_{\text{reward}}(\theta) - \beta \mathcal{J}_{\text{KL}}(\theta). \quad (81)$$

1220 This decoupled design affords greater flexibility, such as using asymmetric clipping for the reward head (e.g., $\epsilon_2 > \epsilon_1$)
 1221 to accelerate learning, while retaining conservative, symmetric clipping for the KL head to ensure stable regularization.
 1222 Baselines or normalization should be applied only to A_{reward} , not to $k_{n'}(x, y)$, to avoid implicitly altering the regularization
 1223 strength β .
 1224

1225 Implementation Notes.

- 1226 • **Token vs. Sequence Level:** Our derivations use sequence-level probabilities. In token-level PPO, it is often more
 1227 stable to compute the ratio as $\rho_k = \exp(\sum_t \log \pi_{\theta}(y_t | \cdot) - \sum_t \log \pi_{\theta_k}(y_t | \cdot))$ and apply clipping at the sequence
 1228 level; per-token clipping can be overly conservative.
- 1229 • **Masking Consistency:** Sum log-probabilities only over action tokens that contribute to the reward/KL (exclude prompt,
 1230 padding, or masked tokens) to keep ρ_k aligned with the heads being optimized.
- 1231 • **Numerical Stability and Support:** Ensure $\pi_{\theta_k}(y | x) > 0$ for all sampled (x, y) and consider numerically capping ρ_k
 1232 to prevent overflows under extreme ratios.
- 1233 • **Adaptive β :** If targeting a desired KL via an adaptive schedule, update β outside the gradient path (detached) and
 1234 avoid mixing it with advantage normalization; adaptation is orthogonal to IS/clipping and works for both combined and
 1235 decoupled forms.

1236 **Sign Convention.** We present objectives for *maximization*. Implementations that *minimize* a loss should negate these
 1237 expressions, e.g., by minimizing $-\mathcal{J}_{\text{combined}}$ or $-(\mathcal{J}_{\text{reward}} - \beta \mathcal{J}_{\text{KL}})$.
 1238

1239 H. Visualization of KL Regularization Gradient Coefficients

1240 The theoretical analysis in Section 5 showed that different KL implementations induce different scalar coefficients multiplying
 1241 the score function. Visualizing these coefficients as a function of the policy’s log-probability provides valuable intuition
 1242 about their behavior.

1243 This section provides the Python code used to generate Figure 1. The plot compares four coefficient functions:

- 1244 • The principled RKL coefficient $\log(\pi_{\theta}/\pi_{\text{ref}})$ (achieved by k_1 in reward and k_2 as loss);
- 1245 • The biased GRPO coefficient $1 - \pi_{\text{ref}}/\pi_{\theta}$ (induced by k_3 as loss);
- 1246 • The bounded MiniMax-01 coefficient $\pi_{\theta} - \pi_{\text{ref}}$;
- 1247 • The constant coefficient 1 (induced by k_1 as loss, which provides no regularization).

1248 The visualization highlights the linear behavior of the principled RKL coefficient versus the asymmetric saturation/unbounded
 1249 growth of the k_3 surrogate in the tails.

```
1250 import torch
1251 import matplotlib.pyplot as plt
1252 # --- Plotting Style ---
1253
```

```

1265 plt.style.use('seaborn-v0_8-whitegrid')
1266 plt.rcParams.update({
1267     "text.usetex": False, # Disable LaTeX rendering
1268     "font.family": "serif", # Use a generic serif font
1269     "font.serif": ["Times New Roman"], # Specify Times New Roman as the serif font
1270     "font.size": 14,
1271     "axes.labelsize": 16,
1272     "legend.fontsize": 12,
1273     "xtick.labelsize": 12,
1274     "ytick.labelsize": 12,
1275 })
1276 # --- Data Generation ---
1277 log_pi_actor = torch.linspace(-5, 0, steps=400)
1278 pi_actor = torch.exp(log_pi_actor)
1279 pi_ref_val = 0.25
1280 log_pi_ref = torch.log(torch.tensor(pi_ref_val))
1281 # --- Coefficients Calculation ---
1282 coeff_k1_loss = torch.ones_like(log_pi_actor)
1283 coeff_k1_reward = log_pi_actor - log_pi_ref
1284 coeff_k3_loss = 1 - pi_ref_val / pi_actor
1285 coeff_minimax = pi_actor - pi_ref_val
1286 # --- Plotting ---
1287 plt.figure(figsize=(10, 6.5))
1288 plt.plot(log_pi_actor, coeff_k1_reward,
1289     label=r'$\log\pi_{\theta} - \log\pi_{\text{ref}}$ ($k_1$ in reward / $k_2$ as
1290     loss) - Principled',
1291     color='#808000', linewidth=3, zorder=10)
1292 plt.plot(log_pi_actor, coeff_k3_loss,
1293     label=r'$1 - \pi_{\text{ref}}/\pi_{\theta}$ ($k_{3^{\prime}}$ in reward / $k_3$
1294     as loss) - Biased Approximation',
1295     color='Firebrick', linestyle='--', linewidth=2)
1296 plt.plot(log_pi_actor, coeff_minimax,
1297     label=r'$\pi_{\theta} - \pi_{\text{ref}}$ (MiniMax-01)',
1298     color='RoyalBlue', linestyle='-.', linewidth=2)
1299 plt.plot(log_pi_actor, coeff_k1_loss,
1300     label=r'$1$ ($k_1$ as loss) - Zero Expected Gradient',
1301     color='Gray', linestyle=':', linewidth=2)
1302 plt.axvline(x=log_pi_ref.item(), color='black', linestyle='--', linewidth=1,
1303     label=r'$\log\pi_{\theta} = \log\pi_{\text{ref}}$')
1304 plt.xlabel(r'Actor Log-Probability: $\log \pi_{\theta}(y|x)$')
1305 plt.ylabel(r'Coefficient of Score Function')
1306 plt.title(r'Comparison of KL Regularization Coefficients ($\pi_{\text{ref}}=0.25$)',
1307     fontsize=18)
1308 plt.legend(loc='upper left')
1309 plt.ylim(-4, 4)
1310 plt.xlim(-5, 0)
1311 plt.tight_layout()
1312 plt.savefig('comparison_kl_regularization_coefficients.png', dpi=300, bbox_inches='tight')
1313 plt.show()

```

Listing 1. Python code to generate the comparison plot of KL gradient coefficients.

1320 I. On the Statistical Instability of the k_3 Value Estimator

1321 The main text argues that KL implementations should be evaluated by their gradient properties, not their value-estimation
 1322 properties. However, even when viewed purely as a value estimator, k_3 has significant limitations that are often overlooked.
 1323

1324 This section analyzes the statistical properties of the k_3 value estimator. We show that its advertised “unbiasedness” relies
 1325 on regularity conditions—specifically, absolute continuity of the reference with respect to the sampling distribution—that
 1326 can fail in practice. When these conditions are violated, k_3 can be biased, and its variance can become infinite. These failure
 1327 modes provide additional justification for preferring the principled k_1 in reward or k_2 as loss formulations.
 1328

1329 I.1. Precondition for Unbiasedness

1330 An estimator is unbiased if its expectation equals the true value. For k_3 :

$$1331 \mathbb{E}_q[k_3] = \mathbb{E}_q[\delta(x) - 1 - \log \delta(x)] = (\mathbb{E}_q[\delta(x)] - 1) + D_{\text{KL}}(q \parallel p) \quad (82)$$

1332 For k_3 to be unbiased, it is necessary that $\mathbb{E}_q[\delta(x)] = 1$. This condition is met if p is absolutely continuous with respect to q
 1333 ($p \ll q$), which means that the support of p must be contained within the support of q .
 1336

1337 The condition of a finite KL divergence ($D_{\text{KL}}(q \parallel p) < \infty$) is **not sufficient** to guarantee unbiasedness. For example, let q
 1338 be the uniform distribution in $[0, 1]$ and p be the uniform distribution on $[0, 2]$.

- 1339 • The KL divergence $D_{\text{KL}}(q \parallel p) = \int_0^1 1 \cdot \log(\frac{1}{0.5}) dx = \log 2$, which is finite.
- 1340 • However, $\mathbb{E}_q[\delta(x)] = \int_0^1 1 \cdot \frac{p(x)}{q(x)} dx = \int_0^1 \frac{0.5}{1} dx = 0.5$.
- 1341 • The estimator expectation is therefore $\mathbb{E}_q[k_3] = (0.5 - 1) + \log 2 = \log 2 - 0.5$, which is biased.

1342 I.2. Infinite Variance and the Chi-Squared Divergence

1343 The variance of k_3 is dominated by the second moment of the importance ratio, $\mathbb{E}_q[\delta(x)^2]$. This term is directly related to
 1344 the Chi-squared divergence.

1345 When $p \ll q$, the identity holds: $\chi^2(p \parallel q) = \mathbb{E}_q[(\delta(x) - 1)^2] = \mathbb{E}_q[\delta(x)^2] - 1$. If p is not absolutely continuous with
 1346 respect to q ($p \not\ll q$), $\chi^2(p \parallel q)$ is defined to be infinite.

1347 Therefore, the variance of k_3 will be infinite if $\mathbb{E}_q[\delta(x)^2]$ is infinite. This occurs if $p \not\ll q$ or if $p \ll q$ but the tails of q are
 1348 sufficiently lighter than the tails of p . While the divergence of $\mathbb{E}_q[\delta(x)^2]$ is the primary cause of instability, the finiteness of
 1349 $\text{Var}(k_3)$ also technically requires the finiteness of $\mathbb{E}_q[(\log \delta(x))^2]$.

1350 I.3. The Gaussian Case and an Empirical Demonstration

1351 For two Gaussian distributions, $p \sim \mathcal{N}(\mu_p, \sigma_p^2)$ and $q \sim \mathcal{N}(\mu_q, \sigma_q^2)$, the variance of k_3 is finite if and only if $\sigma_q^2 > \sigma_p^2/2$.
 1352 This condition illustrates that the sampling distribution q must be sufficiently “wide” relative to the reference distribution p .
 1353 This condition generalizes for multivariate Gaussians with covariance matrices Σ_p and Σ_q . **The expectation $\mathbb{E}_q[r(x)^2]$ is
 1354 calculated via an integral involving the ratio of two Gaussian probability densities. For this integral to converge (and
 1355 thus for the variance to be finite), it is required that the matrix $2\Sigma_q - \Sigma_p$ be positive definite.**

1356 This failure mode is empirically illustrated below, where a narrow Gaussian $q(x)$ ($\sigma_q = 0.2$) is used to estimate the KL
 1357 divergence to a standard Gaussian $p(x)$ ($\sigma_p = 1$). This configuration violates the condition, since $0.2^2 \not> 1^2/2$.

```

1358 import torch
1359 import torch.distributions as dist
1360 # p: reference distribution, q: sampling distribution
1361 p = dist.Normal(loc=0, scale=1)
1362 q = dist.Normal(loc=0.1, scale=0.2) # A narrow distribution where Var[k3] is infinite
1363 # Sample from the narrow distribution q
1364 x = q.sample(sample_shape=(10_000,))
1365

```

```

137511 # Ground truth KL divergence D_KL(q || p)
137612 true_kl = dist.kl_divergence(q, p)
137713
137814 # Compute the log-ratio log(p(x)/q(x))
137915 log_r = p.log_prob(x) - q.log_prob(x)
138016 r = torch.exp(log_r)
138117
138118 # Define estimators
138219 k1 = -log_r
138320 k2 = log_r.pow(2) / 2
138421 k3 = r - 1 - log_r
138522
138523 # --- Code to generate output ---
138624 print(f"True KL Divergence: {true_kl:.4f}\n")
138725 print("Estimator          | Sample Mean   | Sample Std. Dev.")
138826 print("-----|-----|-----")
138927 estimators = {"k1": k1, "k2": k2, "k3": k3}
139028
139129 for name, k in estimators.items():
139230     mean = k.mean()
139331     std = k.std()
139432     print(f"{name:<17} | {mean:>13.4f} | {std:>16.4f}")
139533
139534 # --- Actual Output 1 ---
139635 True KL Divergence: 1.1344
139736
139737 Estimator          | Sample Mean   | Sample Std. Dev.
139838 -----|-----|-----
139939 k1                  |          1.1272 |          0.6912
140040 k2                  |          0.8742 |          0.6006
140141 k3                  |          0.8136 |          8.8244
140242 # --- Actual Output 2 ---
140343 True KL Divergence: 1.1344
140444
140445 Estimator          | Sample Mean   | Sample Std. Dev.
140546 -----|-----|-----
140647 k1                  |          1.1336 |          0.6611
140748 k2                  |          0.8611 |          0.5210
140849 k3                  |          0.6817 |          4.1082
140950 # --- Actual Output 3 ---
141051 True KL Divergence: 1.1344
141152
141153 Estimator          | Sample Mean   | Sample Std. Dev.
141254 -----|-----|-----
141355 k1                  |          1.1348 |          0.6709
141456 k2                  |          0.8689 |          0.4968
141557 k3                  |          0.6595 |          1.4925
141658 # --- Actual Output 4 ---
141759 True KL Divergence: 1.1344
141860
141861 Estimator          | Sample Mean   | Sample Std. Dev.
141962 -----|-----|-----
142063 k1                  |          1.1256 |          0.6962
142164 k2                  |          0.8758 |          0.6263
142265 k3                  |          0.9772 |         26.7379

```

Listing 2. Code illustrating the high variance of the k_3 value estimator when the sampling distribution $q(x)$ is too narrow.

The results illustrate the issue: across repeated trials, the sample standard deviation of k_3 can be substantially larger than that of k_1 and k_2 , and the sample mean can deviate noticeably from the true KL value. This discrepancy is not systematic estimator bias, but rather sampling error caused by heavy-tailed variance, implying that far more samples may be needed for reliable estimation in such settings.

J. Group Normalization Stability Issues

Beyond the gradient-level issues analyzed in the main text, GRPO’s per-prompt group normalization can introduce a separate numerical stability problem. When the variance within a group of responses is very small, normalization can dramatically amplify tiny numerical differences, potentially destabilizing training.

This section describes the issue and proposes a simple fix: clipping the group standard deviation to prevent pathological amplification.

The issue. GRPO normalizes rewards within each group: for a prompt with G responses and rewards $\mathbf{r} = \{r_1, \dots, r_G\}$, the advantage is

$$A_i = \frac{r_i - \text{mean}_{\text{group}}(\mathbf{r})}{\text{std}_{\text{group}}(\mathbf{r})}. \quad (83)$$

Stability issue. When the within-group variance is very small (e.g., $\mathbf{r} = [0.99999, 1.00001, 0.99999, 1.00001]$), normalization can substantially amplify tiny numerical differences. For the above example, the resulting advantages become approximately $[-0.8660, 0.8660, -0.8660, 0.8660]$ (using the unbiased sample standard deviation), which can destabilize optimization by turning near-constant rewards into large-magnitude updates.

Proposed solution. Clip the standard deviation to prevent pathological amplification:

$$A_i = \frac{r_i - \text{mean}_{\text{group}}(\mathbf{r})}{\text{clip_std}_{\text{group}}(\mathbf{r})}, \quad (84)$$

$$\text{clip_std}_{\text{group}}(\mathbf{r}) = \max\left(\min(\text{std}_{\text{group}}(\mathbf{r}), \text{std}_{\text{max}}), \text{std}_{\text{min}}\right).$$

Here, $\text{std}_{\text{min}} > 0$ is a small floor that prevents exploding normalization when variance collapses, and std_{max} avoids under-normalization when variance is unusually large. In practice, setting std_{min} as a small constant relative to the reward scale (e.g., 10^{-1}) may be effective.

Why this matters beyond binary rewards. Although binary 0/1 rewards in RLVR can sometimes mitigate extreme cases, more general regression reward models—such as those trained with Bradley–Terry (BT) losses—often produce continuous scores that may become highly concentrated (e.g., near 0 or 1) on easy or very hard prompts. In such regimes, within-group standard deviations can be arbitrarily small even when rewards are bounded in $[0, 1]$, and group normalization will over-scale negligible differences unless a variance floor (or clipping) is used. Therefore, std clipping is important not only for numerical stability but also to avoid over-amplifying noise when reward predictions saturate.

Remark. For reward scores bounded in $[0, 1]$, $\text{std}(\mathbf{r}) < 1$ always holds, but it can be orders of magnitude smaller than 1 in practice; the smaller the variance, the stronger the amplification effect from group normalization. Clipping $\text{std}_{\text{group}}(\mathbf{r})$ preserves the intended scale-invariance when variance is moderate, while guarding against instability when variance collapses.

K. Forward KL vs. Reverse KL: A Gradient-Centric Reinterpretation

The main text focuses on the Reverse KL (RKL) divergence—the standard regularization objective in RLHF—and shows that k_3 as loss is a biased first-order surrogate for the RKL gradient. An alternative perspective, recently noted by Tang & Munos (2025), is that k_3 as loss actually optimizes a Forward KL (FKL) objective, $D_{\text{KL}}(\pi_{\text{ref}} || \pi_{\theta})$.

This section refines this observation using our gradient framework and provides a deeper understanding of k_3 as loss. We show that it is a *variance-reduced estimator* of the FKL gradient with an implicit baseline of -1 . This explains two empirical observations:

- **Low variance near the reference:** When $\pi_{\theta} \approx \pi_{\text{ref}}$, the coefficient $1 - \delta$ fluctuates near zero, giving low variance.
- **Weak constraint under drift:** The FKL geometry is “mean-seeking” rather than “mode-seeking,” imposing only a finite penalty for generating out-of-distribution tokens.

These properties make k_3 as loss behave well near the reference but provide weaker regularization when the policy drifts—consistent with the training instabilities observed in Section 6.

Derivation of the FKL Gradient. The Forward KL objective is defined as:

$$\mathcal{J}_{\text{FKL}}(\theta) = \mathbb{E}_{x \sim \mathcal{D}} \left[D_{\text{KL}}(\pi_{\text{ref}}(\cdot|x) \parallel \pi_{\theta}(\cdot|x)) \right] = \mathbb{E}_{x \sim \mathcal{D}, y \sim \pi_{\text{ref}}(\cdot|x)} \left[\log \pi_{\text{ref}}(y|x) - \log \pi_{\theta}(y|x) \right]. \quad (85)$$

Rewriting the gradient expectation using importance sampling over the current policy π_{θ} (to allow on-policy estimation) yields the standard FKL policy gradient:

$$\nabla_{\theta} \mathcal{J}_{\text{FKL}}(\theta) = \mathbb{E}_{x \sim \mathcal{D}, y \sim \pi_{\theta}(\cdot|x)} \left[\underbrace{-\frac{\pi_{\text{ref}}(y|x)}{\pi_{\theta}(y|x)}}_{\text{Standard FKL coeff. } (-\delta)} \nabla_{\theta} \log \pi_{\theta}(y|x) \right], \quad \text{where } \delta := \frac{\pi_{\text{ref}}(y|x)}{\pi_{\theta}(y|x)}. \quad (86)$$

k_3 as Loss: **FKL with an Implicit Baseline.** Recalling Equation (18), the gradient induced by the k_3 as loss formulation uses the coefficient $1 - \delta$. Comparing this with the standard FKL coefficient in Equation (86) reveals an implicit decomposition:

$$\nabla_{\theta} \mathcal{J}_{k_3 \text{ as loss}}(\theta) = \mathbb{E}_{x \sim \mathcal{D}, y \sim \pi_{\theta}(\cdot|x)} \left[\left(\underbrace{(-\delta)}_{\text{Standard FKL}} - \underbrace{(-1)}_{\text{Implicit Baseline } b} \right) \nabla_{\theta} \log \pi_{\theta}(y|x) \right]. \quad (87)$$

Although mathematically equivalent in expectation (since $\mathbb{E}[\nabla \log \pi] = 0$), the implicit baseline $b = -1$ acts as a **control variate**:

- **Variance Reduction at Convergence:** Near the reference policy ($\pi_{\theta} \approx \pi_{\text{ref}}$), we have $\delta \approx 1$. The standard FKL estimator $-\delta$ fluctuates around -1 (high variance), whereas the k_3 coefficient $(1 - \delta)$ fluctuates around 0. This explains why k_3 exhibits low variance specifically in the low-KL regime: it is effectively a zero-variance estimator of the FKL gradient at the optimum.

Geometric Implications: Mean-Seeking vs. Mode-Seeking. While the implicit baseline stabilizes estimation, optimizing FKL fundamentally alters the regularization geometry compared to the principled RKL:

- **RKL (Mode-seeking):** As derived in Section 5.2, RKL uses $-\log \delta$. As $\delta \rightarrow 0$ (policy places mass where reference does not), $-\log \delta \rightarrow \infty$. This creates a “barrier” that forces the policy to stay within the reference’s support.
- **FKL (Mean-seeking):** The k_3 coefficient $1 - \delta$ saturates at 1 as $\delta \rightarrow 0$. This “mean-seeking” behavior imposes only a finite penalty for generating out-of-distribution tokens. Consequently, under distribution shift, k_3 as loss fails to strongly penalize drift into regions unsupported by the reference model, consistent with the higher variability and weaker constraint compliance observed in our training curves.

In summary, k_3 as loss can be seen as a statistically coherent, baseline-corrected estimator of the FKL gradient, with favorable variance properties near the reference policy. However, its underlying geometry (mean-seeking and mode-covering) remains fundamentally different from that of RKL (mode-seeking). For RLHF applications where keeping the policy tightly constrained within the support of the reference model is a primary concern, the principled RKL implementations discussed in Section 5.2 may offer stronger and more reliable regularization.

k_3 in reward: **a mixed (and counterproductive) gradient.** A natural question is whether the nonnegative single-sample term k_3 can be used safely as an *in-reward coefficient* (i.e., a scalar multiplier without gradients). Our framework shows it cannot: unlike k_1 , it does not induce the RKL gradient.

Let $\delta = \pi_{\text{ref}}(y|x)/\pi_{\theta}(y|x)$. By definition, $k_3 = \delta - 1 - \log \delta = k_1 + (\delta - 1)$ with $k_1 = -\log \delta$. Therefore, under on-policy sampling,

$$\begin{aligned} \nabla_{\theta} \mathcal{J}_{k_3 \text{ in reward}}(\theta) &= \mathbb{E}_{x \sim \mathcal{D}, y \sim \pi_{\theta}(\cdot|x)} \left[(k_1 + \delta - 1) \nabla_{\theta} \log \pi_{\theta}(y|x) \right] \\ &= \nabla_{\theta} \mathcal{J}_{\text{RKL}}(\theta) + \mathbb{E}_{x \sim \mathcal{D}, y \sim \pi_{\theta}(\cdot|x)} \left[\delta \nabla_{\theta} \log \pi_{\theta}(y|x) \right], \end{aligned} \quad (88)$$

where the constant term -1 drops by the zero-mean score identity. The remaining extra term rewrites as

$$\mathbb{E}_{x \sim \mathcal{D}, y \sim \pi_\theta(\cdot|x)} \left[\delta \nabla_{\theta} \log \pi_\theta(y|x) \right] = \mathbb{E}_{x \sim \mathcal{D}, y \sim \pi_{\text{ref}}(\cdot|x)} \left[\nabla_{\theta} \log \pi_\theta(y|x) \right] = -\nabla_{\theta} \mathcal{J}_{\text{FKL}}(\theta), \quad (89)$$

where the last equality follows from Equation (86). Hence,

$$\nabla_{\theta} \mathcal{J}_{k_3 \text{ in reward}}(\theta) = \nabla_{\theta} \mathcal{J}_{\text{RKL}}(\theta) - \nabla_{\theta} \mathcal{J}_{\text{FKL}}(\theta). \quad (90)$$

Equivalently, using k_3 as a reward penalty corresponds to optimizing a *difference* of KL objectives, RKL $-$ FKL, not their sum. In RLHF, this adds an *anti-FKL* component that can counteract the intended constraint and exacerbate instability. We therefore recommend using k_1 (or its gradient-equivalent surrogates) for reward shaping instead.

L. Analysis of the DeepSeek-V3.2 Importance-Weighted KL Loss

The off-policy correction framework in Section G showed that k_n as loss implementations require explicit importance sampling (IS) and clipping when samples are drawn from a stale policy. A natural question is whether alternative IS constructions can avoid these complications.

DeepSeek-V3.2 (Liu et al., 2025a) proposes an IS-corrected k_3 estimator for off-policy KL regularization. Interestingly, PPO-style clipping is applied only to the reward head while the KL term is left **unclipped**. This section analyzes this construction through our gradient-centric lens and highlights a practical trade-off: under full differentiation, the estimator recovers the principled k_1 coefficient through cancellation; however, leaving the IS-weighted KL term unclipped can introduce numerical instability when the policy diverges from the behavior policy.

The estimator. DeepSeek defines the following IS-weighted k_3 loss (we use $\mathcal{L}_{\text{KL,DS}}$ to distinguish it from the true KL divergence D_{KL}):

$$\mathcal{L}_{\text{KL,DS}}(\theta) = \frac{\pi_\theta}{\pi_{\text{old}}} \left(\frac{\pi_{\text{ref}}}{\pi_\theta} - \log \frac{\pi_{\text{ref}}}{\pi_\theta} - 1 \right). \quad (91)$$

Our gradient-centric analysis shows that, under full differentiation, this IS-corrected k_3 construction is mathematically equivalent to an in-reward k_1 coefficient, while leaving the IS-weighted KL term unclipped can be numerically fragile. We highlight two issues:

Mathematical equivalence. As derived below, minimizing the fully differentiated expectation of $\mathbb{E}[\rho(\theta) \cdot k_3(\theta)]$ analytically recovers the exact RKL gradient:

$$\nabla_{\theta} \mathcal{J}_{\text{RKL}} = \mathbb{E} \left[\underbrace{\rho(\theta) \cdot \left(-\log \frac{\pi_{\text{ref}}}{\pi_\theta} \right)}_{\text{Evaluated Scalar Coefficient}} \cdot \nabla_{\theta} \log \pi_\theta \right]. \quad (92)$$

DeepSeek’s estimator achieves this through cancellation of terms, effectively reconstructing the k_1 coefficient $-\log(\pi_\theta/\pi_{\text{ref}})$. Thus, under full differentiation, it is mathematically identical to the standard k_1 **in reward** formulation (without clipping), while requiring additional computation.

Practical trade-off with unclipped IS weights. Leaving the IS-weighted KL term **unclipped** creates a practical trade-off:

- **Full-gradient case:** If the gradient path through $\rho(\theta)$ is preserved (so that Part A + Part B terms are both present), the resulting coefficient scales with the potentially large importance weight $\rho(\theta)$. Without clipping, this can yield very large gradients when the policy diverges from π_{old} .
- **Detached-weight case:** If ρ is detached for stability, the estimator drops the sampling-distribution gradient term and reduces to the biased GRPO direction $(1 - \pi_{\text{ref}}/\pi_\theta)$, losing the equivalence to the RKL gradient.

L.1. Mathematical Derivation and Dilemma Analysis

We analyze the gradient of the KL estimator in Equation (91). Let π_θ denote the trainable policy carrying gradients, and π_{old} denote the behavior policy (detached).

Full Gradient Derivation (Equivalence to k_1). The objective function corresponding to Equation (91) is:

$$\mathcal{L}_{\text{KL,DS}}(\theta) = \mathbb{E}_{y \sim \pi_{\text{old}}} \left[\underbrace{\frac{\pi_{\theta}(y)}{\pi_{\text{old}}(y)}}_{\rho(\theta)} \cdot \underbrace{\left(\frac{\pi_{\text{ref}}(y)}{\pi_{\theta}(y)} - \log \frac{\pi_{\text{ref}}(y)}{\pi_{\theta}(y)} - 1 \right)}_{k_3(\theta)} \right]. \quad (93)$$

Applying the product rule $\nabla_{\theta}(f \cdot g) = g \nabla_{\theta} f + f \nabla_{\theta} g$:

$$\nabla_{\theta}(\rho(\theta) \cdot k_3(\theta)) = \underbrace{(\nabla_{\theta} \rho(\theta)) \cdot k_3(\theta)}_{\text{Part A}} + \underbrace{\rho(\theta) \cdot (\nabla_{\theta} k_3(\theta))}_{\text{Part B}}. \quad (94)$$

Note that in the product rule expansion, the term not being differentiated is treated as an evaluated value (hence black θ).

Part A (from $\nabla_{\theta} \rho$): Using the log-derivative trick $\nabla_{\theta} \rho(\theta) = \rho(\theta) \cdot \nabla_{\theta} \log \pi_{\theta}$:

$$\text{Part A} = \rho(\theta) \cdot \underbrace{\left(\frac{\pi_{\text{ref}}}{\pi_{\theta}} - \log \frac{\pi_{\text{ref}}}{\pi_{\theta}} - 1 \right)}_{k_3 \text{ scalar value}} \cdot \nabla_{\theta} \log \pi_{\theta}. \quad (95)$$

Part B (from $\nabla_{\theta} k_3$): Differentiating k_3 with respect to θ yields the biased GRPO coefficient:

$$\nabla_{\theta} k_3(\theta) = \left(1 - \frac{\pi_{\text{ref}}}{\pi_{\theta}} \right) \nabla_{\theta} \log \pi_{\theta}, \quad (96)$$

and therefore:

$$\text{Part B} = \rho(\theta) \cdot \left(1 - \frac{\pi_{\text{ref}}}{\pi_{\theta}} \right) \cdot \nabla_{\theta} \log \pi_{\theta}. \quad (97)$$

Total Gradient (Cancellation): Summing Part A and Part B, the bias terms in the coefficients cancel perfectly:

$$\begin{aligned} \nabla_{\theta} \mathcal{L}_{\text{KL,DS}} &= \rho(\theta) \cdot \left[\left(\frac{\pi_{\text{ref}}}{\pi_{\theta}} - \log \frac{\pi_{\text{ref}}}{\pi_{\theta}} - 1 \right) + \left(1 - \frac{\pi_{\text{ref}}}{\pi_{\theta}} \right) \right] \cdot \nabla_{\theta} \log \pi_{\theta} \\ &= \rho(\theta) \cdot \underbrace{\left(-\log \frac{\pi_{\text{ref}}}{\pi_{\theta}} \right)}_{\text{evaluated } k_1 \text{ coefficient}} \cdot \nabla_{\theta} \log \pi_{\theta}. \end{aligned} \quad (98)$$

Conclusion. Under full differentiation, this construction recovers the same score-function coefficient as k_1 in reward (up to the IS weight). In this sense, it does not provide a different regularization direction; it realizes the k_1 coefficient through cancellation within a more complex expression.

A practical trade-off with unclipped IS weights. DeepSeek-V3.2 separates this term and leaves it **unclipped**. This leads to the following trade-off:

Case 1 (full gradient path). If $\rho(\theta)$ is kept active (i.e., both Part A and Part B are included), the gradient magnitude scales with $\rho(\theta)$. When the policy diverges from π_{old} (e.g., $\pi_{\theta} \gg \pi_{\text{old}}$), $\rho(\theta)$ can become large; without clipping, this can produce very large gradients and numerical instability.

Case 2 (detached weight). If ρ is detached for stability, then $\nabla_{\theta} \rho \equiv 0$ and Part A vanishes. The estimator reduces to Part B:

$$\nabla_{\theta} \mathcal{J}_{\text{detached}} = \rho(\theta) \cdot \left(1 - \frac{\pi_{\text{ref}}}{\pi_{\theta}} \right) \cdot \nabla_{\theta} \log \pi_{\theta}. \quad (99)$$

This recovers the **biased GRPO coefficient** $1 - \pi_{\text{ref}}/\pi_{\theta}$, rather than the principled k_1 coefficient.

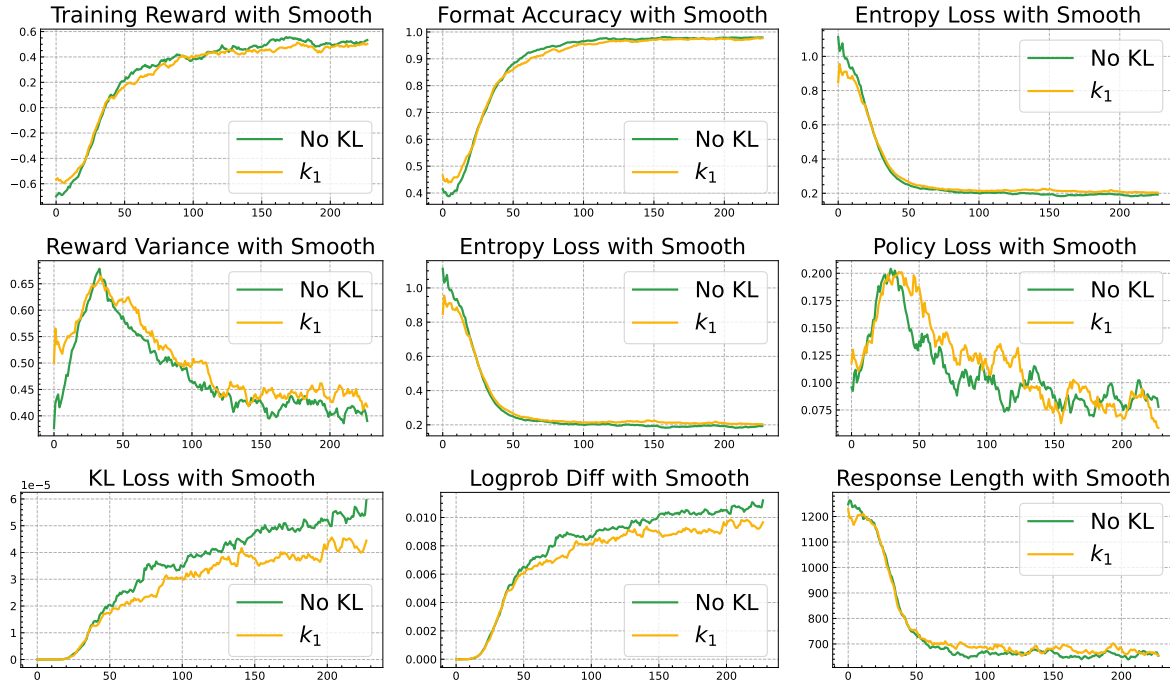


Figure 5. 1.5B-scale comparison of k_1 as loss versus no KL regularization. The gap is smaller than at 7B, but k_1 as loss still does not provide reliable regularization and can add stochasticity, consistent with its zero-mean, reference-independent expected gradient (Section 5).

Summary. This analysis highlights a practical tension when using unclipped, IS-weighted decoupled KL terms: preserving the full gradient path can lead to large updates when ρ is large, while detaching ρ changes the regularization direction. The principled approaches in Section 5.2— k_1 in reward with standard PPO correction, or the gradient-equivalent k_2 as loss—avoid this particular failure mode.

M. Additional 1.5B-Scale Experiments

The main text presents experimental results at 7B scale. To assess the generality of our findings, this section provides complementary experiments at 1.5B scale using the Qwen2.5-Math-1.5B model.

The 1.5B experiments test the same hypotheses as the main text: (i) k_1 as loss should provide no meaningful regularization, and (ii) k_2 as loss should enforce a stronger constraint than k_3 as loss. Consistent with our theoretical analysis, we find that these patterns hold at 1.5B scale, though the differences between methods are somewhat less pronounced than at 7B.

Compared to the 7B-scale behavior in the main text (Figure 2), the 1.5B runs separate less. This matches our analysis: k_1 as loss mainly behaves as a zero-mean perturbation, so its qualitative effect is smaller at 1.5B and more pronounced at 7B in our runs.

Weaker KL weight ($\beta = 0.001$). We also conduct experiments at 1.5B scale with a weaker and more realistic KL weight ($\beta = 0.001$) to test the sensitivity of different surrogates under larger drift. As shown in Figure 7, the **KL Loss** panel reveals a clear loss spike for k_3 as loss (orange curve), while k_2 as loss (green curve) remains stable throughout training. This pattern mirrors the 7B results in Figure 4, confirming that the instability of the k_3 surrogate under distribution shift is consistent across model scales.

N. Downstream Benchmark Performance

The main text (Table 3) reports a subset of downstream results for 7B models. This section provides the full evaluation across both model scales and additional benchmarks.

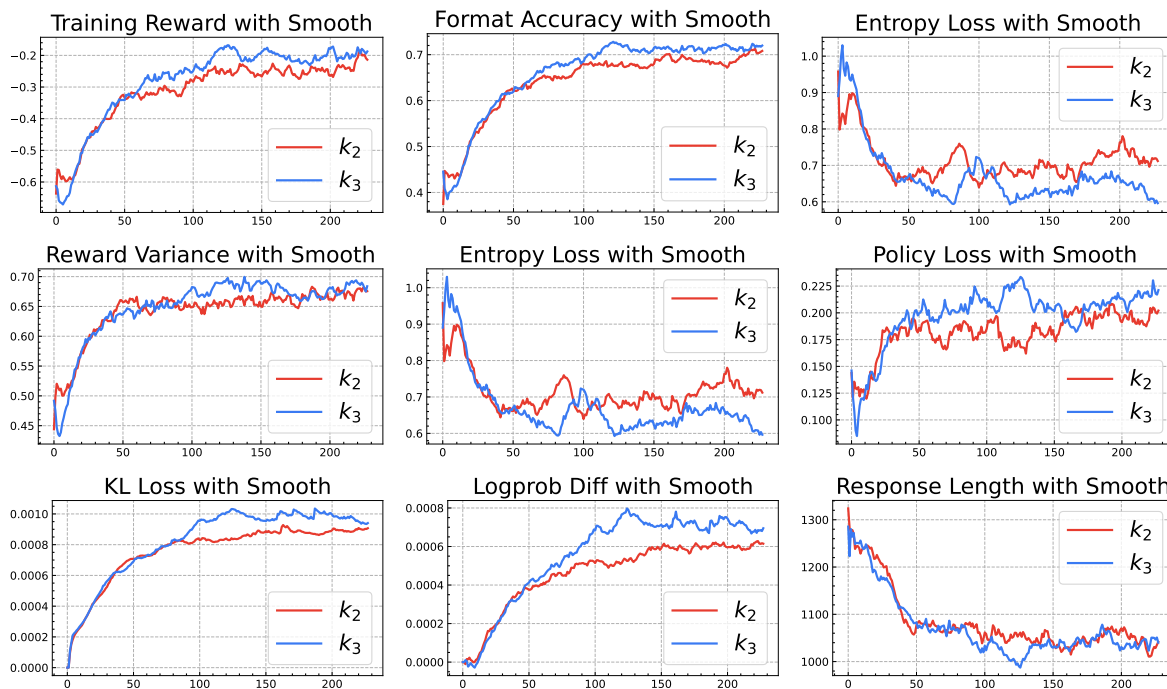


Figure 6. 1.5B-scale comparison of the principled k_2 as loss against its first-order surrogate k_3 as loss. Both variants constrain the policy, but k_2 as loss enforces a stronger, more stable regularization effect.

Table 4 summarizes downstream performance for 7B and 1.5B models across six math benchmarks (AIME 24/25, AMC, MATH-500, Minerva, Olympiad) and three general reasoning benchmarks (ARC-c, GPQA, MMLU-Pro). Several patterns emerge that are consistent with our theoretical analysis:

k_1 as loss provides no regularization benefit. Performance is close to the no-KL baseline, consistent with our analysis: its expected gradient vanishes and is independent of the reference model.

k_2 as loss enforces a stronger constraint than k_3 as loss. The principled k_2 as loss maintains tighter coupling to the reference, which in this setting correlates with lower downstream scores on both math benchmarks (AIME (Li et al., 2024), AMC, MATH-500 (Hendrycks et al., 2021)) and general reasoning benchmarks (ARC-c (Clark et al., 2018), GPQA* (Rein et al., 2024), MMLU-Pro (Wang et al., 2024)). In contrast, k_3 as loss allows larger divergence, which can yield higher reward but less stable training.

Summary. These results align with our theoretical analysis: k_1 as loss does not constrain; k_2 as loss implements the principled RKL and enforces a strong constraint; k_3 as loss is a weaker surrogate that permits larger drift.

O. Statement on the Use of Large Language Models

We used LLMs solely for language polishing and editing. All retrieval of related work, algorithmic design, and theoretical derivations are carried out by the authors.

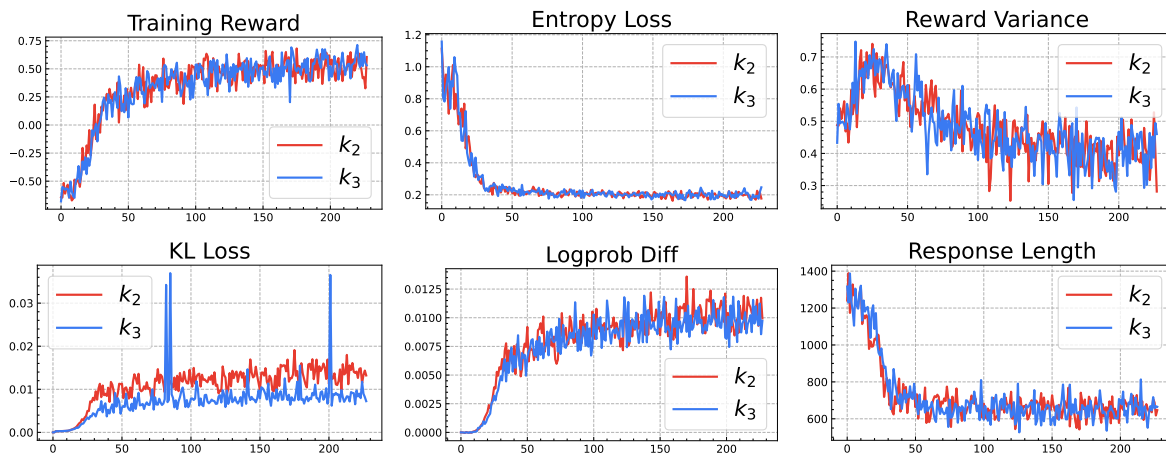


Figure 7. 1.5B-scale training diagnostics for k_2 as loss and k_3 as loss at a weaker KL weight ($\beta = 0.001$). The **KL Loss** panel (bottom-left) shows a visible spike for k_3 as loss that is absent for k_2 as loss, consistent with the 7B-scale results in Figure 4. This demonstrates that the k_3 surrogate’s sensitivity to distribution shift persists across model scales.

Table 4. Main experiment results ($\beta = 0.5$) on math and general reasoning benchmarks based on **Qwen2.5-Math-7B** and **Qwen2.5-Math-1.5B**.

Model	Math Reasoning Performance						General Domain Reasoning Performance			
	AIME 24/25	AMC	MATH-500	Minerva	Olympiad	Avg.	ARC-c	GPQA*	MMLU-Pro	Avg.
Qwen2.5-Math-7B	11.5/4.9	31.3	43.6	7.4	15.6	19.0	18.2	11.1	16.9	15.4
RL w/o KL	20.5/14.4	55.6	78.6	36.8	42.4	41.4	81.7	33.8	46.9	54.1
RL w/. k_1 as loss	19.1/11.6	56.0	80.6	40.8	43.0	41.8	79.7	29.8	45.1	51.5
RL w/. k_2 as loss	15.4/7.5	48.5	64.2	16.9	24.9	29.6	31.3	15.2	27.1	24.5
RL w/. k_3 as loss	19.0/7.3	48.9	65.4	18.8	29.0	31.4	29.6	19.2	27.7	25.5
Qwen2.5-Math-1.5B	7.2/3.6	26.4	28.0	9.6	21.2	16.0	3.5	4.0	2.5	3.3
RL w/o KL	12.5/4.8	43.7	66.8	28.3	31.9	31.3	43.7	19.2	23.1	28.7
RL w/. k_1 as loss	13.8/4.7	41.5	68.0	25.7	31.9	30.9	36.6	18.2	21.0	25.3
RL w/. k_2 as loss	7.0/5.5	35.2	52.8	14.7	29.0	24.0	7.8	7.6	4.9	6.8
RL w/. k_3 as loss	7.7/3.8	34.9	54.2	15.8	28.0	24.1	11.3	8.1	5.5	8.3