

# Relative Score Policy Optimization for Diffusion Language Models

Zichao Yu<sup>1</sup> Shengze Xu<sup>2</sup> Bingqing Jiang<sup>3</sup> Wenyi Zhang<sup>1</sup> Difan Zou<sup>3</sup>

## Abstract

Diffusion large language models (dLLMs) offer a promising route to parallel and efficient text generation, but improving their reasoning ability requires effective post-training. Reinforcement learning with verifiable rewards (RLVR) is a natural choice for this purpose, yet its application to dLLMs is hindered by the absence of tractable sequence-level log-ratios, which are central to standard policy optimization. The lack of tractable sequence-level log-ratios forces existing methods to rely on high-variance ELBO-based approximations, where high verifier rewards can amplify inaccurate score estimates and destabilize RL training. To overcome this issue, we propose **Relative Score Policy Optimization (RSPO)**, a simple RLVR method that uses verifiable rewards to calibrate noisy likelihood estimates in dLLMs. The core of our algorithm relies on a key observation: a reward advantage can be interpreted not only as an update direction, but also as a target for the relative log-ratio between the current and reference policies. Accordingly, RSPO calibrates this noisy relative log-ratio estimate by comparing its reward advantage with the reward-implied target relative log-ratio, updating the policy according to the gap between the current estimate and the target rather than the raw advantage alone. Experiments on mathematical reasoning and planning benchmarks show that RSPO yields especially strong gains on planning tasks and competitive mathematical-reasoning performance.

## 1. Introduction

Reinforcement learning with verifiable rewards (RLVR) (He et al., 2025; Shao et al., 2024; Yu et al., 2025; Zheng et al., 2025a) has become an effective post-training paradigm for improving language models on reasoning tasks: sample candidate answers, verify their correctness, and update the model so that correct answers become more likely (Sutton et al., 1998; Ouyang et al., 2022; Jaech et al., 2024; Guo et al., 2025). For autoregressive language models, this update is naturally expressed through policy log-ratios, namely the change in log-likelihood of a sampled response under the current policy relative to a reference or previous policy. Because autoregressive models admit an explicit left-to-right factorization, these log-ratios are tractable, enabling methods such as Proximal Policy Optimization (PPO) (Schulman et al., 2017) and Group Relative Policy Optimization (GRPO) (Shao et al., 2024) to combine reward improvement with Kullback–Leibler (KL) control.

Diffusion large language models (dLLMs) (Sahoo et al., 2024; Shi et al., 2024; Nie et al.; Yang et al., 2025; Bie et al., 2025) break this convenient interface. Their generation process proceeds through iterative denoising rather than left-to-right prediction, so the sequence likelihood of a completed response is not directly available. As a result, the sequence-level policy log-ratios used by standard RL objectives become difficult to compute. Recent RL methods for dLLMs therefore replace exact log-ratios with tractable surrogates, such as one-step mean-field estimates (Zhao et al., 2025), ELBO-based token or sequence scores (Ou et al., 2025), trajectory-level formulations (Wang et al., 2025b), and evidence-bound objectives (Wang et al., 2025a; Lin et al., 2025). These substitutes make policy optimization possible, but they also introduce a new source of uncertainty: in RLVR, the verifier reward is a reliable task signal, whereas the model-side relative score used to apply that reward is only an approximate and noisy estimate of the true likelihood change.

<sup>1</sup> University of Science and Technology of China, Hefei, Anhui, China

<sup>2</sup> The Chinese University of Hong Kong, Hong Kong, China

<sup>3</sup> The University of Hong Kong, Hong Kong, China

Correspondence to: Difan Zou <dzou@hku.hk>.

Accepted to FoGen 2026: Foundations of Deep Generative Models: Understanding Memorization, Generalization, and Reasoning, an ICML 2026 workshop (non-archival).

This mismatch raises a basic question: how should dLLM RLVR combine a reliable verifier reward with a noisy likelihood-ratio surrogate? Standard advantage-weighted objectives address only part of this question. They multiply a model-side score by a reward or advantage coefficient, thereby specifying which responses should be encouraged. When the score is an ELBO-based log-ratio surrogate, however, the coefficient does not depend on the current relative score. A high-advantage response can therefore continue to amplify a noisy likelihood estimate even after its estimated relative likelihood has moved far from the reference. Thus, the reward provides a preference direction, but the objective lacks a calibrated target for the relative score itself.

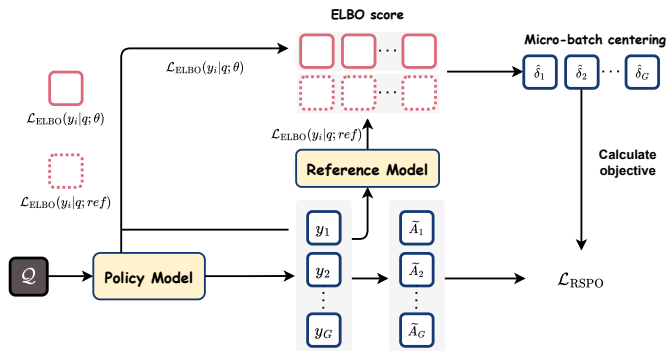


Figure 1. Overview of RSPO: verifier rewards define centered relative-score targets, and the residual to the current ELBO-based score provides policy-update feedback.

Our key observation is that KL-regularized policy improvement provides exactly this target. The optimal KL-regularized policy assigns each response a current-reference log-ratio determined by its reward, up to a prompt-level normalization term. Within a comparison group, this normalization term is shared and disappears after centering. The centered verifier rewards therefore specify the relative likelihood changes that the improved policy should assign to sampled responses. Based on this principle, we propose **Relative Score Policy Optimization (RSPO)**, a simple RLVR objective for diffusion language models. RSPO estimates a centered current-reference score with ELBO differences, derives a centered relative-score target from verifier rewards, and uses the residual gap as feedback in the policy update. In this way, reliable rewards calibrate noisy likelihood-ratio estimates instead of serving only as fixed multiplicative weights. Figure 1 gives an overview of this feedback loop: the policy generates a group of responses for each prompt, the reference model is used to form ELBO-based log-ratio scores, and micro-batch centering turns these scores into relative feedback. RSPO compares this current centered score with the verifier-induced target and optimizes the resulting residual objective; details are deferred to Section 3. Experiments on mathematical reasoning and planning benchmarks show that RSPO yields large planning gains and competitive mathematical-reasoning performance against existing RL baselines.

Our contributions are summarized as follows:

- We identify a calibration problem in dLLM RLVR. Verifier rewards provide reliable task-level feedback, but the ELBO-based log-ratios used for policy optimization are noisy likelihood-ratio surrogates. Standard advantage-weighted objectives use these rewards as fixed multiplicative coefficients, which can keep pushing already high-scoring responses and over-amplify noisy relative-score estimates.
- We derive a principled relative-score target from KL-regularized policy improvement. Although the optimal policy log-ratio contains an unknown normalization term, this term is shared within a comparison group and can be removed by centering. The resulting centered form shows that verifier rewards specify not only which responses should be preferred, but also the relative likelihood changes an improved policy should assign to them.
- We propose RSPO, a simple stop-gradient objective for diffusion language models. Instead of using rewards only as fixed advantages, RSPO compares the current centered ELBO-based relative score with the reward-implied target and updates according to the remaining gap. We further provide a local first-order analysis showing that the RSPO update admits a finite centered-score target and is first-order equivalent to a matched quadratic target objective.
- We validate RSPO on mathematical reasoning and planning benchmarks. Across these settings, RSPO improves diffusion language models over existing dLLM RL baselines, supporting our key observations and demonstrating the superior performance of the proposed algorithm.

## 2. Preliminaries

**Masked Diffusion Language Models.** Masked diffusion models (MDMs) (Sahoo et al., 2024) generate sequences by reversing a gradual masking process. Given a prompt  $q$  and a clean response  $y = (y^1, \dots, y^L)$ , the forward process corrupts  $y$  into  $z_t = (z_t^1, \dots, z_t^L)$  at timestep  $t \in [0, 1]$ , while keeping  $q$  fixed. Each token is independently replaced by the mask token [MASK], whose one-hot representation is denoted by  $m$ :  $q_{t|0}(z_t^i | y^i) = \text{Cat}(z_t^i; \alpha_t y^i + (1 - \alpha_t)m)$ ,

where  $\text{Cat}$  denotes a categorical distribution and  $\alpha_t$  is a decreasing noise schedule with  $\alpha_0 = 1$  and  $\alpha_1 = 0$ . Thus  $\mathbf{z}_0 = \mathbf{y}$  and  $\mathbf{z}_1$  is fully masked; for the linear schedule  $\alpha_t = 1 - t$ , each token is masked with probability  $t$ . The reverse process is parameterized by a denoising model  $\pi_\theta$ . For  $0 \leq s < t \leq 1$ , the transition from  $\mathbf{z}_t$  to  $\mathbf{z}_s$  keeps unmasked tokens fixed and predicts masked tokens from the corrupted context:

$$p_\theta(z_s^i | \mathbf{z}_t, \mathbf{q}) = \begin{cases} \text{Cat}(z_s^i; z_t^i), & z_t^i \neq \mathbf{m}, \\ \text{Cat}\left(z_s^i; \frac{(1 - \alpha_s)\mathbf{m} + (\alpha_s - \alpha_t)\pi_\theta(\cdot | \mathbf{z}_t, \mathbf{q})}{1 - \alpha_t}\right), & z_t^i = \mathbf{m}. \end{cases}$$

Generation starts from a fully masked response and iteratively denoises masked positions, often in parallel. Since MDMs do not use an autoregressive factorization, their likelihood is commonly optimized through an evidence lower bound (ELBO) (Ou et al.; Shi et al., 2024):

$$\mathcal{L}_{\text{ELBO}}(\mathbf{y} | \mathbf{q}; \theta) = \mathbb{E}_{t \sim \mathcal{U}[0,1], \mathbf{z}_t \sim q_{t|0}(\cdot | \mathbf{y})} \left[ \sum_{i=1}^L w(t) \mathbf{1}[z_t^i = \mathbf{m}] \log \pi_\theta(y^i | \mathbf{z}_t, \mathbf{q}) \right] \leq \log \pi_\theta(\mathbf{y} | \mathbf{q}).$$

Here the model is trained to recover clean tokens only at masked positions. The time-dependent weight  $w(t) = \frac{\alpha_t'}{\alpha_t - 1}$  is determined by the noise schedule, which reduces to  $w(t) = 1/t$  under the linear schedule  $\alpha_t = 1 - t$ .

**Reinforcement Learning with Verifiable Rewards.** In reinforcement learning with verifiable rewards (RLVR), a language model is treated as a policy  $\pi_\theta$  that generates a response  $\mathbf{y}$  conditioned on a prompt  $\mathbf{q}$ . A verifier assigns a scalar reward  $r(\mathbf{y}, \mathbf{q})$  to the completed response, and the goal is to improve the policy toward high-reward outputs while controlling its deviation from a reference policy  $\pi_{\text{ref}}$ . The standard reward-maximization objective is  $J(\theta) = \mathbb{E}_{\mathbf{y} \sim \pi_\theta(\cdot | \mathbf{q})} [r(\mathbf{y}, \mathbf{q})]$ , whose policy-gradient form is

$$\nabla_\theta J(\theta) = \mathbb{E}_{\mathbf{y} \sim \pi_\theta(\cdot | \mathbf{q})} [r(\mathbf{y}, \mathbf{q}) \nabla_\theta \log \pi_\theta(\mathbf{y} | \mathbf{q})].$$

For autoregressive language models, the sequence log-likelihood  $\log \pi_\theta(\mathbf{y} | \mathbf{q})$  can be computed by summing token-level log-probabilities. However, this is no longer tractable for dLLMs, whose generation marginalizes over iterative denoising paths rather than following a fixed left-to-right factorization. Consequently, recent dLLM RL methods use an ELBO-based score derived from the denoising objective as a practical substitute (Zhao et al., 2025; Ou et al., 2025; Wang et al., 2025b;a; Lin et al., 2025)  $\log \pi_\theta(\mathbf{y} | \mathbf{q}) \approx \mathcal{L}_{\text{ELBO}}(\mathbf{y} | \mathbf{q}; \theta)$ . Here the ELBO score is tractable because it only requires evaluating denoising probabilities at sampled timesteps and masked positions, rather than summing over all possible denoising trajectories. **The Intractability Challenge.** This substitution changes what the RL objective can reliably measure. The verifier reward remains a direct task-level signal: it tells us whether a completed response is correct or useful. By contrast, the ELBO-based score is only a tractable proxy for the sequence log-likelihood. It is a lower-bound-style estimate, and in practice it is evaluated through sampled timesteps, sampled corruptions, and token-level denoising predictions. Thus, the reward signal and the model-side score have different reliability: the former can be exact for verifiable tasks, while the latter is approximate and can be noisy.

This distinction matters when the ELBO score is used inside a reward-weighted policy objective. A high reward can increase the weight placed on a sampled response, but the objective does not know whether the corresponding ELBO-based relative score is still too small, already sufficiently large, or large mainly because of estimation noise. As a result, reward weighting provides a direction of preference but not a calibrated target for the relative score itself. The goal of this work is to develop an RLVR objective for dLLMs that uses reliable verifier rewards to calibrate ELBO-based relative scores, rather than merely using rewards as fixed weights on approximate likelihood surrogates.

## 3. Method

### 3.1. Motivation

Section 2 motivates the following objective-design issue: when dLLM RLVR relies on an ELBO-based likelihood surrogate, the surrogate requires a well-defined optimization target. We derive this target from the standard KL-regularized policy-improvement problem. In particular, for a fixed prompt  $\mathbf{q}$  and reference policy  $\pi_{\text{ref}}$ , the solution of the standard KL-regularized policy-improvement problem is

$$\pi^*(\mathbf{y} | \mathbf{q}) = Z(\mathbf{q})^{-1} \pi_{\text{ref}}(\mathbf{y} | \mathbf{q}) \exp(r(\mathbf{y}, \mathbf{q})/\beta).$$

where  $Z(\mathbf{q})$  is the prompt-dependent normalizer. Therefore, the ideal current-reference log-ratio is  $\Delta^*(\mathbf{y}, \mathbf{q}) := r(\mathbf{y}, \mathbf{q})/\beta - \log Z(\mathbf{q})$ . For a group  $\{\mathbf{y}_i\}_{i=1}^G$  sampled from the same prompt, let  $r_i = r(\mathbf{y}_i, \mathbf{q})$  and  $\Delta_i^* = \Delta^*(\mathbf{y}_i, \mathbf{q})$ . Then define  $\bar{\Delta}^*$  by the mean of  $\Delta_i^*$  over all  $i$ 's, then the normalization factor  $Z(\mathbf{q})$  eliminates under centering:

$$\Delta_i^* - \bar{\Delta}^* = \frac{1}{\beta}(r_i - \bar{r}). \quad (1)$$

In practice, we use a group-relative advantage  $\tilde{A}_i$  as the normalized version of the centered reward signal. The scaling by  $1/\beta$  and any group-shared reward normalization are absorbed into  $\tilde{A}_i$ , with  $\sum_{i=1}^G \tilde{A}_i = 0$  for a complete prompt group.<sup>1</sup> Moreover, we introduce a fixed feedback coefficient  $\lambda > 0$  and use  $\tau_i := \tilde{A}_i/\lambda$  as the reward-implied target for the centered relative score. The coefficient  $\lambda$  controls the scale of the target: larger  $\lambda$  makes the target more conservative, while smaller  $\lambda$  asks the model to move farther from the reference for the same group-relative advantage.

It remains to define a tractable relative score for dLLMs. Since exact sequence likelihoods are unavailable, RSPO uses current-reference differences of ELBO scores. Let  $\mathcal{E}_i^\theta = \mathcal{L}_{\text{ELBO}}(\mathbf{y}_i | \mathbf{q}; \theta)$  and  $\mathcal{E}_i^{\text{ref}} = \mathcal{L}_{\text{ELBO}}(\mathbf{y}_i | \mathbf{q}; \text{ref})$ . For completion length  $L_c$ ,

$$\delta_i(\theta) = \frac{\mathcal{E}_i^\theta - \mathcal{E}_i^{\text{ref}}}{L_c}. \quad (2)$$

Positive  $\delta_i(\theta)$  means that the current model assigns a higher ELBO score to response  $\mathbf{y}_i$  than the reference model. If an implementation stores cross-entropy or negative ELBO rather than ELBO, the current-reference difference must be reversed.

*Remark 3.1* (Sign convention). Eq. (2) must increase with completion likelihood; otherwise high-reward samples receive the wrong update direction.

We then center this score within the same prompt group. For a group  $\mathcal{B}$  of size  $G$ , define  $\bar{\delta}_{\mathcal{B}}(\theta) = N^{-1} \sum_{j \in \mathcal{B}} \delta_j(\theta)$  and  $\hat{\delta}_i(\theta) = \delta_i(\theta) - \text{sg}(\bar{\delta}_{\mathcal{B}}(\theta))$ , where  $\text{sg}(\cdot)$  denotes detachment. Then, The centered score  $\hat{\delta}_i(\theta)$  is RSPO's tractable counterpart of the centered log-ratio in Eq. (1). When the micro-batch advantages sum to zero, the verifier supplies the target  $\hat{\delta}_i(\theta) \rightarrow \tilde{A}_i/\lambda$ , so the update should follow the residual error between the current relative score and this reward-implied target.

### 3.2. Relative Score Policy Optimization

The name *Relative Score Policy Optimization* reflects the two ingredients above. First, the quantity being controlled is the centered current-reference relative score  $\hat{\delta}_i(\theta)$ . Second, the objective keeps the policy-optimization form: a differentiable model-side score is multiplied by a scalar feedback coefficient.

To see the difference from standard advantage weighting, consider a loss micro-batch  $\mathcal{B}$  with  $N = |\mathcal{B}|$ . The advantage-weighted (AW) surrogate uses the group-relative advantage as a fixed coefficient,  $\ell_{\text{AW}}(\theta; \mathcal{B}) = -N^{-1} \sum_{i \in \mathcal{B}} \tilde{A}_i \hat{\delta}_i(\theta)$ . This objective always pushes a positive-advantage response upward, regardless of its current centered relative score. It therefore encodes a preference direction, but not a stopping rule tied to the reward-implied target.

RSPO replaces this fixed coefficient with the remaining calibration gap. Recall that the target relative score is  $\tilde{A}_i/\lambda$ , or equivalently that the scaled residual is  $e_i(\theta) = \tilde{A}_i - \lambda \hat{\delta}_i(\theta)$ . RSPO uses a detached version of this residual as the feedback coefficient,  $w_i = \tilde{A}_i - \lambda \text{sg}(\hat{\delta}_i(\theta))$ , and defines

$$\ell_{\text{RSPO}}(\theta; \mathcal{B}) = -\frac{1}{N} \sum_{i \in \mathcal{B}} w_i \hat{\delta}_i(\theta). \quad (3)$$

This coefficient adapts to the current relative score. If  $\hat{\delta}_i(\theta) < \tilde{A}_i/\lambda$ , then  $w_i > 0$  and the update increases the score. If  $\hat{\delta}_i(\theta) > \tilde{A}_i/\lambda$ , then  $w_i < 0$  and the update decreases the score. Near the target, the coefficient becomes small, so the update naturally weakens. Thus RSPO keeps the familiar policy-gradient-style objective while turning verifier rewards into calibrated feedback toward a finite relative-score target. The exact first-order gradient of Eq. (3) is analyzed in Section 4.

<sup>1</sup>The analysis only requires the advantages in a complete prompt group to sum to zero; any group-shared reward scaling preserves this property.

*Remark 3.2.* A direct squared calibration loss between  $\widehat{\delta}_i(\theta)$  and  $\widetilde{A}_i/\lambda$  would induce the same local first-order coefficient. RSPO deliberately keeps a policy-optimization form instead: the calibration gap is detached and used as a scalar coefficient, while differentiation passes through a single current relative-score factor.

### 3.3. Training Procedure

Operationally, RSPO can be implemented with the same outer loop as standard group-relative RLVR, as illustrated in Figure 1. For each prompt, we sample a group of completions, evaluate them with a verifier, and convert the resulting rewards into group-relative advantages. We then estimate the current-reference ELBO difference for each completion, center these relative scores within the loss micro-batch, and minimize the RSPO loss in Eq. (3). When computing the current and reference ELBO scores for the same completion, we use the same diffusion masks for both models to reduce Monte Carlo variance in their difference. The reference model is kept fixed throughout training, and  $\lambda$  is treated as a fixed feedback-scale hyperparameter. Full score-estimation and algorithmic details are deferred to Appendix B.

## 4. Theory

The previous section introduced RSPO as a simple change to the usual group-relative RLVR update: keep the policy-score form, but replace the fixed advantage coefficient with a residual to a reward-implied relative-score target. This section explains what that change means mathematically by analyzing the local update RSPO applies in one backward pass. The point is not to prove global convergence for dLLMs, and it is not to assume that we can recover the exact sequence likelihood that Section 2 identified as intractable. Instead, we fix the sampled micro-batch: completions, verifier rewards, group-relative advantages, reference ELBO scores, diffusion masks, and the feedback scale  $\lambda$  are all treated as constants. The batch center is detached, as in Eq. (3), and the only differentiated quantity is the current-model score through  $\delta_i(\theta)$ . Under this local view, the theory answers three questions that mirror the design choices in Section 3: what centering removes, what coefficient multiplies the current score derivative, and what finite target the update is trying to reach. Estimator-level details and proofs are deferred to Appendices D and E.

### 4.1. Relative-Score Feedback and Fixed Point

We first make precise why RSPO is a *relative-score* method. The ELBO values used by dLLMs are only tractable surrogates for sequence log-likelihoods, and their absolute level can contain shifts shared by many samples in a batch. Such shared shifts should not determine which completion is preferred within a prompt group. Detached centering removes them in the forward pass, while preserving the per-sample derivative of the current score:  $\sum_{i \in \mathcal{B}} \widehat{\delta}_i = 0$ , and  $\nabla_{\theta} \widehat{\delta}_i = \nabla_{\theta} \delta_i(\theta)$  because the center is not differentiated.

**Lemma 4.1** (Forward centering and zero-sum weights). *For a fixed sampled micro-batch,  $\sum_{i \in \mathcal{B}} \widehat{\delta}_i = 0$  and  $\sum_{i \in \mathcal{B}} w_i = \sum_{i \in \mathcal{B}} \widetilde{A}_i$ . If  $\mathcal{B}$  is a union of complete prompt groups and  $\widetilde{A}_i$  is a zero-sum group-relative advantage, optionally divided by a group-shared reward scale, then the RSPO weights are also zero-sum.*

The lemma is a useful consistency check. When the advantages are group-relative, the feedback term does not introduce a new batch-level push. It continues to compare responses within the same prompt group, which is exactly the comparison structure used to derive the centered target in Section 3.

We next look at the actual gradient of the RSPO loss. Since the residual coefficient is detached, differentiation passes through only one copy of the current relative score.

**Proposition 4.2** (RSPO relative-score feedback gradient). *For fixed  $\lambda$  and a fixed sampled micro-batch,*

$$\nabla_{\theta} \ell_{\text{RSPO}}(\theta; \mathcal{B}) = -\frac{1}{N} \sum_{i \in \mathcal{B}} (\widetilde{A}_i - \lambda \widehat{\delta}_i) \nabla_{\theta} \delta_i(\theta). \quad (4)$$

*The policy-gradient coefficient is the relative-score feedback error, rather than the advantage alone.*

This formula is the core local interpretation of RSPO. The usual advantage-weighted update, as used in policy-gradient-style RL objectives such as PPO or GRPO (Schulman et al., 2017; Shao et al., 2024), multiplies the score derivative by a fixed advantage. RSPO instead multiplies it by  $\widetilde{A}_i - \lambda \widehat{\delta}_i$ . If the centered score is still below its reward-implied target  $\widetilde{A}_i/\lambda$ , the coefficient is positive and the update increases the score. If the score has already passed that target,

the coefficient becomes negative and pulls it back. Near the target, the coefficient is small. So the update is still a policy-score update, but the coefficient now measures how much target error remains.

**Corollary 4.3** (Standard advantage-weighted objective at the reference point). *If the coupled current and reference ELBO scores match for every sampled completion in the micro-batch, then  $\hat{\delta}_i = 0$  for all  $i \in \mathcal{B}$  and  $\nabla_{\theta} \ell_{\text{RSPO}}(\theta; \mathcal{B}) = \nabla_{\theta} \ell_{\text{AW}}(\theta; \mathcal{B})$ . If additionally  $\sum_{i \in \mathcal{B}} \tilde{A}_i = 0$ , then centering does not change the forward value of this standard surrogate.*

This corollary explains why RSPO should be viewed as a modification of advantage-weighted RLVR, not a replacement for the policy-improvement direction. At the reference point, before the current model has moved away from the reference on the sampled completions, RSPO gives the same first-order update as the standard advantage-weighted surrogate. The distinction appears after the relative scores start moving: advantage weighting keeps applying the same coefficient, whereas RSPO asks whether the current score is still below, near, or above its target.

**Proposition 4.4** (Fixed point under zero-sum advantages). *Assume  $\lambda > 0$  and  $\sum_{i \in \mathcal{B}} \tilde{A}_i = 0$ . If*

$$\hat{\delta}_i = \tilde{A}_i / \lambda, \quad i \in \mathcal{B}, \quad (5)$$

*then every sample weight is zero and  $\nabla_{\theta} \ell_{\text{RSPO}}(\theta; \mathcal{B}) = 0$ . Conversely, if the Jacobian of  $\delta(\theta)$  has full row rank on the micro-batch, every first-order stationary point satisfies Eq. (5).*

The fixed point is the formal version of the target-tracking story. A reward advantage is not treated as a coefficient to apply indefinitely. It specifies a finite centered current-reference score, with  $\lambda$  setting the scale of that target. The full-row-rank assumption is only needed for the converse direction; even without it, Eq. (5) is the condition under which all RSPO feedback weights vanish.

## 4.2. First-Order Equivalence of RSPO

The preceding results describe the coefficient used by RSPO. There is still a natural question: should RSPO be understood as a calibration objective, or as a policy-optimization update? The answer is that, locally, it is both. It tracks the same target as a simple quadratic calibration loss, while keeping the implementation in the familiar form of a detached coefficient times a differentiable policy score.

Let  $\delta$  and  $\hat{\delta}$  denote the micro-batch vectors, and write  $\langle a, b \rangle_{\mathcal{B}} = N^{-1} \sum_{i \in \mathcal{B}} a_i b_i$  and  $\|a\|_{\mathcal{B}}^2 = \langle a, a \rangle_{\mathcal{B}}$ . Consider the matched quadratic objective

$$\ell_{\text{quad}, \lambda}(\theta; \mathcal{B}) = -\langle \tilde{A}, \delta \rangle_{\mathcal{B}} + \frac{\lambda}{2} \|\hat{\delta}\|_{\mathcal{B}}^2. \quad (6)$$

When the advantages are zero-sum, completing the square shows that Eq. (6) has the same target  $\hat{\delta}_i = \tilde{A}_i / \lambda$  as Proposition 4.4, up to terms whose derivative vanishes under detached centering; see Appendix D.1. This comparison separates the target being tracked from the computational graph used to track it.

**Theorem 4.5** (First-order equivalence to a matched quadratic objective). *For fixed  $\lambda$  and a fixed sampled micro-batch,*

$$\nabla_{\theta} \ell_{\text{RSPO}}(\theta; \mathcal{B}) = \nabla_{\theta} \ell_{\text{quad}, \lambda}(\theta; \mathcal{B}). \quad (7)$$

*Thus RSPO follows the same first-order update vector as Eq. (6), while detaching one copy of  $\hat{\delta}$  in the feedback term.*

The theorem clarifies the role of the stop-gradient construction. The finite target can be written as a quadratic calibration target, but RSPO implements the corresponding first-order direction as a policy-score update using the quantities already computed during training: a centered current-reference ELBO score and a detached residual coefficient. The equivalence is only first-order. The scalar objectives differ for higher-order differentiation, and this section does not claim a general variance reduction or global convergence guarantee. Instead, it gives the local mechanism that the experiments probe next: whether residual feedback, reference subtraction, and centering improve dLLM RLVR when the likelihood signal is an ELBO-based surrogate.

## 5. Experiments

The experiments are designed to test the two claims suggested by the method and theory. First, if verifier rewards can serve as targets for centered relative scores, then RSPO should be especially useful in settings where the reward is

reliable but sparse: the verifier can judge a completed answer, but it does not provide token-level or denoising-step supervision. Second, if the gains come from the relative-score feedback mechanism, then removing feedback, reference subtraction, or centering should change the optimization behavior in predictable ways. We therefore evaluate RSPO on both planning and mathematical reasoning benchmarks, and then use ablations to isolate the roles of feedback, current-reference scoring, and centering.

## 5.1. Experimental Setup

**Models and Datasets.** We follow the experimental protocols of d1 (Zhao et al., 2025) and wd1 (Tang et al., 2025). We use LLaDA-8B-Instruct (Nie et al.) as the base diffusion language models. Our evaluation covers two categories of reasoning tasks: mathematical reasoning, including GSM8K (Cobbe et al., 2021) and MATH500 (Lightman et al., 2023), and planning, including Countdown (Pan et al., 2025a) and Sudoku (cor). We follow prior work (Zhao et al., 2025; Tang et al., 2025) for train-test splits, reward functions, and evaluation protocols, except for Sudoku. For Sudoku, to avoid train-test leakage, we split the data by puzzle solutions so that the test set contains unseen answers, preventing models from solving test instances through memorization.

**Baselines.** We compare RSPO with representative RL fine-tuning methods for diffusion language models, including d1/Diffu-GRPO (Zhao et al., 2025), VRPO on LLaDA-1.5 (Zhu et al., 2025), wd1 (Tang et al., 2025), SAPO (Xie et al., 2025), and TraceRL (Wang et al., 2025b). These baselines cover both GRPO-style adaptations and recent dLLM-specific policy optimization methods, providing a broad comparison across mathematical and planning benchmarks.

**Training Details.** For RSPO training, we apply Low-Rank Adaptation (LoRA) with rank  $r = 128$  and scaling factor  $\alpha = 64$ , and set  $\lambda = 0.01$ . All experiments are conducted in the zero-shot setting. For both RL rollouts and evaluation, we use the semi-autoregressive confidence-based decoding strategy following LLaDA, d1, and wd1. We apply the same generation setup as d1: the denoising timestep is set to half of the total sequence length, and the sequence is divided into blocks of 32 tokens. At each diffusion step, we unmask the two tokens with the highest confidence, measured by the probability of the sampled token, within the current incomplete block. During RL rollouts, we use a generation length of 256 and a sampling temperature of 0.9 across all benchmarks, except for Sudoku, where the temperature is set to 0.3 following d1. During evaluation, the sampling temperature is set to 0.0. We evaluate checkpoints every 100 training steps and report results from the checkpoint with the highest average test accuracy across generation lengths of 256 and 512.

## 5.2. Main Results

Table 1 reports final accuracy on planning and mathematical reasoning benchmarks. RSPO gives its clearest gains on sparse-reward planning tasks: on Sudoku, it reaches 92.1/90.8 accuracy at generation lengths 256/512, improving over the strongest non-RSPO result by 66.5/65.4 points; on Countdown, it reaches 78.83/73.83, improving by 26.83/17.53 points. These gains are consistent across generation budgets and align with the motivation of RSPO: reliable response-level rewards can serve as calibrated targets for noisy ELBO-based current-reference scores. On mathematical reasoning, the pattern is more mixed but still favorable. RSPO obtains the best length-256 results on both GSM8K and MATH500, and the best length-512 result on MATH500. On GSM8K at length 512, RSPO is below the strongest baseline, indicating that the benefit is task-dependent in high-performing mathematical-reasoning regimes. *Overall, RSPO is strongest on planning benchmarks while remaining competitive on standard mathematical reasoning tasks.*

Figure 2 reports training reward trajectories on the planning benchmarks. On both Countdown and Sudoku, RSPO improves verifier rewards faster than d1 and wd1 and reaches higher reward levels within the plotted training window. This trajectory evidence suggests that RSPO changes the optimization path during training, rather than merely selecting a better final checkpoint. It also matches the local mechanism in Section 4: as a response approaches its reward-implied relative-score target, RSPO reduces or reverses the update coefficient instead of repeatedly applying the same fixed advantage. Additional reasoning reward trajectories are provided in Figure 4 in Appendix G.

## 5.3. Ablation Studies

The ablations focus on Sudoku\_new (Wang et al., 2025a), a sparse correctness-reward planning task, to isolate the three RSPO components introduced in Sections 3 and 4: residual feedback through  $\lambda$ , current-reference subtraction in  $\delta_i$ , and centering of the relative scores.

Table 1. Benchmark accuracy of RSPO and dLLM RL baselines across planning and mathematical reasoning tasks. Columns 256 and 512 denote generation length; best results are in **bold**.

Methods	Sudoku		Countdown		GSM8K		Math500	
	256	512	256	512	256	512	256	512
LLaDA-8B-Instruct (Nie et al.)	16.2	6.0	19.5	16.0	76.7	78.2	32.4	36.2
+ d1 (Zhao et al., 2025)	24.1	15.9	31.3	37.1	78.1	81.2	34.1	39.0
+ VRPO (Zhu et al., 2025)	12.8	9.6	22.3	18.0	80.1	81.5	35.6	34.8
+ wd1 (Tang et al., 2025)	22.0	24.6	51.2	46.1	80.8	82.3	34.4	39.0
+ SAPO (Xie et al., 2025)	20.3	16.1	52.0	56.3	80.6	82.1	33.8	38.4
+ TraceRL (Wang et al., 2025b)	25.6	25.4	50.4	52.6	81.3	<b>82.4</b>	35.6	39.1
+ RSPO (ours)	<b>92.1</b>	<b>90.8</b>	<b>78.83</b>	<b>73.83</b>	<b>82.63</b>	80.55	<b>42.07</b>	<b>40.25</b>

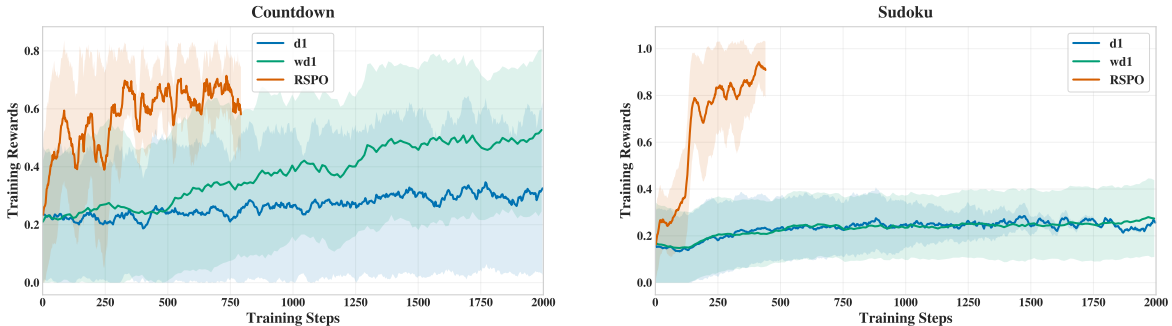


Figure 2. Training reward dynamics of RSPO and baselines on planning benchmarks. Shaded regions indicate variation across runs.

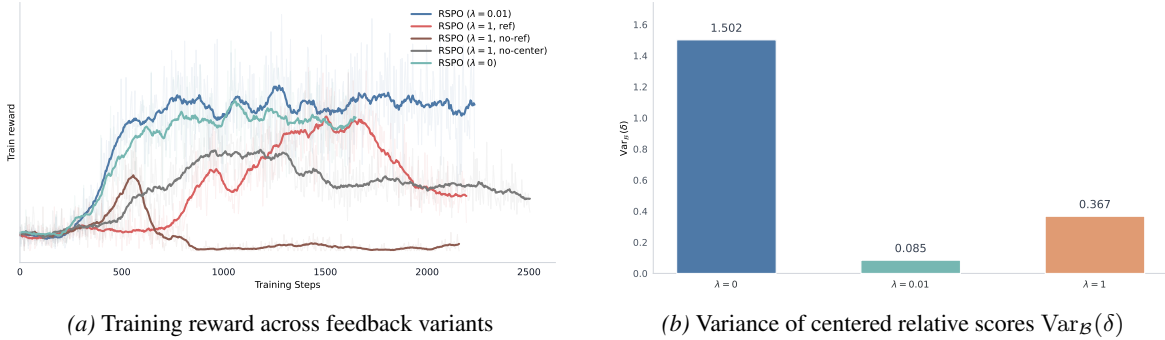
**Relative-Score Feedback.** We first ablate the feedback coefficient. Setting  $\lambda = 0$  removes the residual feedback term and reduces the update to the standard advantage-weighted form over the centered current-reference score. Figure 3a shows that this variant improves early but converges to a lower reward level, supporting Section 3: fixed advantages rank samples but do not define finite relative-score targets. RSPO with  $\lambda = 0.01$  gives the best reward dynamics, whereas stronger feedback at  $\lambda = 1$  is less effective and removing reference-model subtraction at  $\lambda = 1$  severely degrades reward. Thus RSPO works best when feedback calibrates current-reference scores without letting noisy ELBO estimates dominate reward-driven improvement.

**Relative-Score Variance Diagnostics.** We next examine the spread of the centered current-reference scores. Figure 3b reports  $\text{Var}_{\mathcal{B}}(\delta)$ , the batch variance of the ELBO-based relative scores after centering. Because centering removes shared score offsets, this quantity measures relative spread within the response group rather than a batch-level shift. Without feedback ( $\lambda = 0$ ),  $\text{Var}_{\mathcal{B}}(\delta)$  is large; at  $\lambda = 1$ , the spread is reduced but rewards are worse. The moderate setting  $\lambda = 0.01$  gives both the strongest reward dynamics and the smallest measured  $\text{Var}_{\mathcal{B}}(\delta)$ , consistent with Proposition 4.2. Thus the diagnostic supports feedback calibration, but remains only a proxy related to Appendix D.2 and is not used to adapt  $\lambda$ .

**Batch-Mean Score-Offset Diagnostics.** Finally, we check the centering operation itself. RSPO applies feedback to  $\hat{\delta}_i$ , not to the raw ELBO difference  $\delta_i$ , so shared ELBO-score shifts do not affect within-group preferences. Table 2 reports that the mean absolute batch-mean offset stays near zero with centering, including  $1.32 \times 10^{-9}$  at  $\lambda = 0.01$  and  $9.91 \times 10^{-9}$  at  $\lambda = 1$ , while removing centering raises it to  $1.14 \times 10^{-1}$  with much larger variation. This confirms the centering condition in Lemma 4.1: RSPO feedback should act on relative likelihood changes within the response group, rather than on a shared ELBO-score shift.

## 6. Conclusion and Limitations

This paper introduced RSPO, a relative-score objective for RLVR post-training of diffusion language models. The central idea is to use group-relative verifier rewards not only to rank sampled completions, but also to define where their centered current-reference scores should move. RSPO implements this idea by converting the gap between the reward-implied target and the current centered ELBO score into a detached feedback coefficient for a policy-optimization update. This gives a direct way to use noisy dLLM likelihood surrogates without treating reward advantages as indefinitely



(a) Training reward across feedback variants

(b) Variance of centered relative scores  $\text{Var}_B(\delta)$ 

Figure 3. Sudoku ablations of RSPO feedback components. The two panels compare reward dynamics and centered relative-score variance under different feedback settings.

Table 2. Batch-mean relative-score offset in Sudoku ablations. The table checks whether centering removes shared ELBO-score shifts across feedback and reference settings.

$\lambda$	Centering	Reference	Mean  offset	Std.	Range
0	✓	✓	$1.80 \times 10^{-8}$	$2.40 \times 10^{-8}$	$[-1.14, 1.23] \times 10^{-7}$
0.01	✓	✓	$1.32 \times 10^{-9}$	$2.41 \times 10^{-9}$	$[-1.68, 1.49] \times 10^{-8}$
1	✓	✓	$9.91 \times 10^{-9}$	$1.28 \times 10^{-8}$	$[-5.59, 5.59] \times 10^{-8}$
1	✓	–	$3.94 \times 10^{-8}$	$4.94 \times 10^{-8}$	$[-1.49, 1.49] \times 10^{-7}$
1	–	✓	$1.14 \times 10^{-1}$	$1.42 \times 10^{-1}$	$[-4.11, 5.08] \times 10^{-1}$

applied update weights. Empirically, RSPO shows its clearest gains on sparse-reward planning tasks and remains competitive on mathematical reasoning benchmarks, with ablations supporting the roles of residual feedback, reference subtraction, and centering.

We also note several limitations that point to useful future directions. The current configurations of RSPO keep RSPO simple and close to standard group-relative RLVR, but they leave open how to adapt  $\lambda$  automatically or vary the reference policy during training. The theory in Section 4 is local and first-order: it explains the implemented update by matching RSPO to a quadratic target objective at the gradient level, but does not establish global convergence or higher-order equivalence. It is interesting and important to develop a global theory for dLLM RL with ELBO-based score estimates. Finally, the score-spread diagnostic  $\text{Var}_B(\delta)$  is only a local proxy rather than an exact KL estimator, and implementations must respect the ELBO sign convention in Remark 3.1; more precise trust-region diagnostics and robust score-estimation procedures are promising directions for future work.

## References

- Arel’s sudoku generator. <https://www.ocf.berkeley.edu/~arel/sudoku/main.html>. Accessed: 2026-05-03.
- AI, I., Bie, T., Chen, H., Chen, T., Cheng, Z., Cui, L., Gan, K., Huang, Z., Lan, Z., Li, H., et al. Llada2. 0-uni: Unifying multimodal understanding and generation with diffusion large language model. *arXiv preprint arXiv:2604.20796*, 2026.
- Arriola, M., Gokaslan, A., Chiu, J. T., Yang, Z., Qi, Z., Han, J., Sahoo, S. S., and Kuleshov, V. Block diffusion: Interpolating between autoregressive and diffusion language models. In *The Thirteenth International Conference on Learning Representations*.
- Austin, J., Johnson, D. D., Ho, J., Tarlow, D., and Van Den Berg, R. Structured denoising diffusion models in discrete state-spaces. *Advances in neural information processing systems*, 34:17981–17993, 2021.
- Bie, T., Cao, M., Chen, K., Du, L., Gong, M., Gong, Z., Gu, Y., Hu, J., Huang, Z., Lan, Z., et al. Llada2. 0: Scaling up diffusion language models to 100b. *arXiv preprint arXiv:2512.15745*, 2025.
- Campbell, A., Benton, J., De Bortoli, V., Rainforth, T., Deligiannidis, G., and Doucet, A. A continuous time framework for discrete denoising models. *Advances in Neural Information Processing Systems*, 35:28266–28279, 2022.
- Cheng, S., Bian, Y., Liu, D., Zhang, L., Yao, Q., Tian, Z., Wang, W., Guo, Q., Chen, K., Qi, B., et al. Sdar: A synergistic diffusion-autoregression paradigm for scalable sequence generation. *arXiv preprint arXiv:2510.06303*, 2025.
- Cobbe, K., Kosaraju, V., Bavarian, M., Chen, M., Jun, H., Kaiser, L., Plappert, M., Tworek, J., Hilton, J., Nakano, R., et al. Training verifiers to solve math word problems. *arXiv preprint arXiv:2110.14168*, 2021.
- Dhariwal, P. and Nichol, A. Diffusion models beat gans on image synthesis. *Advances in neural information processing systems*, 34:8780–8794, 2021.
- Guo, D., Yang, D., Zhang, H., Song, J., Wang, P., Zhu, Q., Xu, R., Zhang, R., Ma, S., Bi, X., et al. Deepseek-r1: Incentivizing reasoning capability in llms via reinforcement learning. *arXiv preprint arXiv:2501.12948*, 2025.
- HaCohen, Y., Chiprut, N., Brazowski, B., Shalem, D., Moshe, D., Richardson, E., Levin, E., Shiran, G., Zabari, N., Gordon, O., et al. Ltx-video: Realtime video latent diffusion. *arXiv preprint arXiv:2501.00103*, 2024.
- HaCohen, Y., Brazowski, B., Chiprut, N., Bitterman, Y., Kvochko, A., Berkowitz, A., Shalem, D., Lifschitz, D., Moshe, D., Porat, E., et al. Ltx-2: Efficient joint audio-visual foundation model. *arXiv preprint arXiv:2601.03233*, 2026.
- He, S., Xia, T., Zhou, X., and Wei, H. Response-level rewards are all you need for online reinforcement learning in llms: A mathematical perspective. *arXiv preprint arXiv:2506.02553*, 2025.
- Ho, J., Jain, A., and Abbeel, P. Denoising diffusion probabilistic models. *Advances in neural information processing systems*, 33:6840–6851, 2020.
- Hu, E. J., Shen, Y., Wallis, P., Allen-Zhu, Z., Li, Y., Wang, S., Wang, L., Chen, W., et al. Lora: Low-rank adaptation of large language models. *Iclr*, 1(2):3, 2022.
- Jaech, A., Kalai, A., Lerer, A., Richardson, A., El-Kishky, A., Low, A., Helyar, A., Madry, A., Beutel, A., Carney, A., et al. Openai o1 system card. *arXiv preprint arXiv:2412.16720*, 2024.
- Lightman, H., Kosaraju, V., Burda, Y., Edwards, H., Baker, B., Lee, T., Leike, J., Schulman, J., Sutskever, I., and Cobbe, K. Let’s verify step by step. In *The twelfth international conference on learning representations*, 2023.
- Lin, N., Zhang, J., Hou, L., and Li, J. Boundary-guided policy optimization for memory-efficient rl of diffusion large language models. *arXiv preprint arXiv:2510.11683*, 2025.
- Loshchilov, I. and Hutter, F. Decoupled weight decay regularization. *arXiv preprint arXiv:1711.05101*, 2017.
- Lou, A., Meng, C., and Ermon, S. Discrete diffusion modeling by estimating the ratios of the data distribution. In *Forty-first International Conference on Machine Learning*.

- Nie, S., Zhu, F., You, Z., Zhang, X., Ou, J., Hu, J., ZHOU, J., Lin, Y., Wen, J.-R., and Li, C. Large language diffusion models. In *The Thirty-ninth Annual Conference on Neural Information Processing Systems*.
- Ou, J., Nie, S., Xue, K., Zhu, F., Sun, J., Li, Z., and Li, C. Your absorbing discrete diffusion secretly models the conditional distributions of clean data. In *The Thirteenth International Conference on Learning Representations*.
- Ou, J., Han, J., Xu, M., Xu, S., Xie, J., Ermon, S., Wu, Y., and Li, C. Principled rl for diffusion llms emerges from a sequence-level perspective. *arXiv preprint arXiv:2512.03759*, 2025.
- Ouyang, L., Wu, J., Jiang, X., Almeida, D., Wainwright, C. L., Mishkin, P., Zhang, C., Agarwal, S., Slama, K., Ray, A., et al. Training language models to follow instructions with human feedback. In *Advances in Neural Information Processing Systems*, 2022.
- Pan, J., Zhang, J., Wang, X., Yuan, L., Peng, H., and Suhr, A. Tinyzero. <https://github.com/Jiayi-Pan/TinyZero>, 2025a. Accessed: 2025-01-24.
- Pan, L., Tao, S., Zhai, Y., Fu, Z., Fang, L., He, M., Zhang, L., Liu, Z., Ding, B., Liu, A., et al. d-treerpo: Towards more reliable policy optimization for diffusion language models. *arXiv preprint arXiv:2512.09675*, 2025b.
- Sahoo, S. S., Arriola, M., Schiff, Y., Gokaslan, A., Marroquin, E., Chiu, J. T., Rush, A., and Kuleshov, V. Simple and effective masked diffusion language models. *Advances in Neural Information Processing Systems*, 37:130136–130184, 2024.
- Schulman, J., Wolski, F., Dhariwal, P., Radford, A., and Klimov, O. Proximal policy optimization algorithms. *arXiv preprint arXiv:1707.06347*, 2017.
- Seedance, T., Chen, D., Chen, L., Chen, X., Chen, Y., Chen, Z., Chen, Z., Cheng, F., Cheng, T., Cheng, Y., et al. Seedance 2.0: Advancing video generation for world complexity. *arXiv preprint arXiv:2604.14148*, 2026.
- Shao, Z., Wang, P., Zhu, Q., Xu, R., Song, J., Bi, X., Zhang, H., Zhang, M., Li, Y., et al. Deepseekmath: Pushing the limits of mathematical reasoning in open language models. *arXiv preprint arXiv:2402.03300*, 2024.
- Shi, J., Han, K., Wang, Z., Doucet, A., and Titsias, M. Simplified and generalized masked diffusion for discrete data. *Advances in neural information processing systems*, 37:103131–103167, 2024.
- Song, Y., Sohl-Dickstein, J., Kingma, D. P., Kumar, A., Ermon, S., and Poole, B. Score-based generative modeling through stochastic differential equations. In *International Conference on Learning Representations*.
- Sutton, R. S., Barto, A. G., et al. *Reinforcement learning: An introduction*, volume 1. MIT press Cambridge, 1998.
- Tang, X., Dolga, R., Yoon, S., and Bogunovic, I. wd1: Weighted policy optimization for reasoning in diffusion language models. *arXiv preprint arXiv:2507.08838*, 2025.
- von Werra, L., Belkada, Y., Tunstall, L., Beeching, E., Thrush, T., Lambert, N., Huang, S., Rasul, K., and Gallouédec, Q. TRL: Transformers Reinforcement Learning, 2020. URL <https://github.com/huggingface/trl>.
- Wan, T., Wang, A., Ai, B., Wen, B., Mao, C., Xie, C.-W., Chen, D., Yu, F., Zhao, H., Yang, J., et al. Wan: Open and advanced large-scale video generative models. *arXiv preprint arXiv:2503.20314*, 2025.
- Wang, C., Rashidinejad, P., Su, D., Jiang, S., Wang, S., Zhao, S., Zhou, C., Shen, S. Z., Chen, F., Jaakkola, T., et al. Spg: Sandwiched policy gradient for masked diffusion language models. *arXiv preprint arXiv:2510.09541*, 2025a.
- Wang, Y., Yang, L., Li, B., Tian, Y., Shen, K., and Wang, M. Revolutionizing reinforcement learning framework for diffusion large language models. *arXiv preprint arXiv:2509.06949*, 2025b.
- Xie, S., Kong, L., Song, X., Dong, X., Chen, G., Xing, E. P., and Zhang, K. Step-aware policy optimization for reasoning in diffusion large language models. *arXiv preprint arXiv:2510.01544*, 2025.
- Yang, L., Tian, Y., Li, B., Zhang, X., Shen, K., Tong, Y., and Wang, M. Mmada: Multimodal large diffusion language models. *arXiv preprint arXiv:2505.15809*, 2025.

- Ye, J., Xie, Z., Zheng, L., Gao, J., Wu, Z., Jiang, X., Li, Z., and Kong, L. Dream 7b: Diffusion large language models. *arXiv preprint arXiv:2508.15487*, 2025.
- You, Z., Nie, S., Zhang, X., Hu, J., Zhou, J., Lu, Z., Wen, J.-R., and Li, C. Llada-v: Large language diffusion models with visual instruction tuning. *arXiv preprint arXiv:2505.16933*, 2025.
- You, Z., Zhang, X., Zhou, J., Li, C., and Wen, J.-R. Llada-o: An effective and length-adaptive omni diffusion model. *arXiv preprint arXiv:2603.01068*, 2026.
- Yu, Q., Zhang, Z., Zhu, R., Yuan, Y., Zuo, X., Yue, Y., Dai, W., Fan, T., Liu, G., Liu, L., et al. Dapo: An open-source llm reinforcement learning system at scale. *arXiv preprint arXiv:2503.14476*, 2025.
- Zhao, S., Gupta, D., Zheng, Q., and Grover, A. d1: Scaling reasoning in diffusion large language models via reinforcement learning. *arXiv preprint arXiv:2504.12216*, 2025.
- Zheng, C., Liu, S., Li, M., Chen, X.-H., Yu, B., Gao, C., Dang, K., Liu, Y., Men, R., Yang, A., et al. Group sequence policy optimization. *arXiv preprint arXiv:2507.18071*, 2025a.
- Zheng, Z., Peng, X., Lou, Y., Shen, C., Young, T., Guo, X., Wang, B., Xu, H., Liu, H., Jiang, M., et al. Open-sora 2.0: Training a commercial-level video generation model in \$200 k. *arXiv preprint arXiv:2503.09642*, 2025b.
- Zhong, J., Wang, K., Ding, D., Feng, Z., Bai, H., Xiang, Y., Sun, J., and Xu, Q. Stabilizing reinforcement learning for diffusion language models. *arXiv preprint arXiv:2603.06743*, 2026.
- Zhu, F., Wang, R., Nie, S., Zhang, X., Wu, C., Hu, J., Zhou, J., Chen, J., Lin, Y., Wen, J.-R., et al. Llada 1.5: Variance-reduced preference optimization for large language diffusion models. *arXiv preprint arXiv:2505.19223*, 2025.

## Contents

### A. Notation Summary

Table 3. Notation used in the RSPO formulation.

Symbol	Meaning
$q$	Prompt
$y$	Generic completion used in policy-improvement discussion
$o_i$	The $i$ -th sampled completion in a group
$G$	Group size
$K$	Number of ELBO Monte Carlo mask samples per completion
$\mathcal{B}$	Current loss micro-batch
$N$	Micro-batch size, $ \mathcal{B} $
$\pi_\theta$	Current model or policy parameterized by $\theta$
$\pi_{\text{ref}}$	Frozen reference model or policy
$\pi^*$	Ideal KL-regularized improved policy in the motivating calculation
$\beta$	KL temperature in the motivating policy-improvement calculation
$Z(q)$	Prompt-level normalizer in the motivating policy-improvement calculation
$r_i$	Reward of completion $o_i$
$\bar{r}$	Group mean reward, $G^{-1} \sum_i r_i$
$\sigma_r$	Group reward standard deviation for the normalized advantage variant
$s_q$	Prompt-level scale used by the normalized advantage variant
$\tilde{A}_i$	Zero-sum group-relative advantage passed to the loss, optionally divided by $\sigma_r + 10^{-4}$
$L_c$	Number of completion tokens used to normalize score differences
$\hat{\mathcal{E}}_\theta$	Likelihood-oriented ELBO or score estimate for the current model
$\hat{\mathcal{E}}_{\text{ref}}$	Matched reference-model ELBO or score estimate
$M_{i,k}$	The $k$ -th sampled diffusion mask set for completion $o_i$
$\delta_i$	Per-token ELBO log-ratio surrogate, $(\hat{\mathcal{E}}_\theta(o_i; q) - \hat{\mathcal{E}}_{\text{ref}}(o_i; q))/L_c$
$\bar{\delta}_{\mathcal{B}}$	Current micro-batch mean of $\delta$
$\hat{\delta}_i$	Centered relative-score surrogate, $\delta_i - \text{sg}(\bar{\delta}_{\mathcal{B}})$
$\lambda$	Fixed feedback coefficient; the main RSPO experiments set $\lambda = 0.01$ unless otherwise stated
$\text{sg}(\cdot)$	Detachment operator
$e_i$	Remaining calibration gap, $\tilde{A}_i - \lambda \hat{\delta}_i$
$w_i$	RSPO weight, $\tilde{A}_i - \lambda \text{sg}(\hat{\delta}_i)$
$\ell_{\text{AW}}$	Standard group-relative advantage-weighted surrogate
$\ell_{\text{RSPO}}$	RSPO objective
$\ell_{\text{quad}, \lambda}$	Matched quadratic objective with RSPO’s relative-score target
$\langle a, b \rangle_{\mathcal{B}}$	Micro-batch inner product, $N^{-1} \sum_{i \in \mathcal{B}} a_i b_i$
$\ a\ _{\mathcal{B}}$	Micro-batch norm, with $\ a\ _{\mathcal{B}}^2 = \langle a, a \rangle_{\mathcal{B}}$
$\text{Var}_{\mathcal{B}}(\delta)$	Within-micro-batch variance of the relative-score surrogate

### B. Training Procedure and Score Estimation

This appendix specifies the score-estimation and training details used by RSPO. The main text abstracts these details into an ELBO-based current-reference relative score  $\delta_i$  for each completion, micro-batch centering, and the relative-score feedback loss in Eq. (3).

### B.1. Masked diffusion ELBO estimator

For a prompt-completion pair  $(q, o_i)$ , let  $L_c$  be the number of completion tokens. We estimate a likelihood-oriented sequence score by sampling  $K$  masks and averaging denoising log-probabilities on completion tokens:

$$\widehat{\mathcal{E}}_{\theta}(o_i; q) = \frac{1}{K} \sum_{k=1}^K \frac{L_c}{|M_{i,k} \cap o_i|} \sum_{t \in M_{i,k} \cap o_i} \log p_{\theta}(x_t | x_{\setminus M_{i,k}}, q, M_{i,k}). \quad (8)$$

Here  $M_{i,k}$  is the  $k$ -th sampled mask set for completion  $o_i$ ,  $x_t$  is the token at position  $t$ ,  $x_{\setminus M_{i,k}}$  denotes the sequence context outside the masked positions, and  $p_{\theta}$  is the current model’s denoising distribution. Masks with no completion token are resampled or omitted from the Monte Carlo average. For each completion, the same mask set is reused for the current model and the reference model, giving the coupled per-token score difference

$$\delta_i = \frac{\widehat{\mathcal{E}}_{\theta}(o_i; q) - \widehat{\mathcal{E}}_{\text{ref}}(o_i; q)}{L_c}.$$

The factor  $L_c/|M_{i,k} \cap o_i|$  makes  $\widehat{\mathcal{E}}_{\theta}$  a sequence-level completion score, so the division by  $L_c$  in  $\delta_i$  gives the per-token score difference used in Eq. (2). If an implementation stores a negative ELBO or cross-entropy instead, the sign must follow Remark 3.1.

### B.2. Algorithm

---

#### Algorithm 1 RSPO for diffusion language model reinforcement learning

---

**Require:** Current dLLM  $\pi_{\theta}$ , fixed reference dLLM  $\pi_{\text{ref}}$ , reward function  $R$ , group size  $G$ , ELBO samples  $K$ , feedback coefficient  $\lambda > 0$ .

- 1: **for** each training step **do**
  - 2:   Sample prompts and generate  $G$  completions  $\{o_i\}_{i=1}^G$  per prompt from  $\pi_{\theta}$ .
  - 3:   Compute rewards  $r_i = R(q, o_i)$  and group-relative advantages within each prompt group.
  - 4:   Optionally normalize the advantages within the group to obtain  $\tilde{A}_i$ .
  - 5:   Estimate  $\widehat{\mathcal{E}}_{\theta}(o_i; q)$  and  $\widehat{\mathcal{E}}_{\text{ref}}(o_i; q)$  with the shared masks in Eq. (8).
  - 6:   Form  $\delta_i = (\widehat{\mathcal{E}}_{\theta}(o_i; q) - \widehat{\mathcal{E}}_{\text{ref}}(o_i; q))/L_c$  and center it within the loss micro-batch:  $\widehat{\delta}_i = \delta_i - \text{sg}(\bar{\delta}_{\mathcal{B}})$ .
  - 7:   Set the relative-score feedback weight  $w_i = \tilde{A}_i - \lambda \text{sg}(\widehat{\delta}_i)$  and descend  $\ell_{\text{RSPO}}(\theta; \mathcal{B}) = -N^{-1} \sum_{i \in \mathcal{B}} w_i \widehat{\delta}_i$ .
  - 8:   Record  $\text{Var}_{\mathcal{B}}(\delta)$  as a stability diagnostic.
  - 9: **end for**
- 

## C. Related Work

**Diffusion Language Models.** Diffusion models were originally developed as a powerful generative framework for continuous data, especially high-fidelity image generation (Song et al., 2020; Ho et al., 2020; Dhariwal & Nichol, 2021) and video generation (Wan et al., 2025; Seedance et al., 2026; HaCohen et al., 2024; 2026; Zheng et al., 2025b). More recently, this paradigm has been extended to discrete text generation, giving rise to diffusion language models that generate token sequences through iterative denoising (Austin et al., 2021; Campbell et al., 2022; Lou et al.; Sahoo et al., 2024; Shi et al., 2024). Unlike autoregressive models, dLLMs are not tied to a fixed left-to-right generation order, enabling bidirectional context modeling, flexible-order refinement, and parallel decoding. Recent large-scale dLLMs, such as LLaDA (Nie et al.; Bie et al., 2025), Dream (Ye et al., 2025), multimodal diffusion language models (Yang et al., 2025; You et al., 2025; 2026; AI et al., 2026), and block-diffusion variants (Arriola et al.; Cheng et al., 2025), have demonstrated competitive generation quality and promising efficiency compared with autoregressive models. These advances motivate post-training methods that can further improve dLLMs on reasoning, alignment, and task-specific objectives.

**Reinforcement Learning for Diffusion Language Models.** RL for dLLMs is challenging because PPO/GRPO-style objectives rely on policy log-probabilities or importance ratios that are tractable for autoregressive models but unavailable for diffusion generation. Existing methods mainly differ in how they approximate or avoid these ratios. diffu-GRPO (Zhao et al., 2025) uses one-step mean-field log-probability estimates to adapt GRPO to masked

dLLMs, while wd1 (Tang et al., 2025) removes explicit importance ratios by reformulating policy optimization as a weighted log-likelihood objective. ELBO-based methods approximate likelihoods or log-ratios with denoising objectives, including VRPO in LLaDA 1.5 (Zhu et al., 2025), sequence-level ELBO policy optimization (Ou et al., 2025), sandwiched evidence-bound policy gradients (Wang et al., 2025a), and memory-efficient large Monte Carlo (large-MC) ELBO optimization (Lin et al., 2025). Other works exploit diffusion trajectories or step structure, such as TraceRL (Wang et al., 2025b) and d-TreeRPO (Pan et al., 2025b), or stabilize noisy ratio-based training with clipping and self-normalization (Zhong et al., 2026). Our work is complementary: instead of changing the likelihood surrogate, trajectory estimator, or clipping rule, we derive the centered log-ratio target implied by KL-regularized group-relative policy improvement and optimize noisy ELBO log-ratios toward this target with relative-score feedback.

## D. Additional Theory Details

### D.1. Coefficient convention for the matched quadratic objective

Expanding the RSPO loss makes the coefficient convention explicit:

$$\ell_{\text{RSPO}}(\theta; \mathcal{B}) = -\langle \tilde{A}, \hat{\delta} \rangle_{\mathcal{B}} + \lambda \langle \text{sg}(\hat{\delta}), \hat{\delta} \rangle_{\mathcal{B}}. \quad (9)$$

The forward value of the detached feedback term in Eq. (9) is  $\lambda \|\hat{\delta}\|_{\mathcal{B}}^2$ , but its first-order gradient is the gradient of  $\frac{\lambda}{2} \|\hat{\delta}\|_{\mathcal{B}}^2$ . If one removes the detachment operator from the written scalar expression  $\lambda \langle \text{sg}(\hat{\delta}), \hat{\delta} \rangle_{\mathcal{B}}$  without changing the coefficient, the penalty gradient doubles. Therefore, comparisons with a fully differentiable matched quadratic objective use the first-order matched coefficient, i.e.,  $\frac{\lambda}{2} \|\hat{\delta}\|_{\mathcal{B}}^2$ .

Let  $c_{\mathcal{B}} = \text{sg}(N^{-1} \sum_{j \in \mathcal{B}} \delta_j)$  be the detached micro-batch center and let  $\mathbf{1}$  be the all-ones vector, so  $\hat{\delta} = \delta - c_{\mathcal{B}} \mathbf{1}$ . Completing the square in Eq. (6) gives the local relative-score target form

$$\ell_{\text{quad}, \lambda}(\theta; \mathcal{B}) = \frac{\lambda}{2} \|\hat{\delta}\|_{\mathcal{B}}^2 - \frac{1}{\lambda} \tilde{A} \|\hat{\delta}\|_{\mathcal{B}}^2 - \frac{1}{2\lambda} \|\tilde{A}\|_{\mathcal{B}}^2 - c_{\mathcal{B}} \langle \tilde{A}, \mathbf{1} \rangle_{\mathcal{B}}. \quad (10)$$

When the micro-batch advantages sum to zero, the final term vanishes and the visible target is  $\hat{\delta} = \tilde{A}/\lambda$ . If the advantages are not zero-sum, the final term is not a forward-value constant, but it has zero derivative in the local backward pass because  $c_{\mathcal{B}}$  is detached.

Theorem 4.5 is first-order only. RSPO is not the same scalar objective as Eq. (6) for higher-order differentiation, and the theorem does not by itself establish lower gradient variance. The precise claim is that the update uses the relative-score feedback coefficient  $(\tilde{A}_i - \lambda \hat{\delta}_i)$  while preventing one copy of the noisy ELBO difference from contributing a differentiable path in the backward graph. Variance reduction is therefore an empirical property to be evaluated in experiments.

### D.2. Local KL interpretation of the monitored variance

The centered log-ratio variance is interpreted as a local divergence proxy rather than as an exact KL estimator. The following second-order calculation gives its scale.

**Proposition D.1** (Second-order KL proxy). *Let  $P$  and  $Q$  be two nearby completion distributions for the same prompt with common support, let  $Y$  denote a completion sampled from the distribution indicated in the expectation, and define*

$$\delta(y) = \log \frac{Q(y)}{P(y)}.$$

*Assume  $\delta$  is uniformly small, where  $\|\delta\|_{\infty}$  denotes the supremum norm. Then*

$$\text{KL}(P\|Q) = \frac{1}{2} \text{Var}_{Y \sim P}(\delta(Y)) + O(\|\delta\|_{\infty}^3),$$

*and, to the same second order,*

$$\text{KL}(Q\|P) = \frac{1}{2} \text{Var}_{Y \sim P}(\delta(Y)) + O(\|\delta\|_{\infty}^3).$$

*Consequently, when the current and reference policies are close and the sampled micro-batch is representative of the local completion distribution, the monitored quantity  $\text{Var}_{\mathcal{B}}(\delta)$  is a second-order proxy for approximately twice a local KL divergence.*

Proposition D.1 supports monitoring  $\text{Var}_{\mathcal{B}}(\delta)$  as a trust-region-like diagnostic, but the evaluated method does not use it to update  $\lambda$ . It does not imply that a finite sampled micro-batch variance is an unbiased or exact estimator of either forward or reverse KL. In particular, the approximation depends on the sampling distribution, the sign convention for  $\delta$ , and the current policy being close to the reference policy.

## E. Additional Proof Details

### E.1. Proof of Lemma 4.1

The forward-centered surrogate is

$$\widehat{\delta}_i = \delta_i - c_{\mathcal{B}}, \quad c_{\mathcal{B}} = \text{sg}\left(\frac{1}{N} \sum_{j \in \mathcal{B}} \delta_j\right).$$

The detachment operator changes derivatives but not forward values, so in the forward pass

$$\begin{aligned} \sum_{i \in \mathcal{B}} \widehat{\delta}_i &= \sum_{i \in \mathcal{B}} \delta_i - \sum_{i \in \mathcal{B}} \frac{1}{N} \sum_{j \in \mathcal{B}} \delta_j \\ &= \sum_{i \in \mathcal{B}} \delta_i - N \cdot \frac{1}{N} \sum_{j \in \mathcal{B}} \delta_j = 0. \end{aligned}$$

Using  $w_i = \widetilde{A}_i - \lambda \text{sg}(\widehat{\delta}_i)$  and again using that detachment preserves forward values,

$$\begin{aligned} \sum_{i \in \mathcal{B}} w_i &= \sum_{i \in \mathcal{B}} \widetilde{A}_i - \lambda \sum_{i \in \mathcal{B}} \text{sg}(\widehat{\delta}_i) \\ &= \sum_{i \in \mathcal{B}} \widetilde{A}_i - \lambda \sum_{i \in \mathcal{B}} \widehat{\delta}_i = \sum_{i \in \mathcal{B}} \widetilde{A}_i. \end{aligned}$$

If the group-relative advantages sum to zero within each complete prompt group, then any micro-batch that is a union of complete prompt groups has  $\sum_{i \in \mathcal{B}} \widetilde{A}_i = 0$ . Any group-shared reward scaling preserves this zero-sum property, and the weights are zero-sum.

### E.2. Proof of Proposition 4.2

The implemented loss is

$$\ell_{\text{RSPO}}(\theta; \mathcal{B}) = -\frac{1}{N} \sum_{i \in \mathcal{B}} (\widetilde{A}_i - \lambda \text{sg}(\widehat{\delta}_i)) \widehat{\delta}_i.$$

Rewards, advantages, the reference scores,  $\lambda$ ,  $\text{sg}(\widehat{\delta}_i)$ , and the batch center  $c_{\mathcal{B}}$  are constants with respect to the derivative. Since

$$\widehat{\delta}_i = \delta_i - c_{\mathcal{B}}, \quad \nabla_{\theta} c_{\mathcal{B}} = 0,$$

we have  $\nabla_{\theta} \widehat{\delta}_i = \nabla_{\theta} \delta_i$ . Therefore

$$\begin{aligned} \nabla_{\theta} \ell_{\text{RSPO}}(\theta; \mathcal{B}) &= -\frac{1}{N} \sum_{i \in \mathcal{B}} (\widetilde{A}_i - \lambda \text{sg}(\widehat{\delta}_i)) \nabla_{\theta} \widehat{\delta}_i \\ &= -\frac{1}{N} \sum_{i \in \mathcal{B}} (\widetilde{A}_i - \lambda \widehat{\delta}_i) \nabla_{\theta} \delta_i(\theta), \end{aligned}$$

which is the stated formula.

### E.3. Proof of Corollary 4.3

If the current model and reference model have identical coupled ELBO scores on the sampled completions, then

$$\widehat{\mathcal{E}}_{\theta}(o_i; q) = \widehat{\mathcal{E}}_{\text{ref}}(o_i; q) \quad \text{for all } i \in \mathcal{B}.$$

The per-token surrogate definition gives

$$\delta_i = \frac{\widehat{\mathcal{E}}_\theta(o_i; q) - \widehat{\mathcal{E}}_{\text{ref}}(o_i; q)}{L_c} = 0.$$

Therefore the detached batch center is also zero,

$$c_{\mathcal{B}} = \text{sg}\left(\frac{1}{N} \sum_{j \in \mathcal{B}} \delta_j\right) = 0,$$

and hence  $\widehat{\delta}_i = \delta_i - c_{\mathcal{B}} = 0$  for every  $i \in \mathcal{B}$ . The RSPO weight becomes

$$w_i = \widetilde{A}_i - \lambda \text{sg}(\widehat{\delta}_i) = \widetilde{A}_i.$$

Using Proposition 4.2,

$$\nabla_\theta \ell_{\text{RSPO}}(\theta; \mathcal{B}) = -\frac{1}{N} \sum_{i \in \mathcal{B}} \widetilde{A}_i \nabla_\theta \delta_i(\theta).$$

For the standard group-relative advantage-weighted surrogate,

$$\ell_{\text{AW}}(\theta; \mathcal{B}) = -\frac{1}{N} \sum_{i \in \mathcal{B}} \widetilde{A}_i \widehat{\delta}_i(\theta),$$

and  $\nabla_\theta \widehat{\delta}_i = \nabla_\theta \delta_i$  because the center is detached. Hence

$$\nabla_\theta \ell_{\text{AW}}(\theta; \mathcal{B}) = -\frac{1}{N} \sum_{i \in \mathcal{B}} \widetilde{A}_i \nabla_\theta \delta_i(\theta) = \nabla_\theta \ell_{\text{RSPO}}(\theta; \mathcal{B}).$$

If  $\sum_{i \in \mathcal{B}} \widetilde{A}_i = 0$ , then

$$\begin{aligned} -\frac{1}{N} \sum_{i \in \mathcal{B}} \widetilde{A}_i \widehat{\delta}_i &= -\frac{1}{N} \sum_{i \in \mathcal{B}} \widetilde{A}_i (\delta_i - c_{\mathcal{B}}) \\ &= -\frac{1}{N} \sum_{i \in \mathcal{B}} \widetilde{A}_i \delta_i + \frac{c_{\mathcal{B}}}{N} \sum_{i \in \mathcal{B}} \widetilde{A}_i \\ &= -\frac{1}{N} \sum_{i \in \mathcal{B}} \widetilde{A}_i \delta_i. \end{aligned}$$

This proves the forward-value identity.

#### E.4. Proof of Proposition 4.4

If  $\widehat{\delta}_i = \widetilde{A}_i/\lambda$  for all  $i \in \mathcal{B}$ , then

$$w_i = \widetilde{A}_i - \lambda \text{sg}(\widehat{\delta}_i) = \widetilde{A}_i - \lambda \widehat{\delta}_i = 0$$

for every sample in the micro-batch. Substituting  $w_i = 0$  into the loss gradient in Proposition 4.2 gives

$$\nabla_\theta \ell_{\text{RSPO}}(\theta; \mathcal{B}) = -\frac{1}{N} \sum_{i \in \mathcal{B}} 0 \cdot \nabla_\theta \delta_i(\theta) = 0.$$

The zero-sum assumption  $\sum_i \widetilde{A}_i = 0$  is needed for this target to be compatible with the forward centering identity  $\sum_i \widehat{\delta}_i = 0$ :

$$\sum_{i \in \mathcal{B}} \frac{\widetilde{A}_i}{\lambda} = \frac{1}{\lambda} \sum_{i \in \mathcal{B}} \widetilde{A}_i = 0.$$

For the converse, let  $J_\theta = \nabla_\theta \delta(\theta) \in \mathbb{R}^{N \times d}$  be the Jacobian of the micro-batch relative-score vector, where  $d$  is the number of trainable parameters. Proposition 4.2 can be written as

$$\nabla_\theta \ell_{\text{RSPO}}(\theta; \mathcal{B}) = -\frac{1}{N} J_\theta^\top (\tilde{A} - \lambda \hat{\delta}).$$

At a stationary point this gradient is zero, so

$$J_\theta^\top (\tilde{A} - \lambda \hat{\delta}) = 0.$$

If the Jacobian can realize arbitrary first-order changes in the sampled relative-score vector, equivalently if  $J_\theta$  has full row rank on these sampled directions, then the nullspace of  $J_\theta^\top$  is trivial. Hence

$$\tilde{A} - \lambda \hat{\delta} = 0,$$

which is exactly  $\hat{\delta}_i = \tilde{A}_i / \lambda$  for all  $i \in \mathcal{B}$ .

### E.5. Proof of Theorem 4.5

Recall the matched quadratic objective:

$$\ell_{\text{quad}, \lambda}(\theta; \mathcal{B}) = -\langle \tilde{A}, \delta \rangle_{\mathcal{B}} + \frac{\lambda}{2} \|\hat{\delta}\|_{\mathcal{B}}^2, \quad \hat{\delta}_i = \delta_i - c_{\mathcal{B}}, \quad \nabla_\theta c_{\mathcal{B}} = 0.$$

Its gradient is

$$\begin{aligned} \nabla_\theta \ell_{\text{quad}, \lambda} &= -\frac{1}{N} \sum_{i \in \mathcal{B}} \tilde{A}_i \nabla_\theta \delta_i + \frac{\lambda}{2} \nabla_\theta \left( \frac{1}{N} \sum_{i \in \mathcal{B}} \hat{\delta}_i^2 \right) \\ &= -\frac{1}{N} \sum_{i \in \mathcal{B}} \tilde{A}_i \nabla_\theta \delta_i + \lambda \frac{1}{N} \sum_{i \in \mathcal{B}} \hat{\delta}_i \nabla_\theta \hat{\delta}_i \\ &= -\frac{1}{N} \sum_{i \in \mathcal{B}} \tilde{A}_i \nabla_\theta \delta_i + \lambda \frac{1}{N} \sum_{i \in \mathcal{B}} \hat{\delta}_i \nabla_\theta \delta_i \\ &= -\frac{1}{N} \sum_{i \in \mathcal{B}} (\tilde{A}_i - \lambda \hat{\delta}_i) \nabla_\theta \delta_i. \end{aligned}$$

By Proposition 4.2, this equals  $\nabla_\theta \ell_{\text{RSPO}}(\theta; \mathcal{B})$ . Therefore the two objectives induce the same first-order update vector for fixed  $\lambda$  and fixed sampled data.

The rewriting in Eq. (10) follows by completing the square:

$$\begin{aligned} -\langle \tilde{A}, \delta \rangle_{\mathcal{B}} + \frac{\lambda}{2} \|\hat{\delta}\|_{\mathcal{B}}^2 &= -\langle \tilde{A}, \hat{\delta} + c_{\mathcal{B}} \mathbf{1} \rangle_{\mathcal{B}} + \frac{\lambda}{2} \|\hat{\delta}\|_{\mathcal{B}}^2 \\ &= -\langle \tilde{A}, \hat{\delta} \rangle_{\mathcal{B}} - c_{\mathcal{B}} \langle \tilde{A}, \mathbf{1} \rangle_{\mathcal{B}} + \frac{\lambda}{2} \|\hat{\delta}\|_{\mathcal{B}}^2 \\ &= \frac{\lambda}{2} \|\hat{\delta} - \frac{1}{\lambda} \tilde{A}\|_{\mathcal{B}}^2 - \frac{1}{2\lambda} \|\tilde{A}\|_{\mathcal{B}}^2 - c_{\mathcal{B}} \langle \tilde{A}, \mathbf{1} \rangle_{\mathcal{B}}. \end{aligned}$$

### E.6. Proof of Proposition D.1

Because  $Q$  is normalized and  $\delta(y) = \log(Q(y)/P(y))$ , we have

$$\mathbb{E}_{Y \sim P}[e^{\delta(Y)}] = \sum_y P(y) \frac{Q(y)}{P(y)} = \sum_y Q(y) = 1,$$

with the same argument written as an integral for continuous spaces. Since  $\delta$  is uniformly small, Taylor expansion gives

$$e^\delta = 1 + \delta + \frac{1}{2} \delta^2 + O(\|\delta\|_\infty^3).$$

Taking expectation under  $P$  and using  $\mathbb{E}_P[e^\delta] = 1$  yields

$$0 = \mathbb{E}_P[\delta] + \frac{1}{2}\mathbb{E}_P[\delta^2] + O(\|\delta\|_\infty^3).$$

Thus

$$\mathbb{E}_P[\delta] = -\frac{1}{2}\mathbb{E}_P[\delta^2] + O(\|\delta\|_\infty^3).$$

The forward KL is

$$\text{KL}(P\|Q) = \mathbb{E}_P[\log \frac{P(Y)}{Q(Y)}] = -\mathbb{E}_P[\delta(Y)] = \frac{1}{2}\mathbb{E}_P[\delta(Y)^2] + O(\|\delta\|_\infty^3).$$

Moreover  $\mathbb{E}_P[\delta] = O(\|\delta\|_\infty^2)$ , so

$$\text{Var}_{Y \sim P}(\delta(Y)) = \mathbb{E}_P[\delta(Y)^2] - \mathbb{E}_P[\delta(Y)]^2 = \mathbb{E}_P[\delta(Y)^2] + O(\|\delta\|_\infty^4).$$

Substituting this into the previous display gives

$$\text{KL}(P\|Q) = \frac{1}{2} \text{Var}_{Y \sim P}(\delta(Y)) + O(\|\delta\|_\infty^3).$$

For the reverse KL,

$$\begin{aligned} \text{KL}(Q\|P) &= \mathbb{E}_Q[\log \frac{Q(Y)}{P(Y)}] = \mathbb{E}_Q[\delta(Y)] \\ &= \mathbb{E}_P[e^{\delta(Y)}\delta(Y)] \\ &= \mathbb{E}_P[(1 + \delta(Y) + O(\|\delta\|_\infty^2))\delta(Y)] \\ &= \mathbb{E}_P[\delta(Y)] + \mathbb{E}_P[\delta(Y)^2] + O(\|\delta\|_\infty^3) \\ &= \frac{1}{2}\mathbb{E}_P[\delta(Y)^2] + O(\|\delta\|_\infty^3) \\ &= \frac{1}{2} \text{Var}_{Y \sim P}(\delta(Y)) + O(\|\delta\|_\infty^3). \end{aligned}$$

This proves both claims.

### E.7. Perturbation from surrogate error

The theory above treats  $\delta_i$  as a log-ratio surrogate. Let  $R_i$  denote an ideal uncentered log-ratio value for sample  $i$ . Suppose the corresponding ideal centered relative-score vector on the current micro-batch is  $\hat{R}_i = R_i - N^{-1} \sum_{j \in \mathcal{B}} R_j$ , and the implemented centered relative-score surrogate is

$$\hat{\delta}_i = \hat{R}_i + \hat{\xi}_i,$$

where  $\xi_i$  is the uncentered surrogate error,  $\hat{\xi}_i$  is the same error after centering, and  $\epsilon$  is a uniform error bound satisfying  $|\xi_i| \leq \epsilon$  before centering. Then  $|\hat{\xi}_i| \leq 2\epsilon$  and  $\|\hat{\xi}\|_{\mathcal{B}} \leq 2\epsilon$ .

Let  $\ell_{\text{AW}}^\delta$  and  $\ell_{\text{AW}}^R$  denote the standard group-relative advantage-weighted surrogate evaluated with the implemented surrogate  $\hat{\delta}$  and the ideal centered relative score  $\hat{R}$ , respectively. Define  $\ell_{\text{RSPO}}^\delta$  and  $\ell_{\text{RSPO}}^R$  analogously for the RSPO objective.

For the standard group-relative advantage-weighted surrogate,

$$|\ell_{\text{AW}}^\delta - \ell_{\text{AW}}^R| = |\langle \tilde{A}, \hat{\xi} \rangle_{\mathcal{B}}| \leq \|\tilde{A}\|_{\mathcal{B}} \|\hat{\xi}\|_{\mathcal{B}} \leq 2\epsilon \|\tilde{A}\|_{\mathcal{B}}.$$

For the variance term,

$$\begin{aligned} \|\hat{\delta}\|_{\mathcal{B}}^2 - \|\hat{R}\|_{\mathcal{B}}^2 &\leq 2\|\hat{R}\|_{\mathcal{B}}\|\hat{\xi}\|_{\mathcal{B}} + \|\hat{\xi}\|_{\mathcal{B}}^2 \\ &\leq 4\epsilon\|\hat{R}\|_{\mathcal{B}} + 4\epsilon^2. \end{aligned}$$

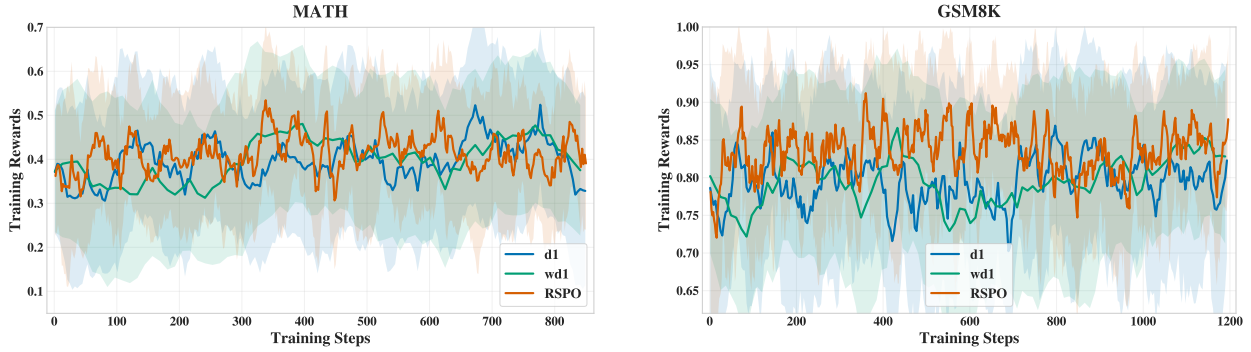


Figure 4. Additional training reward dynamics on mathematical reasoning benchmarks. Shaded regions indicate variation across runs.

Consequently, for a fixed  $\lambda$ , the forward RSPO loss values satisfy

$$|\ell_{\text{RSPO}}^{\delta} - \ell_{\text{RSPO}}^R| \leq 2\epsilon \|\tilde{A}\|_{\mathcal{B}} + \lambda(4\epsilon \|\hat{R}\|_{\mathcal{B}} + 4\epsilon^2).$$

This bound controls the scalar surrogate error. Gradient error additionally depends on the Jacobian error of the ELBO estimator and is therefore implementation-dependent.

## F. Implementation Notes

**Advantage scaling.** RSPO uses zero-sum group-relative advantages. Here  $s_q$  denotes the prompt-level advantage scale and  $\sigma_r$  denotes the reward standard deviation within the prompt group. The main experiments use the unscaled advantage  $s_q = 1$ , while the normalized ablation uses  $s_q = \sigma_r + 10^{-4}$ . When reward scaling is enabled, the sample standard deviation is computed using the training framework’s standard convention.

**Fixed feedback coefficient and inactive components.** The value  $\lambda$  used in every backward pass is fixed rather than adaptively updated; the main experiments use  $\lambda = 0.01$ , while ablations also test other fixed values. Adaptive updates of  $\lambda$ , reference-model exponential-moving-average (EMA) updates, and removal of zero-standard-deviation reward groups are not used in the reported experiments.

**Detached batch centering.** The batch mean subtracted from  $\delta$  is detached before subtraction. This makes the forward values centered but leaves  $\nabla_{\theta} \hat{\delta}_i = \nabla_{\theta} \delta_i$  in the backward pass.

**Reference-model coupling.** For each completion, the current and reference scores share the same mask time and mask set. Without this coupling, the difference  $\delta_i$  includes unnecessary independent Monte Carlo noise.

**Zero-variance reward groups.** If the group reward standard deviation is numerically zero, then all rewards in the group are equal and the group has no relative preference signal. Such groups are retained in the reported training runs; the zero-standard-deviation ratio is recorded only as a diagnostic.

## G. Additional Experimental Details

All reinforcement learning training is implemented with the TRL library (von Werra et al., 2020). We conduct experiments on mathematical reasoning and planning benchmarks, including GSM8K, Math500, Countdown, and Sudoku. All experiments were conducted on 8 NVIDIA H100-80G GPUs

**Parameter-Efficient Fine-Tuning.** For all experiments, we apply Low-Rank Adaptation (LoRA) (Hu et al., 2022) to the base diffusion language model. The LoRA rank is set to  $r = 128$  and the scaling factor is set to  $\alpha = 64$ . All reported results are obtained with LoRA fine-tuning rather than full-parameter fine-tuning.

**Optimization.** We optimize the policy with AdamW (Loshchilov & Hutter, 2017) using  $\beta_1 = 0.9$  and  $\beta_2 = 0.99$ . A constant learning-rate schedule is used throughout training. Unless otherwise specified, the learning rate is set to  $3 \times 10^{-6}$ , weight decay to 0.01, and gradient clipping to 0.2. For estimating ELBO-based relative scores, we use  $K = 2$  Monte Carlo samples for computational efficiency.

**Batching.** We set the group size according to task difficulty. For GSM8K, Countdown, and Sudoku, we use group size  $G = 6$  and total batch size 96. For MATH500, we use group size  $G = 16$  and total batch size 256. Gradient accumulation is applied to reach the target effective batch size.

**Training Steps and Checkpoint Selection.** Models are trained until the reward curve stabilizes. For mathematical reasoning, we train GSM8K and MATH500 for up to 1.5k steps. For planning tasks, we train models for up to 2k steps. During training, checkpoints are evaluated periodically, and we report the checkpoint with the highest average test accuracy across generation lengths 256 and 512.

**Decoding and Evaluation.** All experiments are conducted in the zero-shot setting. For both RL rollouts and evaluation, we use semi-autoregressive confidence-based decoding following prior dLLM work. The denoising timestep is set to half of the total generation length, and the sequence is divided into blocks of 32 tokens. At each diffusion step, the model unmask the two tokens with the highest confidence within the current incomplete block. During RL rollouts, the generation length is set to 256 for all tasks. The sampling temperature is set to 0.9 for GSM8K, MATH500, and Countdown, and to 0.3 for Sudoku. During evaluation, the temperature is set to 0.

**Special Notes on Sudoku.** Sudoku evaluation follows the verifier-based setting used in prior dLLM planning work. A generated solution is considered correct only if it satisfies the Sudoku constraints and is consistent with the given puzzle. Although some prior work uses few-shot prompting for Sudoku, our reported experiments use zero-shot prompting consistently for both training and evaluation.