On the Effect of Negative Gradient in Group Relative Deep Reinforcement Optimization

Wenlong Deng 1,2 , Yi Ren 1 , Muchen Li 1 Danica J. Sutherland 1,3 , Xiaoxiao Li 1,2† , Christos Thrampoulidis 1†

¹University of British Columbia, ²Vector Institute, ³Amii

†Corresponding author

Abstract

Reinforcement learning (RL) has become popular in enhancing the reasoning capabilities of large language models (LLMs), with Group Relative Policy Optimization (GRPO) emerging as a widely used algorithm in recent systems. Despite GRPO's widespread adoption, we identify a previously unrecognized phenomenon we term Lazy Likelihood Displacement (LLD), wherein the likelihood of correct responses marginally increases or even decreases during training. This behavior mirrors a recently discovered misalignment issue in Direct Preference Optimization (DPO), attributed to the influence of negative gradients. We provide a theoretical analysis of GRPO's learning dynamic, identifying the source of LLD as the naive penalization of all tokens in incorrect responses with the same strength. To address this, we develop a method called NTHR, which downweights penalties on tokens contributing to the LLD. Unlike prior DPO-based approaches, NTHR takes advantage of GRPO's group-based structure, using correct responses as anchors to identify influential tokens. Experiments on math reasoning benchmarks demonstrate that NTHR effectively mitigates LLD, yielding consistent performance gains across models ranging from 0.5B to 3B parameters.

1 Introduction

Reinforcement learning (RL) has become increasingly popular for improving reasoning capabilities of large language models (LLMs) [4, 9, 23]. In particular, Group Relative Policy Optimization (GRPO) [22] has emerged as a widely adopted algorithm in RL training pipelines for reasoning tasks. Models such as DeepSeek-R1 [4], DeepSeek-Math [22], Med-R1 [11], and Search-R1 [10] have successfully employed GRPO to achieve notable performance gains in code generation, mathematical problem solving, medical reasoning, and retrieval-augmented generation. These successes highlight GRPO's growing importance as a tool for aligning models to task-specific behaviors through rule-based or heuristic rewards.

Despite its empirical success, the optimization behavior of GRPO remains insufficiently studied, in part due to its recency. In this work, we focus on the impact of *negative gradients* in GRPO, which arise when the advantage function is negative—indicating that the selected action is worse than the average action at a given state (e.g., an incorrect or suboptimal response). In such cases, the learning algorithm decreases the probability of selecting this action. We draw a conceptual connection to Direct Preference Optimization (DPO) [18], which introduces negative gradients when penalizing dis-preferred responses. Recent studies [19, 17, 30, 20] have revealed that negative gradients can lead to an arguably counterintuitive behavior: a simultaneous decrease in the probabilities of preferred responses (y^+). Motivated by this insight, we examine whether similar gradient dynamics in GRPO might cause analogous issues and investigate their impact on policy learning. This effect,

termed misalignment [20] or likelihood displacement [19], is undesirable as it often degrades model performance by diverting probability mass away from optimal responses. In this work, we empirically observe (see Fig. 1) that GRPO can suffer from what we call Lazy Likelihood Displacement (LLD): a failure to sufficiently increase, or even a decrease in, the likelihood of correct answers during training. To explain this phenomenon, we provide an analysis of GRPO's update dynamics and identify that LLD stems from penalizing certain tokens in dispreferred responses. These updates inadvertently reduce the likelihood of correct answers due to shared structural or semantic features between positive and negative responses.

While several approaches have been proposed to mitigate the reduced probabilities of preferred responses (y^+) in DPO [19, 20], they either compromise data efficiency or rely on training procedures that are not directly compatible with GRPO's online learning paradigm. For instance, Razin et al. [19] recommend discarding training examples where the likelihood of positive responses y^+ decreases, thereby avoiding harmful updates—but this strategy reduces data utilization and does not address cases where the increase in the likelihood of y^+ is insufficient. Similarly, Ren and Sutherland [20] propose modifying DPO into a more online-compatible setup by training on both y^+ and y^- , which effectively mitigates the "squeezing effect" imposed by off-policy negative gradient. However, GRPO already employs online response sampling, and our empirical analysis reveals that it still suffers from the LLD effect in a subset of training examples (see Fig. 1).

To effectively address LLD in GRPO, we introduce a negative token hidden reward (NTHR) for *selective token penalization*, which reduces the penalty on tokens in incorrect responses that most strongly contribute to lowering the likelihood of correct ones. Interestingly, we find that these influential tokens often align closely with key reasoning steps in the correct responses (see Fig. 3). Our approach NTHR provides a fine-grained and targeted mitigation strategy that mitigates LLD during training. We evaluate GRPO+NTHR through per-sample training and compare it against baseline methods such as random token dropping and standard GRPO across multiple model sizes. We demonstrate that GRPO+NTHR successfully mitigates LLD and leads to improved generalization on held-out mathematical reasoning problems. In summary, our contributions are as follows:

- **GRPO** exhibits LLD for correct responses. We empirically show that GRPO suffers from Lazy Likelihood Displacement (LLD), an extended form of the misalignment phenomenon previously observed only in DPO.
- **Identifying the source of LLD.** We provide a theoretical explanation for LLD in GRPO, identifying its cause as a penalization of shared reasoning or correct tokens in dispreferred responses.
- **NTHR selective token penalization.** We introduce Negative Token Hidden Reward (NTHR) for selective *token* penalization that effectively mitigates LLD without sacrificing data efficiency.
- Empirical validation on math reasoning tasks. We demonstrate that GRPO+NTHR consistently outperforms GRPO on math benchmarks across model sizes ranging from 0.5B to 3B.

2 Preliminaries and Related Work

2.1 GRPO

GRPO loss, introduced in DeepSeek-Math [4] and DeepSeek-R1 [22], enhances fine-tuning by refining how reward and loss are calculated. Concretely, unlike traditional Proximal Policy Optimization (PPO) [21], GRPO eliminates the need for value function estimation, employing group-relative rewards for a more nuanced optimization process.

For a query-answer pair $(\boldsymbol{x}, \boldsymbol{a})$, the policy π_{θ} samples G responses $\{\boldsymbol{y}_i\}_{i=1}^G$. Each \boldsymbol{y}_i consists of a sequence of $|\boldsymbol{y}_i|$ tokens, and we denote $\boldsymbol{y}_{i, < k}$ the subsequence of the first k tokens. Let r_i denote the reward for response \boldsymbol{y}_i . The advantage of the i-th response is computed by normalizing the group-level rewards $\{r_i\}_{i=1}^G$ and is the same for each token $k=1,\ldots,|\boldsymbol{y}_i|$. Concretely, $\hat{A}_{i,k}:=\frac{r_i-\mu}{\sigma}$, with $\mu=\widehat{\mathbb{E}}[\{r_i\}_{i=1}^G]$ and $\sigma=\sqrt{\widehat{\mathrm{Var}}[\{r_i\}_{i=1}^G]}$ being the empirical average and standard deviation of

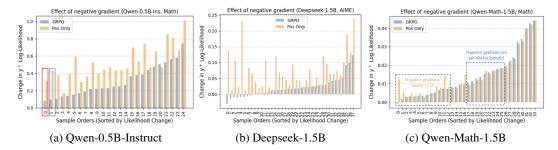


Figure 1: We show that negative gradients can lead to small or reduced likelihood change of positive samples in GRPO. The log-likelihood gains achieved by Pos Only training (orange) are significantly higher than those from GRPO (blue) for Qwen-0.5B-Ins (a) and Deepseek-1.5B (b). In Qwen-Math-1.5B (c), samples with small or reduced $\Delta(x)$ (left) are primarily influenced by negative gradients, as evidenced by their larger $\Delta(x)$ in the Pos Only setup. However, some samples on the right show smaller $\Delta(x)$ than in GRPO, indicating that negative gradients are not always harmful.

the rewards. The GRPO objective $\mathcal{J}_{\text{GRPO}}(\theta)$ is then defined as:

$$\mathbb{E}_{\substack{(\boldsymbol{x}, \boldsymbol{a}) \sim \mathcal{D} \\ \{\boldsymbol{y}_i\}_{i=1}^G \sim \pi_{\theta_{\text{old}}}(\cdot | \boldsymbol{x})}} \left[\frac{1}{\sum_{i=1}^G |\boldsymbol{y}_i|} \sum_{i=1}^G \sum_{k=1}^{|\boldsymbol{y}_i|} \min \left(\gamma_{i,k}(\boldsymbol{\theta}) \hat{A}_{i,k}, \hat{A}_{i,k} \cdot \operatorname{clip} \left(\gamma_{i,k}(\boldsymbol{\theta}), 1 - \varepsilon, 1 + \varepsilon \right) \right) \right]$$
(1)

where ε is a clipping hyperparameter, $\operatorname{clip}(\cdot)$ is the clipping operation, and $\gamma_{i,k}(\theta) = \frac{\pi_{\theta}(\boldsymbol{y}_{i,k}|\boldsymbol{x},\boldsymbol{y}_{i,< k})}{\pi_{\theta_{\text{old}}}(\boldsymbol{y}_{i,k}|\boldsymbol{x},\boldsymbol{y}_{i,< k})}$ is the likelihood ratio between the current policy π_{θ} and the old policy $\pi_{\theta_{\text{old}}}$.

2.2 Positive and Negative Gradient

The GRPO loss includes several regularization mechanisms designed to stabilize training, such as the KL term and the clipping operation. However, it is common to treat the log-likelihood component of the generated responses as responsible for "gathering knowledge," while viewing the regularization components as responsible for "stabilizing the training procedure." Without loss of generality, we focus on the online training setting for GRPO, in which $\gamma_{i,t} \approx 1$ and $\min(\cdot, \text{clip}(\cdot))$ can be safely neglected. This simplification is justified by two observations. First, GRPO inherently operates in a near-online fashion, as it generates new samples on-the-fly and typically uses mini-batches that are comparable in scale to the full batch. Second, prior work [1] has shown that omitting the clipping operation does not degrade performance. Additionally, related studies [1, 8] have demonstrated that the KL term can be omitted when other hyperparameters are carefully tuned. With this simplification, we can better understand the interactions between positive and negative responses in one roll-out. The subtle differences between GRPO and its variants, e.g., DAPO [29], Dr.GRPO [13], GPG [1], etc., could also be well interpreted.

Specifically, the GRPO's objective gradient $\nabla_{\theta} \mathcal{J}_{\text{GRPO}}(\theta)$ with respect to θ can be approximated as $\nabla_{\theta} \mathbb{E}\left[\frac{1}{\sum_{i=1}^{G}|\mathbf{y}_{i}|}\sum_{i=1}^{G}\sum_{k=1}^{|\mathbf{y}_{i}|}\hat{A}_{i,k}\gamma_{i,k}(\theta)\right]$ which using $\nabla_{\theta}\pi_{\theta}=\pi_{\theta}\nabla_{\theta}\log\pi_{\theta}$ further simplifies to

$$\mathbb{E}\Big[\frac{1}{\sum_{i=1}^{G}|\boldsymbol{y}_{i}|}\sum_{i=1}^{G}\sum_{k=1}^{|\boldsymbol{y}_{i}|}\underbrace{\hat{A}_{i,k}\gamma_{i,k}(\boldsymbol{\theta})}_{\text{constant}}\nabla_{\boldsymbol{\theta}}\log\pi_{\boldsymbol{\theta}}(\boldsymbol{y}_{i,k}\mid\boldsymbol{x},\boldsymbol{y}_{i,< k})\Big]\,.$$

This expression allows us to interpret GRPO using a similar framework as proposed by Ren and Sutherland [20], i.e., imposing a positive or a negative pressure on $y_{i,k}$. Since $\gamma_{i,k}$ is a ratio of two probabilities (and hence it must be positive), the sign of the equivalent learning rate on $y_{i,k}$ is determined by $\hat{A}_{i,k}$. Then, following a common practice that the reward is 1 for correct responses and 0 for incorrect ones, it is safe to conclude that all the tokens in the correct responses impose positive gradients while tokens in the wrong responses impose negative ones.

3 The Effect of Negative Gradients

We begin by empirically analyzing the impact of negative gradients in GRPO [4] on the likelihood of generating correct responses. Let each question \boldsymbol{x} be associated with G generated samples, of which $N^+ := N^+(\boldsymbol{x})$ give the correct answer so we call them *positive* and denote by $\boldsymbol{y}_i^+, i \in [N^+]$, and the rest $N^-L = G - N^+$ are incorrect we call them *negative* and denote by $\boldsymbol{y}_j^-, j \in [N^-]$. Positive/negative samples are given rewards $r_i = 1$ and $r_j = 0$ respectively.

We conduct experiments using math-reasoning tasks to assess how negative gradient in GRPO training influences the log-likelihood of correct responses. Specifically, we examine a range of model-dataset combinations that differ in terms of dataset difficulty and model scale: Qwen-2.5-0.5B [25], Qwen-2.5Math-1.5B [26] with the MATH dataset [7], and Deepseek-1.5B [4] with the AIME dataset. For each question, we generate 8 response rollouts, filtering out samples where all responses are either entirely correct or incorrect. We retain only those examples containing a mix of both.

Our analysis focuses on the GRPO baseline (illustrated by the blue bars in Fig. 4). To precisely characterize the learning dynamics of the algorithm, we reinitialize the model parameters θ for each individual sample, perform a single GRPO update to obtain the updated parameters θ' , and evaluate the average change in the log-likelihood of the correct responses:

$$\Delta(\boldsymbol{x}) := \frac{1}{N^{+}} \sum_{i=1}^{N^{+}} \left[\ln \pi_{\theta'}(\boldsymbol{y}_{i}^{+} \mid \boldsymbol{x}) - \ln \pi_{\theta}(\boldsymbol{y}_{i}^{+} \mid \boldsymbol{x}) \right], \tag{2}$$

for an input question x with N^+ correct responses $y_i^+, i \in [N^+]$. As illustrated in Fig. 1, some correct responses experience either a decline or only a marginal improvement in log-likelihood after training. For example, in Deepseek-1.5B (Fig. 1b), a significant portion of samples exhibit reduced likelihoods (negative values in the plot). Similarly, in Qwen-Math-1.5B and Qwen-0.5B, many samples show only minor changes in likelihood. A closer analysis shows the following.

Negative gradients lead to small or reduced likelihood change. In order to isolate the effect of negative gradients, we introduce a variant which we call Pos Only, where negative advantages are masked (i.e., values of $\hat{A}_{i,k} < 0$ are set to 0), effectively removing the influence of negative gradients during training. Fig. 1 shows the Pos Only variant in orange bars. The log-likelihood gains under Pos Only training are substantially higher than those from GRPO—particularly on the left side of the plots, where many samples show several-fold improvements.

Negative gradients are not always harmful. It is important to note that negative gradients are not inherently harmful. In terms of the likelihood change $\Delta(x)$, it is possible, as shown in Fig. 1c, that under Pos Only training, some samples exhibit even smaller $\Delta(x)$ than GRPO, which uses both positive and negative gradients. This decrease may result either from negative gradients "squeezing" the likelihood mass toward confident y^+ [20], or from certain negative gradients increasing $\Delta(x)$, as suggested by our Theorem 4.4 below. Moreover, in terms of performance, we demonstrate that Pos Only actually yields a 1.3% drop in average performance compared to GRPO. The result is shown in Tab. 2 for finetuning Qwen-Math-1.5B on the math dataset using GRPO and Pos Only respectively.

What do negative gradients of questions with small likelihood change penalize? We find that the minimal likelihood change occurs in correct responses of those questions in the dataset that have partially correct (but eventually incorrect) responses or have correct responses. In Fig. 2, we present two representative examples from Qwen-0.5B, where the incorrect responses are partially correct and the corresponding correct responses show the smallest likelihood change. The left example, highlighted in red, shows an answer that is nearly correct and corresponds to the red-boxed sample in Fig. 1a. The right example provides a correct answer but does not follow the expected output format (highlighted in red); it corresponds to the blue-boxed sample in Fig. 1a. Additional examples are provided in Appendix. This suggests that penalizing entire partially correct responses is suboptimal.

Controlling negative gradients at the token level can improve performance. While negative gradients do not necessarily hurt performance and penalizing entire incorrect responses is suboptimal, we show in Sec. 5 that carefully reducing negative gradients at *a token level* can mitigate small likelihood change and lead to improvement boosts. Our approach specifically leverages the previous observation that partial responses might be partly correct by filtering out and selectively penalizing their consistuent tokens.

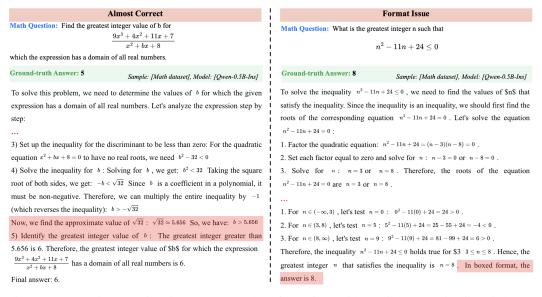


Figure 2: Inspecting negative (incorrect) samples of questions with small average likelihood change $\Delta(x)$ (Eq. (2)) reveals that they are either nearly correct (Left) or get the correct response in a wrong answer format (Right). Thus, penalizing entire negative sample responses might be suboptimal. Red dashed lines denote omitted reasoning steps.

4 Lazy Likelihood Displacement

We have shown that negative gradient in GRPO can cause a marginal or even reduced likelihood of correct responses y^+ . Here, we formalize this observation by introducing the concept of **Lazy Likelihood Displacement (LLD)**—a phenomenon in which the probability of correct answers decreases or increases only marginally after training. Our experiments reveal that LLD frequently emerges in group-based policy optimization settings, indicating that GRPO and its variants are also subject to the same limitations observed in preference optimization frameworks such as DPO [20, 19].

Definition 4.1 Let $\pi_{\theta_{init}}$ and $\pi_{\theta_{fin}}$ denote the initial and final language models, before and after optimizing a preference learning objective \mathcal{J} (e.g., Eq. (1)) over a dataset \mathcal{D} , such that $\mathcal{J}(\theta_{fin}) < \mathcal{J}(\theta_{init})$. We say that **LLD** occurs for a tuple $(\mathbf{x}, \mathbf{y}^+) \in \mathcal{D}$ if, for small nonnegative constant $\epsilon \geq 0$,

$$\ln \pi_{\theta_{fin}}(\mathbf{y}^+|\mathbf{x}) < \ln \pi_{\theta_{init}}(\mathbf{y}^+|\mathbf{x}) + \epsilon. \tag{3}$$

4.1 Understanding LLD

We begin by showing that GRPO effectively performs a weighted group preference optimization between two groups of responses: correct and incorrect ones (detailed proof in the appendix).

Lemma 4.2 When reward is binary, GRPO performs preference optimization between two distinct groups: the group of correct responses $(r_i = 1)$ and the group of incorrect responses $(r_i = 0)$. Specifically, the optimization objective reduces to the following:

$$p^{+} \sum_{i=1}^{N^{+}} \min \left(\frac{\pi_{\theta}(\boldsymbol{y}_{i}^{+}|\boldsymbol{x})}{\pi_{\theta_{old}}(\boldsymbol{y}_{i}^{+}|\boldsymbol{x})}, 1 + \varepsilon \right) - p^{-} \sum_{j=1}^{N^{-}} \max \left(\frac{\pi_{\theta}(\boldsymbol{y}_{j}^{-}|\boldsymbol{x})}{\pi_{\theta_{old}}(\boldsymbol{y}_{j}^{-}|\boldsymbol{x})}, 1 - \varepsilon \right), \tag{4}$$

where $p \stackrel{\Delta}{=} p(\boldsymbol{x}) \stackrel{\Delta}{=} \frac{1}{G} \sum_{i \in [G]} \mathbb{1}[r_i(x) = 1]$ denotes the correctness rate for a given input, $N^+ = pG$ and $N^- = (1-p)G$ are the sizes of the correct and incorrect response groups respectively, and the group-specific weights are defined as $p^+ = \frac{1-p}{\sqrt{p(1-p)}}$ and $p^- = \frac{p}{\sqrt{p(1-p)}}$.

Next, we introduce an assumption of unconstrained features, which allows us to focus our analysis on the final-layer hidden embedding $\mathbf{h_x} \in \mathbb{R}^d$ and the token unembedding matrix $\mathbf{W} \in \mathbb{R}^{|\mathcal{V}| \times d}$, where \mathcal{V} is the vocabulary of tokens.

Assumption 4.3 (Unconstrained Features) Expressive (enough) neural networks can produce unconstrained embeddings $\mathbf{h}_{\mathbf{x}} \in \mathbb{R}^d$ independent of the architecture's specific complexities [27, 16, 34, 19]. These embeddings are subsequently transformed into logits by a token unembedding matrix $\mathbf{W} \in \mathbb{R}^{|\mathcal{V}| \times d}$. The resulting logits are passed through a softmax function to yield a probability distribution over possible next tokens. To assign probabilities to sequences $\mathbf{y} \in \mathcal{V}^*$, the language model π_{θ} operates in an autoregressive manner, i.e., $\pi_{\theta}(\mathbf{y} \mid \mathbf{x}) = \prod_{k=1}^{|\mathbf{y}|} \operatorname{Softmax}(\mathbf{W}\mathbf{h}_{\mathbf{x},\mathbf{y}_{< k}})_{y_k}$.

Here, $\mathbf{h}_{x,y_{< k}}$ is the embedding of sequence $(x,y_{< k})$ and $\operatorname{Softmax}(\cdot)_{y_k}$ is the y_k -th entry of the V-dimensional softmax map. We focus on the online training setting for GRPO, as state in Sec. 2.2, the effect of the clipping operation is effectively removed. We assume that the objective is optimized via gradient flow and theoretically analyze the dynamics of the log-likelihood of a positive response, $\frac{d}{dt} \ln \pi_{\theta(t)}(\mathbf{y}_i^+ \mid \mathbf{x})$, yielding the following characterization (see the appendix for a proof).

Theorem 4.4 For any question \mathbf{x} , at any time $t \geq 0$ of training, and any correct response $\mathbf{y}_i^+, i \in [N^+]$, in addition to the dependence on token unembeddings, the likelihood change $\frac{d}{dt} \ln \pi_{\theta(t)}(\mathbf{y}_i^+|\mathbf{x})$ exhibits increased laziness (that is, has smaller magnitude) as the following quantity increases:

$$\underbrace{p^{-\sum_{k=1}^{|\boldsymbol{y}_{i}^{+}|}\sum_{j=1}^{N^{-}}\sum_{k'=1}^{|\boldsymbol{y}_{j}^{-}|}\alpha_{k,k'}^{-}\cdot\langle\mathbf{h}_{\mathbf{x},\boldsymbol{y}_{i,< k}^{+}},\mathbf{h}_{\mathbf{x},\boldsymbol{y}_{j,< k'}^{-}}\rangle}_{-p^{+\sum_{k=1}^{|\boldsymbol{y}_{i}^{+}|}\sum_{i'=1}^{N^{+}}\sum_{k''=1}^{|\boldsymbol{y}_{i'}^{+}|}\alpha_{k,k''}^{+}\cdot\langle\mathbf{h}_{\mathbf{x},\boldsymbol{y}_{i,< k}^{+}},\mathbf{h}_{\mathbf{x},\boldsymbol{y}_{i',< k''}^{+}}\rangle. (5)$$

Here, $\alpha_{k,k'}^-$ and $\alpha_{k,k''}^+$ are token-level prediction error similarity weights, which quantify the similarity of token-level prediction error across responses (see Appendix for formal definitions.). We refer to the quantity in Eq. (5) as Group Weighted Hidden Embedding Score (GWHES).

The first term in Eq. (5) captures the influence of negative gradients on the likelihood of the correct response. Specifically, this shows that a negative token with a large value of $\alpha_{k,k'}^- \cdot \langle \mathbf{h}_{\mathbf{x}, y_{i,< k}}^+, \mathbf{h}_{\mathbf{x}, y_{j,< k'}}^- \rangle$ is more likely to cause the LLD of the correct responses. This motivates the following corollary.

Corollary 4.5 For any question \mathbf{x} , the negative gradient associated with the hidden embedding of a token k' in any incorrect response \mathbf{y}_j^- , $j \in [N^-]$, will exert a greater adverse effect on the likelihood change $\frac{d}{dt} \ln \pi \theta(t)(\mathbf{y}_i^+|\mathbf{x})$ as the following quantity getslarger:

$$\sum_{k=1}^{|\boldsymbol{y}_{i}^{+}|} \alpha_{k,k'}^{-} \cdot \langle \mathbf{h}_{\mathbf{x},\boldsymbol{y}_{i,< k}^{+}}, \mathbf{h}_{\mathbf{x},\boldsymbol{y}_{j,< k'}^{-}} \rangle.$$
 (6)

This can occur when the negative and positive tokens have high embedding similarity, theoretically supporting the observation in Fig. 2 that LLD samples are often nearly correct.

4.2 Identifying LLD Samples

Theorem 4.4 suggests using $\Delta_{\rm GWHES}$ as a metric to identify samples that cause LLD. Here, we verify the validity of this metric experimentally. For Qwen-1.5B-deepseek, we generate responses for the first 100 questions from the AIME dataset (1983–2023) and compute the GWHES scores for those questions with non-extreme predicted probabilities (0 < p < 1), resulting in 57 valid questions. For Qwen-2.5Math-1.5B, we use questions from the Math dataset and retain 47 valid questions. We then

| Top-K | Qwen-1.5 | B-deepseek | Qwen-1. | 5B-math |
|-------|----------|------------|---------|---------|
| | GWHES | Random | GWHES | Random |
| 10 | 50% | 17.5% | 60% | 21.3% |
| 15 | 75% | 26.3% | 75% | 31.9% |

Table 1: Ranking questions by GWHES results in a significantly higher $\operatorname{Top-}K$ overlap with the likelihood-change ranking compare to ranking randomly. This demonstrates the effectiveness of GWHES in identifying LLD samples.

calculate the Top-K overlap accuracy between two rankings of questions; one based on the ascending

¹Theorem 4.4 can be seen as an extension of [19, Thm.3] for DPO to the more challenging GRPO setting with group-based structure. Unlike DPO, we will exploit this group-structure to selectively penalize tokens rather than responses.

order of likelihood changes $\Delta(x)$ as defined in Eq. (2) (denoted \mathcal{Q}_K^L), and the other based on the descending order of GWHES scores as defined in Eq. (5) (denoted \mathcal{Q}_K^S). The overlap accuracy is computed as $\frac{|\mathcal{Q}_K^L \cap \mathcal{Q}_K^S|}{K}$. The results in Tab. 1 indicate that the overlap achieved using GWHES is significantly greater than that of the random ranking baseline, highlighting a strong alignment between high GWHES scores and LLD samples.

5 Selective Token Penalization in Incorrect Responses

Theorem 4.4 identifies samples prone to causing LLD and Corollary 4.5 quantifies the influence of individual tokens' negative gradients on the likelihood of correct responses. The goal of this section is to explore whether appropriately controlling negative gradients can mitigate LLD and further lead to performance improvements over vanilla GRPO. One approach could be using $\Delta_{\rm GWHES}$ to identify negative samples \boldsymbol{y}_i^- that cause LLD and entirely discard them during training. However, this compromises data efficiency and fails to address cases where samples only marginally improve \boldsymbol{y}^+ . Instead, we propose a more targeted approach: identify tokens that negatively impact $\Delta(\boldsymbol{x})$ and reduce the penalty applied to them.

Math Question: How many elements are in the intersection of the set of all the prime numbers less than 30 and the set of all the odd numbers greater than zero?

Ground-truth Answer: 9

Sample from [Math dataset]

Correct Response

To determine ... Step 1: Identify the set of all prime numbers less than 30. The prime numbers less than 30 are: 2, 3, 5, 7, 11, 13, 17, 19, 23, 29. So, the set of all prime numbers less than 30 is {2, 3, 5, 7, 11, 13, 17, 19, 23, 29}.

Step 2: Identify the set of all odd numbers greater than zero. The set of all odd numbers greater than zero is: $\{1, 3, 5, 7, 9, 11, 13, 15, 17, 19, 21, 23, 25, 27, 29, ...\}$. So, the set of all odd numbers greater than zero is $\{1, 3, 5, 7, 9, 11, 13, 15, 17, 19, 21, 23, 25, 27, 29, ...\}$. Step 3: Find the intersection of the two sets ... From the

two sets, we see that the common elements are 3, 5, 7, 11, 13, 17, 19, 23, and 29. So, the intersection is {3, 5, 7, 11, 13, 17, 19, 23, 29}.

Step 4: Count the number of elements in the intersection. The number of elements in the intersection is 9. Therefore, the number of elements in the intersection of the set of all the prime numbers less than 30 and the set of all the odd numbers greater than zero is \boxed \{9\}.

Incorrect Response Contain Correct Steps

To determine ... First, let 's list the prime numbers less than $3 \ 0 : 2, 3, 5, 7, 11, 13, 17, 19, 23, 29$ So, the set of all prime numbers less than $3 \ 0$ is : $\{2, 3, 5, 7, 11, 13, 17, 19, 23, 29\}$

Next, let 's list the set of all odd numbers greater than zero is: $\{1, 3, 5, 7, 9, 11, 13, 15, 17, 19, 21, 23, 25, 27\}$ Now, we find the intersection of these two sets. The intersection consists of numbers that are both in the set of prime numbers less than 30 and in the set of all odd numbers greater than zero. From the lists above, we can see that the only numbers that are both prime and odd are the odd prime numbers.

Therefore, the intersection is: $\{3,5,7,11,13,17,19,23\}$ To find the number of elements in the intersection, we count the elements in the list: $\{3,5,7,11,13,17,19,23\}$ There are 8 elements in the intersection. Thus, the number of elements in the intersection is boxed $\{8\}$.

Figure 3: Key insight: Tokens of negative samples (incorrect responses) can be logically or step-correct. Tokens with high NTHR tend to strongly correlate with these types of tokens (highlighted in red). yThe bold dots represent omitted reasoning.

5.1 Negative Token Hidden Reward (NTHR)

Motivated by Corollary 4.5, we define the impact that an individual token k' of incorrect response $y_i^-, j \in [N^-]$ has on the likelihood of a group of correct responses as follows:

$$s_{j,< k'}^{-} := \sum_{i=1}^{N^{+}} \sum_{k=1}^{|\mathbf{y}_{i}^{+}|} \alpha_{k,k'}^{-} \cdot \left\langle \mathbf{h}_{\mathbf{x},\mathbf{y}_{i,< k}^{+}}, \mathbf{h}_{\mathbf{x},\mathbf{y}_{j,< k'}^{-}} \right\rangle, \tag{7}$$

which defines a relative influence from imposing a negative gradient on y_j^- to all tokens in all positive responses. Specifically, a token from an incorrect response can negatively affect the likelihood of a group of a group of correct responses when $s_{j,< k'}^- \ge 0$, with the magnitude of $s_{j,< k'}^-$ indicating the severity of this effect. We refer to this value as the **negative token hidden reward** (NTHR):

it quantifies a token's detrimental impact on correct responses' likelihood. In Fig. 3, we visualize tokens with high $s_{j, < k'}$ values—highlighted in red. We find that many of these influential tokens are logically or stepwise correct terms (such as "odd," "prime," and "intersection" in the specific example). These terms exhibit strong semantic alignment with the correct responses.

5.2 NTHR selective token penalization.

As shown in Corollary 4.5, tokens with high positive NTHR values tend to cause significant reductions in the likelihood of generated correct responses. To mitigate this adverse effect, we introduce a selective penalization strategy (detailed in Algorithm 1 and complexity discussion in Appendix) that attenuates the penalty on tokens from negative responses whose NTHR scores exceed a threshold τ . Concretely, we define the resulting set of selected tokens as:

$$\mathbf{V}_{j}^{-} = \left\{ \mathbf{y}_{j,k'}^{-} \middle| s_{j,< k'}^{-} > \tau \right\}. \tag{8}$$

To determine the threshold τ , an effective pratical strategy is computing the minimum average token-level influence exerted by each correct response on all other correct responses. Formally, we set $\tau = \beta \cdot \min_{i' \in [N^+]} \bar{s}_{i'}^+$

Algorithm 1 NTHR selective token penalization

Responses' hidden embeddings 1: **Input:** $\{\mathbf{h}_i^+\}_{i=1}^{N^+}, \{\mathbf{h}_j^-\}_{j=1}^{N^-}, \text{ coefficient } \alpha, \text{ scale factor }$

2: **for** i' = 1 to N^+ **do**

 $\bar{s}_{i'}^+ \leftarrow \text{Eq.}(7)$ 3:

4: end for

5: $\tau \leftarrow \beta \cdot \min_{i' \in [N^+]} \bar{s}_{i'}^+$

6: Select negative tokens:

7: **for** each negative response j **do**

 $\mathbf{V}_{j}^{-} \leftarrow \left\{ \mathbf{y}_{j,k'}^{-} \mid s_{j,< k'}^{-} > \tau \right\}$

10: Apply selective penalization:

11: **for** each retained token $\mathbf{y}_{j,k'}^- \in \mathbf{V}_j^-$ **do**

12: $\hat{A}^-_{j,k',\eta} \leftarrow \eta \cdot \hat{A}^-_{j,k'}$ 13: **end for**

where β is a scale factor and $\bar{s}_{i'}^+$ measures the average impact of the i'-th correct response's tokens on the likelihoods of other correct responses:

$$\bar{s}_{i'}^{+} := \frac{1}{|\mathbf{y}_{i'}^{+}|} \sum_{k''=1}^{|\mathbf{y}_{i'}^{+}|} \sum_{i=1}^{N^{+}} \sum_{k=1}^{|\mathbf{y}_{i}^{+}|} \alpha_{k,k''}^{+}(t) \cdot \left\langle \mathbf{h}_{\mathbf{x},\mathbf{y}_{i,< k}^{+}}(t), \mathbf{h}_{\mathbf{x},\mathbf{y}_{i',< k''}^{+}}(t) \right\rangle. \tag{9}$$

A larger $\bar{s}_{i'}^+$ indicates a stronger mutual influence among positive responses, which could be interpreted as an estimation of the "local elasticity" of a deep neural network [6]. Finally, to apply selective penalization, we define the advantage of each retained token using a scale factor $\eta < 1$, resulting in $\hat{A}_{i,k',\eta}^- := \eta \cdot \hat{A}_{i,k'}^-$. The scale factor will reduce the penalty on selected negative tokens.

NTHR Selective Token Penalization Mitigates LLD

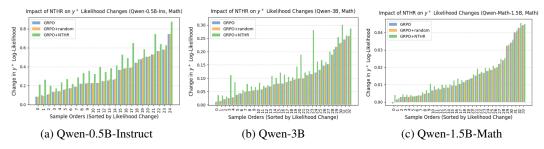


Figure 4: GRPO+NTHR consistently improves likelihood change of correct responses, as indicated by the green bars exceeding the blue bars. While GRPO+Random offers only modest improvements, GRPO+NTHR consistently outperforms it, highlighting the effectiveness of NTHR in identifying LLD tokens.

To evaluate the effectiveness of our NTHR selective token penalization in mitigating LLD, we conduct experiments using the same setup described in Sec. 4. Specifically, we identify the selected token set \mathbf{V}_j^- in incorrect responses using a scale parameter $\beta=1$, and apply no penalization to these tokens by setting $\eta=0$. We refer to this variant as GRPO+NTHR (green bars). We then compare the following two baselines: (1) Standard GRPO (blue bars), and (2) GRPO+Random (orange bars), where we randomly select $|\mathbf{V}_j^-|$ tokens and set their advantage to 0 (thus avoid penalizing them in optimization).

The results in Fig. 4 show that GRPO+NTHR consistently improves likelihood change across all samples, as evidenced by the green bars surpassing the blue bars. Notably, while removing negative gradients entirely—as in the Pos Only in Fig. 1c—can sometimes reduce the likelihood change compared to GRPO in Qwen-Math-1.5B, our GRPO+NTHR approach (Fig. 4c) maintains consistent improvements. Furthermore, unlike GRPO+Random that provides only modest gains across all models, GRPO+NTHR delivers consistent and substantial improvements, highlighting the effectiveness of our proposed NTHR in identifying LLD tokens.

5.4 Performance of NTHR Selective Token Penalization

Having demonstrated the impact of negative gradients on likelihood change and the effectiveness of our method in mitigating the LLD issue in Sec. 5.3, we now study the effect of negative gradient on model performance. For this, we finetune models with different sizes using MATH dataset (levels 3–5) [7] and use greedy decoding to evaluate finetuned models on five math benchmarks: AIME24, AMC, MATH500, Minerva, and Olympiad (detailed implementation details in Appendix).

We find that NTHR consistently boosts model performance across various sizes, with detailed results presented in Tab. 2. Notably, even for Qwen2.5-Math which shows a substantial improvement of over 20% after RL fine-tuning—outpacing gains seen in other models—NTHR still provides an average improvement of 0.8%. For Qwen2.5-0.5B-Ins and Qwen2.5-1.5B-Ins, NTHR consistently outperforms GRPO by 1.1% and 1.5% on average respectively. Similar improvements are observed with a larger model Qwen2.5-3B, where NTHR reaches 36.30 % compared to GRPO's 33.88 % and the base model's 31.36 %. We further fine-tune Qwen2.5-Math on the more challenging DeepScaler [14] dataset and observe that NTHR outperforms GRPO by 1.8%.

| Base model + Method | AIME24 | AMC | MATH500 | Minerva | Olympiad | Avg. |
|--------------------------------|--------|------|---------|---------|----------|-------|
| Qwen2.5-Math-1.5B | | | | | | |
| Base | 3.3 | 20.0 | 39.6 | 7.7 | 24.9 | 19.10 |
| GRPO | 13.3 | 57.5 | 71.8 | 29.0 | 34.1 | 41.14 |
| Pos Only | 10.0 | 57.5 | 70.6 | 30.1 | 31.0 | 39.84 |
| NTHR | 16.7 | 57.5 | 70.8 | 30.5 | 34.2 | 41.94 |
| Qwen2.5-0.5B-Ins | | | | | | |
| Base | 0.0 | 2.5 | 33.4 | 4.4 | 7.0 | 9.46 |
| GRPO | 0.0 | 7.5 | 33.8 | 9.2 | 8.1 | 11.72 |
| NTHR | 0.0 | 10.0 | 36.6 | 8.1 | 8.6 | 12.66 |
| Qwen2.5-1.5B-Ins | | | | | | |
| Base | 0.0 | 22.5 | 53.0 | 19.1 | 20.7 | 23.06 |
| GRPO | 3.3 | 32.5 | 57.2 | 18.8 | 23.0 | 26.96 |
| NTHR | 6.7 | 35.0 | 58.8 | 21.0 | 20.9 | 28.48 |
| Qwen2.5-Math-1.5B (deepscaler) | | | | | | |
| Base | 3.3 | 20.0 | 39.6 | 7.7 | 24.9 | 19.10 |
| GRPO | 10.0 | 42.5 | 72.4 | 32.4 | 31.9 | 37.80 |
| NTHR | 16.7 | 47.5 | 73.2 | 29.4 | 31.4 | 39.60 |
| Qwen2.5-3B | | | | | | |
| Base | 10.0 | 37.5 | 58.6 | 26.1 | 24.6 | 31.36 |
| GRPO | 6.7 | 35.0 | 66.6 | 31.2 | 29.9 | 33.88 |
| NTHR | 10.0 | 47.5 | 65.6 | 31.6 | 26.8 | 36.30 |

Table 2: Results across selected math benchmarks for different Qwen2.5 models and methods. NTHR consistently provides average performance gains on various models.

6 Conclusion

In this work, we study how negative gradients affect the likelihood of correct responses in GRPO. We identify a phenomenon we term Lazy Likelihood Displacement (LLD), where penalization of

incorrect responses inadvertently reduces or lead to small likelihood change of correct ones. To address this, we propose the negative token hidden reward (NTHR) along with a selective token-level penalization strategy. NTHR reduces the penalty on tokens in incorrect responses that contribute most to lowering the likelihood of correct responses, thus successfully mitigating the LLD issue. Through experiments on math reasoning tasks across models ranging from 0.5B to 3B, we demonstrate that NTHR improves GRPO, validating the effectiveness of our approach in addressing LLD and enhancing performance.

Acknowledgments: This work was partially funded by the NSERC Discovery Grant RGPIN-2021-03677, Alliance Grant ALLRP 581098-22, the Natural Science and Engineering Research Council of Canada (NSERC), the Canada CIFAR AI Chairs program, the Canada Research Chair program, an IITP grant funded by MSIT, the BC DRI Group, and the Digital Research Alliance of Canada.

References

- [1] Xiangxiang Chu, Hailang Huang, Xiao Zhang, Fei Wei, and Yong Wang. Gpg: A simple and strong reinforcement learning baseline for model reasoning. *arXiv* preprint arXiv:2504.02546, 2025.
- [2] Wenlong Deng, Yize Zhao, Vala Vakilian, Minghui Chen, Xiaoxiao Li, and Christos Thrampoulidis. Dare the extreme: Revisiting delta-parameter pruning for fine-tuned models. *arXiv* preprint arXiv:2410.09344, 2024.
- [3] Kanishk Gandhi, Ayush Chakravarthy, Anikait Singh, Nathan Lile, and Noah D Goodman. Cognitive behaviors that enable self-improving reasoners, or, four habits of highly effective stars. *arXiv* preprint arXiv:2503.01307, 2025.
- [4] Daya Guo, Dejian Yang, Haowei Zhang, Junxiao Song, Ruoyu Zhang, Runxin Xu, Qihao Zhu, Shirong Ma, Peiyi Wang, Xiao Bi, et al. Deepseek-r1: Incentivizing reasoning capability in llms via reinforcement learning. *arXiv preprint arXiv:2501.12948*, 2025.
- [5] Chaoqun He, Renjie Luo, Yuzhuo Bai, Shengding Hu, Zhen Leng Thai, Junhao Shen, Jinyi Hu, Xu Han, Yujie Huang, Yuxiang Zhang, et al. Olympiadbench: A challenging benchmark for promoting agi with olympiad-level bilingual multimodal scientific problems. *arXiv preprint arXiv:2402.14008*, 2024.
- [6] Hangfeng He and Weijie Su. The local elasticity of neural networks. In *International Conference on Learning Representations*, 2020.
- [7] Dan Hendrycks, Collin Burns, Saurav Kadavath, Akul Arora, Steven Basart, Eric Tang, Dawn Song, and Jacob Steinhardt. Measuring mathematical problem solving with the math dataset. *NeurIPS*, 2021.
- [8] Jingcheng Hu, Yinmin Zhang, Qi Han, Daxin Jiang, Xiangyu Zhang, and Heung-Yeung Shum. Open-reasoner-zero: An open source approach to scaling up reinforcement learning on the base model. *arXiv preprint arXiv:2503.24290*, 2025.
- [9] Aaron Jaech, Adam Kalai, Adam Lerer, Adam Richardson, Ahmed El-Kishky, Aiden Low, Alec Helyar, Aleksander Madry, Alex Beutel, Alex Carney, et al. Openai o1 system card. *arXiv* preprint arXiv:2412.16720, 2024.
- [10] Bowen Jin, Hansi Zeng, Zhenrui Yue, Jinsung Yoon, Sercan Arik, Dong Wang, Hamed Zamani, and Jiawei Han. Search-r1: Training llms to reason and leverage search engines with reinforcement learning. *arXiv preprint arXiv:2503.09516*, 2025.
- [11] Yuxiang Lai, Jike Zhong, Ming Li, Shitian Zhao, and Xiaofeng Yang. Med-r1: Reinforcement learning for generalizable medical reasoning in vision-language models. *arXiv* preprint arXiv:2503.13939, 2025.
- [12] Aitor Lewkowycz, Anders Andreassen, David Dohan, Ethan Dyer, Henryk Michalewski, Vinay Ramasesh, Ambrose Slone, Cem Anil, Imanol Schlag, Theo Gutman-Solo, et al. Solving quantitative reasoning problems with language models. *Advances in Neural Information Processing Systems*, 35:3843–3857, 2022.

- [13] Zichen Liu, Changyu Chen, Wenjun Li, Penghui Qi, Tianyu Pang, Chao Du, Wee Sun Lee, and Min Lin. Understanding r1-zero-like training: A critical perspective. arXiv preprint arXiv:2503.20783, 2025.
- [14] Michael Luo, Sijun Tan, Justin Wong, Xiaoxiang Shi, William Tang, Manan Roongta, Colin Cai, Jeffrey Luo, Tianjun Zhang, Erran Li, Raluca Ada Popa, and Ion Stoica. Deepscaler: Surpassing o1-preview with a 1.5b model by scaling rl, 2025. Notion Blog.
- [15] Yuchun Miao, Sen Zhang, Liang Ding, Yuqi Zhang, Lefei Zhang, and Dacheng Tao. The energy loss phenomenon in rlhf: A new perspective on mitigating reward hacking. *arXiv* preprint *arXiv*:2501.19358, 2025.
- [16] Dustin G Mixon, Hans Parshall, and Jianzong Pi. Neural collapse with unconstrained features. *Sampling Theory, Signal Processing, and Data Analysis*, 20(2):11, 2022.
- [17] Arka Pal, Deep Karkhanis, Samuel Dooley, Manley Roberts, Siddartha Naidu, and Colin White. Smaug: Fixing failure modes of preference optimisation with dpo-positive. arXiv preprint arXiv:2402.13228, 2024.
- [18] Rafael Rafailov, Archit Sharma, Eric Mitchell, Christopher D Manning, Stefano Ermon, and Chelsea Finn. Direct preference optimization: Your language model is secretly a reward model. *Advances in Neural Information Processing Systems*, 36:53728–53741, 2023.
- [19] Noam Razin, Sadhika Malladi, Adithya Bhaskar, Danqi Chen, Sanjeev Arora, and Boris Hanin. Unintentional unalignment: Likelihood displacement in direct preference optimization. *arXiv* preprint arXiv:2410.08847, 2024.
- [20] Yi Ren and Danica J Sutherland. Learning dynamics of llm finetuning. *arXiv preprint* arXiv:2407.10490, 2024.
- [21] John Schulman, Filip Wolski, Prafulla Dhariwal, Alec Radford, and Oleg Klimov. Proximal policy optimization algorithms. *arXiv* preprint arXiv:1707.06347, 2017.
- [22] Zhihong Shao, Peiyi Wang, Qihao Zhu, Runxin Xu, Junxiao Song, Xiao Bi, Haowei Zhang, Mingchuan Zhang, YK Li, Y Wu, et al. Deepseekmath: Pushing the limits of mathematical reasoning in open language models. *arXiv preprint arXiv:2402.03300*, 2024.
- [23] Gemini Team, Rohan Anil, Sebastian Borgeaud, Jean-Baptiste Alayrac, Jiahui Yu, Radu Soricut, Johan Schalkwyk, Andrew M Dai, Anja Hauth, Katie Millican, et al. Gemini: a family of highly capable multimodal models. *arXiv preprint arXiv:2312.11805*, 2023.
- [24] Hemish Veeraboina. Aime problem set 1983-2024, 2023. URL https://www.kaggle.com/datasets/hemishveeraboina/aime-problem-set-1983-2024.
- [25] An Yang, Baosong Yang, Beichen Zhang, Binyuan Hui, Bo Zheng, Bowen Yu, Chengyuan Li, Dayiheng Liu, Fei Huang, Haoran Wei, et al. Qwen2. 5 technical report. *arXiv preprint arXiv:2412.15115*, 2024.
- [26] An Yang, Beichen Zhang, Binyuan Hui, Bofei Gao, Bowen Yu, Chengpeng Li, Dayiheng Liu, Jianhong Tu, Jingren Zhou, Junyang Lin, Keming Lu, Mingfeng Xue, Runji Lin, Tianyu Liu, Xingzhang Ren, and Zhenru Zhang. Qwen2.5-math technical report: Toward mathematical expert model via self-improvement. *arXiv preprint arXiv:2409.12122*, 2024.
- [27] Zhilin Yang, Zihang Dai, Ruslan Salakhutdinov, and William W Cohen. Breaking the softmax bottleneck: A high-rank rnn language model. *arXiv preprint arXiv:1711.03953*, 2017.
- [28] Longhui Yu, Weisen Jiang, Han Shi, Jincheng Yu, Zhengying Liu, Yu Zhang, James T Kwok, Zhenguo Li, Adrian Weller, and Weiyang Liu. Metamath: Bootstrap your own mathematical questions for large language models. *arXiv preprint arXiv:2309.12284*, 2023.
- [29] Qiying Yu, Zheng Zhang, Ruofei Zhu, Yufeng Yuan, Xiaochen Zuo, Yu Yue, Tiantian Fan, Gaohong Liu, Lingjun Liu, Xin Liu, et al. Dapo: An open-source llm reinforcement learning system at scale. *arXiv preprint arXiv:2503.14476*, 2025.

- [30] Lifan Yuan, Ganqu Cui, Hanbin Wang, Ning Ding, Xingyao Wang, Jia Deng, Boji Shan, Huimin Chen, Ruobing Xie, Yankai Lin, et al. Advancing llm reasoning generalists with preference trees. *arXiv preprint arXiv:2404.02078*, 2024.
- [31] Yu Yue, Yufeng Yuan, Qiying Yu, Xiaochen Zuo, Ruofei Zhu, Wenyuan Xu, Jiaze Chen, Chengyi Wang, TianTian Fan, Zhengyin Du, et al. Vapo: Efficient and reliable reinforcement learning for advanced reasoning tasks. *arXiv preprint arXiv:2504.05118*, 2025.
- [32] Weihao Zeng, Yuzhen Huang, Qian Liu, Wei Liu, Keqing He, Zejun Ma, and Junxian He. Simplerl-zoo: Investigating and taming zero reinforcement learning for open base models in the wild. *arXiv preprint arXiv:2503.18892*, 2025.
- [33] Xiaoying Zhang, Jean-Francois Ton, Wei Shen, Hongning Wang, and Yang Liu. Overcoming reward overoptimization via adversarial policy optimization with lightweight uncertainty estimation. *arXiv preprint arXiv:2403.05171*, 2024.
- [34] Yize Zhao, Tina Behnia, Vala Vakilian, and Christos Thrampoulidis. Implicit geometry of next-token prediction: From language sparsity patterns to model representations. *arXiv* preprint *arXiv*:2408.15417, 2024.

7 Appendix

Notation. For any time $t \geq 0$, we use $\boldsymbol{W}(t)$, $\boldsymbol{w}_z(t)$, and $\boldsymbol{h}_{\boldsymbol{z}}(t)$ to denote the token unembedding matrix, unembedding of a token $\boldsymbol{z} \in \mathcal{V}$, and hidden embedding of $z \in \mathcal{V}^*$, respectively. We let \boldsymbol{z}_k be the k-th token in \boldsymbol{z} and $\boldsymbol{z}_{< k}$ be the first k-1 tokens in \boldsymbol{z} . For a question \boldsymbol{x} , the old policy $\pi_{\theta_{\text{old}}}$ generates a group of G samples and resulting $(\boldsymbol{x}, \{\boldsymbol{y}_i^+\}_{N^+}, \{\boldsymbol{y}_j^-\}_{N^-})$, where $N^+ + N^- = G$. Lastly, we denote by $\boldsymbol{e}_z \in \mathbb{R}^{|\mathcal{V}|}$ the standard basis vector corresponding to $z \in \mathcal{V}$.

7.1 Proof of Lemma 4.2: GRPO as Group Preference Optimization

In this section, we demonstrate that training with GRPO constitutes preference optimization. We adopt a binary reward system, assigning r=1 to correct responses and r=0 to incorrect ones, consistent with recent works [13, 29]. It is notable that this also applies to its variants, e.g., DAPO [29], Dr.GRPO [13], GPG [1], etc.

For a single question x, we simplify the expected loss $\mathbb{E}_{\{y_i\}_{i=1}^G \sim \pi_{\theta_{\text{old}}}(\cdot|x)}[\mathcal{J}_x]$ in Eq. (1) by omitting the token normalization term $\frac{1}{\sum_{i=1}^G |y_i|}$, yielding:

$$\mathbb{E}_{\{\boldsymbol{y}_i\}_{i=1}^G \sim \pi_{\theta_{\text{old}}}(\cdot|\boldsymbol{x})} \left[\sum_{i=1}^G \sum_{k=1}^{|\boldsymbol{y}_i|} \min \left(\gamma_{i,k}(\theta) \hat{A}_{i,k}, \hat{A}_{i,k} \cdot \operatorname{clip} \left(\gamma_{i,k}(\theta), 1 - \varepsilon, 1 + \varepsilon \right) \right) \right]. \tag{10}$$

The success probability for question x is $p = \hat{\mathbb{P}}_i(r_{i,x} = 1) \approx \mathbb{P}(r_x = 1)$ where we use a population approximation for large enough G. In this case, we take $\mu = p$ and $\sigma = \sqrt{p(1-p)}$. Consequently, the advantage $\hat{A}_{i,k}$ becomes:

$$\hat{A}_{i,k} = \begin{cases} \frac{1-p}{\sqrt{p(1-p)}} & \text{if } r_i = 1, \\ -\frac{p}{\sqrt{p(1-p)}} & \text{if } r_i = 0. \end{cases}$$
 (11)

Since

$$\min \left(\gamma_{i,k}(\theta) \hat{A}_{i,k}, \ \hat{A}_{i,k} \cdot \operatorname{clip} \left(\gamma_{i,k}(\theta), \ 1 - \varepsilon, \ 1 + \varepsilon \right) \right)$$

is equivalent tos

$$\begin{cases} \hat{A}_{i,k} \cdot \min \left(\gamma_{i,k}(\theta), \ 1 + \varepsilon \right), & \text{if } \hat{A}_{i,k} > 0 \Leftrightarrow r_i = 1, \\ \hat{A}_{i,k} \cdot \max \left(\gamma_{i,k}(\theta), \ 1 - \varepsilon \right), & \text{if } \hat{A}_{i,k} < 0 \Leftrightarrow r_i = 0, \end{cases}$$

the expected loss then becomes:

$$\mathbb{E}_{\{\boldsymbol{y}_{i}\}_{i=1}^{G} \sim \pi_{\theta_{\text{old}}}(\cdot|\boldsymbol{x})}[\mathcal{J}_{\boldsymbol{x}}] = \frac{1-p}{\sqrt{p(1-p)}} \mathbb{E}_{\{\boldsymbol{y}_{i}\}_{i=1}^{G} \sim \pi_{\theta_{\text{old}}}(\cdot|\boldsymbol{x})} \left[\min\left(\frac{\pi_{\theta}(\boldsymbol{y}_{i}|\boldsymbol{x})}{\pi_{\theta_{\text{old}}}(\boldsymbol{y}_{i}|\boldsymbol{x})}, 1+\varepsilon\right) \mathbf{1}_{r_{i}=1} \right] - \frac{p}{\sqrt{p(1-p)}} \mathbb{E}_{\{\boldsymbol{y}_{i}\}_{i=1}^{G} \sim \pi_{\theta_{\text{old}}}(\cdot|\boldsymbol{x})} \left[\max\left(\frac{\pi_{\theta}(\boldsymbol{y}_{i}|\boldsymbol{x})}{\pi_{\theta_{\text{old}}}(\boldsymbol{y}_{i}|\boldsymbol{x})}, 1-\varepsilon\right) \mathbf{1}_{r_{i}=0} \right].$$

$$(12)$$

When generating G sampled responses, the loss becomes:

$$p^{+} \sum_{i=1}^{N^{+}} \min \left(\frac{\pi_{\theta}(\boldsymbol{y}_{i}^{+}|\boldsymbol{x})}{\pi_{\theta_{\text{old}}}(\boldsymbol{y}_{i}^{+}|\boldsymbol{x})}, 1 + \varepsilon \right) - p^{-} \sum_{j=1}^{N^{-}} \max \left(\frac{\pi_{\theta}(\boldsymbol{y}_{j}^{-}|\boldsymbol{x})}{\pi_{\theta_{\text{old}}}(\boldsymbol{y}_{j}^{-}|\boldsymbol{x})}, 1 - \varepsilon \right), \tag{13}$$

where $N^+=pG$, $N^-=(1-p)G$, $p^+=\frac{1-p}{\sqrt{p(1-p)}}$ and $p^-=\frac{p}{\sqrt{p(1-p)}}$. This represents group preference optimization, i.e., increasing the likelihood of correct responses while penalizing incorrect ones.

7.2 Proof of Theorem 4.4

Assume that all responses are sequences whose first tokens are distinct from each other. Then we analyze the likelihood change for a correct response y_i^+ by measuring $\frac{d}{dt} \ln \pi_{\theta(t)}(y_i^+|\mathbf{x})$. According

to the chain rule, using the loss in Eq. (13) and denoting $\pi_{< k} = \pi_{\theta}(\mathbf{y}_{i,k}^+|\mathbf{x},\mathbf{y}_{i,< k}^+)$ and $\pi_{\text{old},< k} = \pi_{\theta_{\text{old}}}(\mathbf{y}_{i,k}^+|\mathbf{x},\mathbf{y}_{i,< k}^+)$ for brevity , we can obtain:

$$\frac{d}{dt} \ln \pi_{\theta(t)}(\boldsymbol{y}_{i}^{+}|\boldsymbol{x}) = \left\langle \nabla \ln \pi_{\theta(t)}(\boldsymbol{y}_{i}^{+}|\boldsymbol{x}), \frac{d\theta(t)}{dt} \right\rangle$$

$$= \left\langle \nabla \ln \pi_{\theta(t)}(\boldsymbol{y}_{i}^{+}|\boldsymbol{x}), \right.$$

$$p^{+} \sum_{i'=1}^{N^{+}} \sum_{k=1}^{|\boldsymbol{y}_{i'}^{+}|} \frac{\pi_{< k}}{\pi_{\text{old},< k}} \delta \left(1 + \epsilon - \frac{\pi_{< k}}{\pi_{\text{old},< k}} \right) \nabla \ln \pi_{\theta(t)}(\boldsymbol{y}_{i',k}^{+}|\boldsymbol{x}, \boldsymbol{y}_{i',< k}^{+})$$

$$- p^{-} \sum_{j=1}^{N^{-}} \sum_{k'=1}^{|\boldsymbol{y}_{j}^{-}|} \frac{\pi_{< k'}}{\pi_{\text{old},< k'}} \delta \left(\frac{\pi_{< k'}}{\pi_{\text{old},< k'}} - 1 + \epsilon \right) \nabla \ln \pi_{\theta(t)}(\boldsymbol{y}_{j,k'}^{-}|\boldsymbol{x}, \boldsymbol{y}_{j,< k'}^{-}) \right\rangle$$
(15)

where $\delta(x) = \begin{cases} 1 & \text{if } x \geq 0 \\ 0 & \text{if } x < 0 \end{cases}$ is the gradient of the clip function. We assume $\pi_{\theta(t)} = \pi_{\text{old}}$ since GRPO is online or nearly online due to sampling new responses at each iteration. Thus we obtain:

$$\frac{d}{dt} \ln \pi_{\theta(t)}(\boldsymbol{y}_i^+|\boldsymbol{x}) = \left\langle \nabla \ln \pi_{\theta(t)}(\boldsymbol{y}_i^+|\boldsymbol{x}), p^+ \sum_{i'=1}^{N^+} \nabla \ln \pi_{\theta(t)}(\boldsymbol{y}_{i'}^+|\boldsymbol{x}) - p^- \sum_{i=1}^{N^-} \nabla \ln \pi_{\theta(t)}(\boldsymbol{y}_j^-|\boldsymbol{x}) \right\rangle.$$

As per the unconstrained features Assumption 4.3, the model's trainable parameters are

$$\theta = \left(\boldsymbol{W} \,, \mathbf{h}_{\boldsymbol{x}} \,,\, \left\{ \mathbf{h}_{\boldsymbol{x},\boldsymbol{y}_{i',< k}^+} \right\}_{i' \in [N^+], k \in \{2,...,|\boldsymbol{y}_{i'}^+|\}} \,,\, \left\{ \mathbf{h}_{\boldsymbol{x},\boldsymbol{y}_{i,< k'}^-} \right\}_{j \in [N^-], k' \in \{2,...,|\boldsymbol{y}_{i}^-|\}} \right).$$

Here, we also used the (mild) assumption that all responses differ in their first token. Unfolding the gradients with respect to these parameters yields:

$$\frac{d}{dt} \ln \pi_{\theta(t)}(\boldsymbol{y}_{i}^{+}|\boldsymbol{x}) = \left\langle \nabla_{\boldsymbol{W}} \ln \pi_{\theta(t)}(\boldsymbol{y}_{i}^{+}|\boldsymbol{x}), p^{+} \sum_{i'=1}^{N^{+}} \nabla_{\boldsymbol{W}} \ln \pi_{\theta(t)}(\boldsymbol{y}_{i'}^{+}|\boldsymbol{x}) - p^{-} \sum_{j}^{N^{-}} \nabla_{\boldsymbol{W}} \ln \pi_{\theta(t)}(\boldsymbol{y}_{j}^{-}|\boldsymbol{x}) \right\rangle
+ \left\langle \nabla_{\mathbf{h}_{\mathbf{x}}} \ln \pi_{\theta(t)}(\boldsymbol{y}_{i,1}^{+}|\boldsymbol{x}), p^{+} \sum_{i'=1}^{N^{+}} \nabla_{\mathbf{h}_{\mathbf{x}}} \ln \pi_{\theta(t)}(\boldsymbol{y}_{i',1}^{+}|\boldsymbol{x}) - p^{-} \sum_{j=1}^{N^{-}} \nabla_{\mathbf{h}_{\mathbf{x}}} \ln \pi_{\theta(t)}(\boldsymbol{y}_{j,1}^{-}|\boldsymbol{x}) \right\rangle
+ \sum_{k=2}^{|\boldsymbol{y}_{i}^{+}|} ||\nabla_{\mathbf{h}_{\boldsymbol{x},\boldsymbol{y}_{i,
(16)$$

For softmax model output the gradients can be easily computed as follows:

$$abla_{oldsymbol{W}} \ln \pi_{ heta(t)}(oldsymbol{z} | oldsymbol{x}) = \sum_{k=1}^{|oldsymbol{z}|} \left(\mathbf{e}_{oldsymbol{z}_k} - \pi_{ heta(t)}(\cdot | oldsymbol{x}, oldsymbol{z}_{< k})
ight) \mathbf{h}_{oldsymbol{z}_{< k}}^{ op}(t)$$

In addition, the gradient with respect to the hidden representation at each position k is:

$$\nabla_{\mathbf{h}_{\boldsymbol{x},\boldsymbol{z}_{< k}}} \ln \pi_{\theta(t)}(\boldsymbol{z}|\boldsymbol{x}) = \boldsymbol{w}_{\boldsymbol{z}_k}(t) - \sum_{\boldsymbol{z} \in \mathcal{V}} \pi_{\theta(t)}(\boldsymbol{z}|\mathbf{x},\boldsymbol{z}_{< k}) \cdot \boldsymbol{w}_{\boldsymbol{z}}(t), \quad k \in \{1,\ldots,|\boldsymbol{z}|\}$$

Putting this back in (16) together with a few algebra steps, yields

$$\frac{d}{dt}\ln \pi_{\theta(t)}(\boldsymbol{y}_i^+|\boldsymbol{x}) = (I) - (II) + (III) + (IV)$$
(17)

where:

$$(I) = p^{+} \sum_{k=1}^{|\boldsymbol{y}_{i}^{+}|} \sum_{i'=1}^{N^{+}} \sum_{k''=1}^{|\boldsymbol{y}_{i'}^{+}|} \alpha_{k,k''}^{+}(t) \cdot \left\langle \mathbf{h}_{\boldsymbol{x},\boldsymbol{y}_{i,< k}^{+}}(t), \mathbf{h}_{\boldsymbol{x},\boldsymbol{y}_{i',< k''}^{+}}(t) \right\rangle$$

$$(18)$$

(II) =
$$p^{-}\sum_{k=1}^{|\boldsymbol{y}_{i}^{+}|}\sum_{i=1}^{N^{-}}\sum_{k'=1}^{|\boldsymbol{y}_{j}^{-}|}\alpha_{k,k'}^{-}(t)\cdot\left\langle \mathbf{h}_{\boldsymbol{x},\boldsymbol{y}_{i,< k}^{+}}(t), \mathbf{h}_{\boldsymbol{x},\boldsymbol{y}_{j,< k'}^{-}}(t)\right\rangle$$
 (19)

$$(III) = \left\langle \boldsymbol{w}_{\boldsymbol{y}_{i,1}^{+}}(t) - \sum_{z \in \mathcal{V}} \pi_{\theta(t)}(z|\boldsymbol{x}) \cdot \boldsymbol{w}_{z}(t), \sum_{i'=1}^{N^{+}} p^{+} \boldsymbol{w}_{\boldsymbol{y}_{i',1}^{+}} - \sum_{j=1}^{N^{-}} p^{-} \boldsymbol{w}_{\boldsymbol{y}_{j,1}^{-}} \right\rangle$$
(20)

$$(IV) = \sum_{k=2}^{|\mathbf{y}^+|} \left\| \mathbf{w}_{y_{i,k}^+}(t) - \sum_{z \in \mathcal{V}} \pi_{\theta(t)}(z|\mathbf{x}, \mathbf{y}_{i,< k}^+) \cdot \mathbf{w}_z(t) \right\|^2$$
(21)

where
$$\alpha^+_{k,k^{\prime\prime}}(t) = \left\langle \mathbf{e}_{oldsymbol{y}_{i,k}^+} - \pi_{\theta(t)}(\cdot|oldsymbol{x}, oldsymbol{y}_{i,< k}^+), \mathbf{e}_{oldsymbol{y}_{i^\prime,k^{\prime\prime}}^+} - \pi_{\theta(t)}(\cdot|oldsymbol{x}, oldsymbol{y}_{i^\prime,< k^{\prime\prime}}^+) \right\rangle$$
 and $\alpha^-_{k,k^\prime}(t) = \left\langle \mathbf{e}_{oldsymbol{y}_{i,k}^+} - \pi_{\theta(t)}(\cdot|oldsymbol{x}, oldsymbol{y}_{i,< k}^-) - \pi_{\theta(t)}(\cdot|oldsymbol{x}, oldsymbol{y}_{i,< k}^-) \right\rangle$.

Specifically, (I) and (II) capture how token embeddings influence changes in likelihood, while (III) and (IV) reflect how the geometry of token unembeddings governs such changes. Our focus is on (I) and (II) because:

- Token embeddings encapsulate the contribution of all network parameters excluding the token unembedding layer.
- Token embeddings are influenced by the words in the sample, which span a broader space than token unembeddings.

Our results in Tab. 1 further validate the effectiveness of these terms. Thus we arrive at Theorem 4.4: as (II) - (I) increases, the likelihood change decreases.

7.3 Implementation Details

Dataset. For training, we use the MATH dataset (levels 3–5) to train the model. Additionally, we include a subset of the DeepScaler dataset [14], which contains more challenging problems. For evaluation, we assess the reasoning capabilities of the fine-tuned models on five standard math benchmarks: AIME 2024 [24], AMC, MATH500, Minerva Math [12], and OlympiadBench [5].

Models. We we select models range from 0.5B-3B [25]. We use 0.5B-ins and 1.5B-ins models as [32] found small base model may not follow the format prompt well. We also use Qwen2.5-Math-1.5B [26] to show the influence of math knowledge. For 3B model, we use Qwen-2.5-3B [25] base model. For all models, we conduct reinforcement fine-tuning using GRPO and NTHR using the same hyperparameters.

Hyperparamters. For the 0.5B model, we use two A6000 GPUs with a batch size of 32, maximum rollout length of 2500 tokens, learning rate $5e^{-7}$ and a mini-batch size of 16, resulting in two iteration updates per training step. We use math dataset to train the model for 105 steps.

For the larger models, we utilize four A100 GPUs with a batch size of 256, learning rate 1e-6 and a mini-batch size of 64, leading to four iteration updates per step, for the math dataset, we train for 40 steps, which approximates one epoch. For the DeepScaler dataset, which contains more data, we train for 45 steps.

Across all models, we generate 8 rollouts per prompt. We use a default sampling temperature of 1.0, a clipping ratio of 0.2, and set the KL loss coefficient to 1×10^{-4} . The Qwen-math model [26] uses its full context length of 3072 tokens for rollouts, while all other models use a maximum rollout length of 4000 tokens. To accelerate training, we employ **dynamic sampling** [29], which filters out samples with zero advantage.

Lastly, we use $\beta=1$ and define the weighting factor as $\eta=2\times|0.5-p|$, where p is the success rate. This formulation penalizes questions with success rates near 0.5 less, avoiding excessive reward for easy samples while not over-relying on difficult questions, whose scarce correct responses may

not yield effective token selection. Finally, we set the temperature to 0 to perform greedy decoding, following standard practice in math evaluation tasks [28, 2, 13]. This enables the assessment of the model's most confident output, which is essential for math tasks requiring exact correctness and providing users with consistent answers to repeated queries.

7.4 Complexity Discussion

In this section, we introduce several techniques to enhance the computational efficiency of NTHR, which is essential for practical implementation.

Output last-layer embeddings with old policy probabilities. Since GRPO relies on the old policy π_{old} to compute the old probabilities used in the $\text{clip}(\cdot)$ operation, we only need to extract the last-layer embeddings during the forward pass of the old policy. This approach introduces no additional forward passes and thus keeps the computational overhead minimal. Moreover, last-layer embeddings have been widely adopted in prior work [33, 15] for constructing robust reward signals.

Calculate summations first. We more conveniently rewrite (7) as a matrix inner product.

$$\left\langle \sum_{i=1}^{N^{+}} \sum_{k=1}^{|\boldsymbol{y}_{i}^{+}|} \left(\mathbf{e}_{\boldsymbol{y}_{i,k}^{+}} - \pi_{\theta(t)}(\cdot | \boldsymbol{x}, \boldsymbol{y}_{i,< k}^{+}) \right) \mathbf{h}_{\mathbf{x}, \boldsymbol{y}_{i,< k}^{+}}^{T}, \left(\mathbf{e}_{\boldsymbol{y}_{j,k'}^{-}} - \pi_{\theta(t)}(\cdot | \boldsymbol{x}, \boldsymbol{y}_{j,< k'}^{-}) \right) \mathbf{h}_{\boldsymbol{x}, \boldsymbol{y}_{j,< k'}^{-}}^{T} \right\rangle$$

Importantly, our reformulation involves calculating the summations over i, k first before taking the inner product. This reduction reduces the overall complexity for compute (7) by 3.

Focus on vocabulary of the responses. As the formulation involves computing the outer product between the prediction error vector (e.g., $\mathbf{e}_{y_{i,k}^+} - \pi_{\theta(t)}(\cdot|\mathbf{x},\mathbf{y}_{i,< k}^+)$) and the hidden embedding, which incurs a computational complexity of $O(|\mathcal{V}|d)$. Since the probability mass is primarily concentrated on the output words, for each question \boldsymbol{x} , we restrict the computation to the vocabulary $\mathcal{V}_{\boldsymbol{x}}^*$ associated with its generated responses. Since $|\mathcal{V}_{\boldsymbol{x}}^*| \ll |\mathcal{V}|$, this significantly lowers the overall cost to $O(|\mathcal{V}_{\boldsymbol{x}}^*|d)$.

Running time of each module. We also track the average time cost of each module during training, as reported in Tab. 3. Notably, the data generation (Data Gen) module that using dynamic sampling accounts for the majority of the total training time. In contrast, the overhead introduced by NTHR is minimal, contributing only a small fraction to the overall cost. Notably, although deepseek-1.5B has a longer average output length of approximately 3,400 tokens, our NTHR still maintains a low time overhead.

| Model+dataset | Data Gen | Model Upd | THR | Ref | Old Prob | Total (Sec) |
|--------------------------|----------|-----------|-----|-----|----------|-------------|
| Qwen2.5-Math-1.5B (Math) | 250 | 140 | 22 | 55 | 55 | 522 |

Table 3: Average running time (per step, in seconds) of each module for different models and tasks.

$$\alpha_{k,k'} = \left\langle \mathbf{e}_{\boldsymbol{y}_{i,k}^{+}} - \pi_{\theta}(\cdot \mid \boldsymbol{x}, \boldsymbol{y}_{i,\leq k}^{+}), \mathbf{e}_{\boldsymbol{y}_{k'}} - \pi_{\theta}(\cdot \mid \boldsymbol{x}, \boldsymbol{y}_{\leq k'}) \right\rangle$$
(22)

As a result, our method does not introduce additional inference overhead, and the overall complexity remains small.

7.5 Performance across training iterations

In this section, we present the performance of various models across training iterations. As illustrated in Fig. 5, although performance fluctuates throughout training, NTHR consistently outperforms GRPO. Notably, for the Qwen-2.5-1.5B-Ins model, we report results at update step 100 (corresponding to 25 training steps), as performance begins to decline beyond this point.

7.6 Ablation study

In this section, we conduct ablation studies to shed more light on the role played by different modules. **Ablation on** β . We conduct ablation study on β , i.e., the scaling factor for the threshold τ in Eq. (8), which is then used to control how bad the negative tokens should be masked out. We set β with $\{-\infty, 0, 0.1, 1, \}$ and train with Qwen2.5-Math-1.5B using Math dataset. Choosing $\beta = -\infty$ corresponds to apply η to all incorrect tokens. As shown in Tab. 4, this removal results in worse

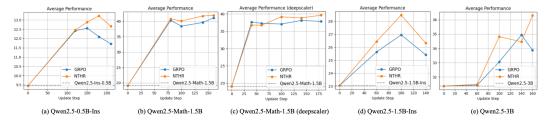


Figure 5: Performance across training iterations for various models, NTHR consistently outperforms GRPO for most of the training process.

performance compared to GRPO. As further shown in Tab. 4, incorporating NTHR consistently improves performance over GRPO. While different β values result in varying degrees of improvement, all settings outperform GRPO, indicating that reducing the influence of identified negative gradients is beneficial.

| Base model + Method | AIME24 | AMC | MATH500 | Minerva | Olympiad | Avg. |
|--------------------------|-------------|--------------|-------------|-------------|-------------|-------|
| GRPO | 13.3 | 57.5 | 71.8 | 29.0 | 34.1 | 41.14 |
| NTHR $(\beta = -\infty)$ | 13.3 | 50.0 | 71.8 | 30.5 | 34.4 | 40.00 |
| NTHR $(\beta = 0)$ | 20.0 | 55.0 | 70.0 | 29.8 | 32.9 | 41.54 |
| NTHR ($\beta = 0.1$) | 13.3 | 57.5 | <u>71.4</u> | 30.9 | 34.4 | 41.50 |
| NTHR ($\beta = 1.0$) | <u>16.7</u> | 57. 5 | 70.8 | <u>30.5</u> | <u>34.2</u> | 41.94 |

Table 4: Ablation study on β . Although changes in β (except ∞) lead to variations in performance, they consistently improve GRPO, indicating that reduce influence of identified negative gradients is beneficial.

Ablation on η . We conduct an ablation study on η using three variants: $\eta = 1 - p$, $\eta = p$, and $\eta = 2 \cdot |0.5 - p|$. As shown in Tab. 5, while the optimal choice of η varies across different models, the variant $\eta = 2 \cdot |0.5 - p|$ consistently yields strong performance across all models. Therefore, we adopt $\eta = 2 \cdot |0.5 - p|$ in our training.

| Base model + Method | AIME24 | AMC | MATH500 | Minerva | Olympiad | Avg. |
|-----------------------------------|--------|------|---------|---------|----------|-------|
| Qwen2.5Math-1.5B | | | | | | |
| GRPO | 13.3 | 57.5 | 71.8 | 29.0 | 34.1 | 41.14 |
| NTHR $(\eta = p)$ | 13.3 | 55.0 | 72.4 | 29.0 | 33.3 | 40.60 |
| NTHR $(\eta = 1 - p)$ | 13.3 | 62.5 | 72.8 | 30.5 | 34.2 | 42.66 |
| NTHR $(\eta = 2 \cdot 0.5 - p)$ | 16.7 | 57.5 | 70.8 | 30.5 | 34.2 | 41.94 |
| Qwen2.5-3B | | | | | | |
| GRPO | 6.7 | 35.0 | 66.6 | 31.2 | 29.9 | 33.88 |
| NTHR $(\eta = p)$ | 10.0 | 47.5 | 64.6 | 33.8 | 26.8 | 36.54 |
| NTHR $(\eta = 1 - p)$ | 6.7 | 42.5 | 63.8 | 32.4 | 28.0 | 34.68 |
| NTHR $(\eta = 2 \cdot 0.5 - p)$ | 10.0 | 47.5 | 65.6 | 31.6 | 26.8 | 36.30 |

Table 5: Ablation study on η .

7.7 Results on Llama

As shown in Fig. 6, the response length of Llama3.2-1B-Instruct declines rapidly after a few epochs, with the average length dropping from about 1.2K tokens to roughly 500. This reduction may stem from the model's limited cognitive behaviors [3], but once equipped with them, they can match Qwen's trajectory of self-improvement [3]. Despite this, We still test NTHR using Llama3.2-1B-Instruct and observed consistent improvements when applying our NTHR method as shown in Tab. 6.

7.8 NLL Loss as a Complementary

Recently, Yue et al. [31] proposed directly applying an NLL loss to positive samples, which aligns closely with the

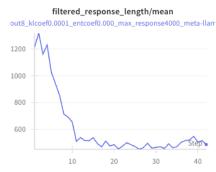


Figure 6: Response length dynamics of Llama3.2-1B-Instruct across different stages of GRPO training.

| Method | math500 | minerva_math | olympiad | aime24 | amc23 | avg |
|------------|---------|--------------|----------|--------|-------|--------------|
| GRPO | 37.0 | 7.4 | 9.3 | 3.3 | 20.0 | 15.4 |
| + NLL Loss | 36.0 | 9.6 | 11.0 | 3.3 | 25.0 | 16.98 |
| + NTHR | 37.4 | 8.1 | 11.0 | 3.3 | 20.0 | <u>15.96</u> |

Table 6: Performance comparison across math reasoning benchmarks on Llama. Bold numbers indicate the best performance per column.

LLD issue we identify, as it increases the likelihood of correct responses. Specifically, incorporating a supervised NLL term reinforces the model's confidence in generating correct outputs. Motivated by this, we augmented GRPO with an additional NLL loss and trained it using Llama-3.2-1B-Instruct. As shown in Tab. 6, this modification leads to improved performance across multiple tasks. The NLL objective tackles LLD from a complementary perspective—strengthening correct responses through positive gradients—while NTHR mitigates harmful negative updates. Both methods consistently yield performance gains.

7.9 Results on Deepseek-1.5B

We utilize four A100 GPUs with a batch size of 256 and a mini-batch size of 64, resulting in four update iterations per training step. Each prompt generates 8 rollouts. We adopt a default sampling temperature of 0.6, set the clipping ratio to 0.2, and use a KL loss coefficient of 1×10^{-4} . The maximum rollout length is capped at 4000 tokens. We set $\beta=0$ and define the weighting factor as $\eta=2\cdot|0.5-p|$. We limit training to 35 steps, as we observed a consistent decline in response length with continued training.

Since the context window length is limited, the results are very unstable, as a results, we calculate the last two checkpoints' average performance and report the results in the table below. Although this setting may not be optimal or robust, we include the results and observe that our NTHR consistently outperforms GRPO on average.

| Base model + Method | AIME24 | AMC | MATH500 | Minerva | Olympiad | Avg. |
|---------------------|--------|------|---------|---------|----------|------|
| Deepseek-1.5B | | | | | | |
| Base | 20.0 | 50.0 | 68.8 | 20.6 | 29.6 | 37.8 |
| GRPO | 25.0 | 65.0 | 81.10 | 34.6 | 40.0 | 49.1 |
| NTHR | 25.0 | 70.0 | 80.60 | 33.7 | 40.7 | 50.0 |

Table 7: Results across selected math benchmarks for deepseek-1.5B.

7.10 Examples of LLD Incorrect Responses

We present the top-ranked incorrect response based on LLD for Qwen-2.5Math-1.5B in Fig. 1c. As shown in Tab. 8, the answer is nearly correct, with the only issue occurring in the final step.

Similarly, Fig. 1b shows the top-ranked incorrect response by LLD for Deepseek-1.5B. As detailed in Tab. 9, the response includes the correct answer multiple times; however, due to continued self-reflection, it becomes truncated by the context window limit, resulting in a final answer that is wrong.

7.11 Words with high NTHR

We present additional examples of words with high NTHR scores. As shown in Tab. 10, using the same response from Tab. 8, we observe that the correct parts are associated with high NTHR values, while the incorrect portion in step 5 receives low NTHR values (in black). We also show response with incorrect output format in Tab. 11, using the same response from Fig. 2, we observe that core steps and correct answer 8 are with high NTHR.

| Model & Data | Problem & Answer | Incorrect Response |
|----------------|--|--|
| Qwen2.5math- | Question: A portion of the | To solve the problem, we need to analyze the given functions |
| 1.5B | graph of a quadratic function | f(x), $g(x)$, and $h(x)$ and determine the points of intersection |
| Math (order 1) | f(x) is shown. | between the graphs of these functions. |
| | Let $g(x) = -f(x)$, $h(x) =$ | Graph of $f(x)$: The graph of $f(x)$ is a quadratic function. From the given |
| | f(-x). If a is the number of | graph, we can observe the following: |
| | points where the graphs of $y = f(x)$ and $y = g(x)$ intersect, | • The roots of $f(x)$ are at $x = 1$ and $x = 3$. |
| | and b is the number of points | • The vertex of the parabola is at $x=2$. |
| | and b is the number of points where the graphs of $y = f(x)$ and $y = h(x)$ intersect, what is 10a + b? | Graph of $g(x) = -f(x)$: The graph of $g(x)$ is the reflection of the graph of $f(x)$ across the x-axis. Therefore: |
| | Answer: 21 | • The roots of $g(x)$ are also at $x = 1$ and $x = 3$. |
| | | • The vertex of the parabola for $g(x)$ is at $x = 2$. |
| | | Graph of $h(x) = f(-x)$: The graph of $h(x)$ is the reflection of the graph of $f(x)$ across the y-axis. Therefore: |
| | | • The roots of $h(x)$ are at $x = -1$ and $x = -3$. |
| | | • The vertex of the parabola for $h(x)$ is at $x = -2$. |
| | | Next, we determine the points of intersection: Intersection of $f(x)$ and $g(x)$: Since $g(x) = -f(x)$, the graphs of $f(x)$ and $g(x)$ intersect at the points where $f(x) = -f(x)$. This simplifies to: |
| | | f(x) = 0 |
| | | The roots of $f(x)$ are at $x=1$ and $x=3$. Therefore, the graphs of $f(x)$ and $g(x)$ intersect at two points: |
| | | a = 2 |
| | | Intersection of $f(x)$ and $h(x)$: Since $h(x) = f(-x)$, the graphs of $f(x)$ and $h(x)$ intersect at the points where $f(x) = f(-x)$. This implies that $f(x)$ is an even function, and its graph is symmetric about the y-axis. The roots of $f(x)$ are at $x = 1$ and $x = 3$, but these are not symmetric about the y-axis. Therefore, the graphs of $f(x)$ and $h(x)$ do not intersect. Thus: $b = 0$ Finally, we calculate $10a + b$: $10a + b = 10 \cdot 2 + 0 = 20$ Thus, the final answer is: |
| | | 20 |

Table 8: Incorrect response example from qwen2.5math-1.5B. The incorrect response is almost correct.

| Model & Data | Problem & Answer | Incorrect Response |
|----------------|---|---|
| Deepseek-1.5B | Question: One commercially | To solve the problem, explain your reasoning step by step, |
| AIME (order 1) | available ten-button lock may be | and put your final answer within boxed. Alright, so I have this |
| | opened by pressing – in any or- | problem about ten-button locks and figuring out how many |
| | der – the correct five buttons. | additional combinations a redesign would allow. Let me try |
| | The sample shown below has | to break it down step by step. First, the original lock. It seems |
| | $\{1, 2, 3, 6, 9\}$ as its combination. | like a ten-button lock where you need to press five buttons |
| | Suppose that these locks are redesigned so that sets of as many | in any order to open it. The combination given is 1,2,3,6,9, which is a set of five buttons. So, the question is asking about |
| | as nine buttons or as few as | how many additional combinations would be possible if we |
| | one button could serve as com- | redesign the locks to allow sets of as many as nine buttons |
| | binations. How many additional | or as few as one button. I think I need to figure out the total |
| | combinations would this allow? | number of possible combinations before and after the redesign |
| | Answer: 770 | and then find the difference. That should give me the number |
| | | of additional combinations. Starting with the original lock |
| | | TIL C d L C LIV. 1 L V . 1000 |
| | | Therefore, the number of additional combinations is 1022 – |
| | | 252 = 770. Wait, that sthe same answer as before. So, perhaps 770 is the correct answer. But wait, let me make sure. Is the |
| | | original number of combinations one, or 252? |
| | | |
| | | Then, the redesigned lock allows for any number of buttons |
| | | from one to nine, so $2^10 - 2 = 1022$ possible combinations, |
| | | but the original combination was just one. So, the number |
| | | of additional combinations would be $1022 - 1 = 1021$. But |
| | | that seems like a lot more. Wait, which is it? Is the original combination one specific combination, or is it 252 unique |
| | | combination one specific combination, or is a 252 unique combinations? In the problem, it says, "one commercially |
| | | available ten-button lock may be opened by pressing – in any |
| | | order – the correct five buttons." |
| | | |
| | | After redesign, the number of possible combinations is 1022. |
| | | So, the number of additional combinations is $1022-1 = 1021$. |
| | | But wait, the redesign allows for any combination from one to |
| | | nine buttons. So, the number of possible combinations is 1022, which includes the original 252 combinations. Therefore, the |
| | | number of additional combinations is $1022 - 252 = 770$. |
| | | indiffice of additional combinations is 1022 202 – 110. |
| | | The redesigned allows for more combinations, but if you |
| | | consider that each combination is a set of buttons, then the |
| | | number of additional combinations is $1022 - 252 = 770$. But |
| | | the problem says, "the sample shown below has 1,2,3,6,9 as |
| | | its combination." So, the sample is a specific combination, so |
| | | that's one combination. So, originally, the number of possible |
| | | combinations was 252, including the correct one. So, if you redesign it, the number of possible combinations is 1022, |
| | | which includes all the original 252 combinations plus 1022 – |
| | | 252 = 770 additional combinations. Therefore, the number |
| | | of additional combinations is 770. But to make sure, let me |
| | | think of another approach. If the original lock had only one |
| | | combination, the number of possible combinations is 252, |
| | | each being a unique set of five buttons. Then, if you redesign |
| | | it so that you can have any number of buttons from one to |
| | | nine, the number of possible combinations is Stop due to limited context window length. |

Table 9: Incorrect response example from deepseek-1.5B on AIME: Although the response includes the correct answers, excessive self-reflection causes the output to exceed the context window, preventing it from being formatted correctly. Green dashes means omitted reasoning.

| Model & Data | Problem & Answer | Incorrect Response |
|------------------|--|--|
| Qwen2.5math- | Question: A portion of the | To solve the problem, we need to analyze the given functions |
| 1.5B | graph of a quadratic function | f(x), $g(x)$, and $h(x)$ and determine the points of intersection |
| Math | f(x) is shown. | between the graphs of these functions. |
| | | 1. Graph of $f(x)$: The graph of $f(x)$ is a quadratic function. |
| | Let $g(x) = -f(x)$, $h(x) =$ | From the given graph, we can observe the following: - The |
| | f(-x). If a is the number of | roots of $f(x)$ are at $x = 1$ and $x = 3$ The vertex of the |
| | points where the graphs of $y =$ | parabola is at $x = 2$. |
| | f(x) and $y = g(x)$ intersect, | 2. Graph of $g(x) = -f(x)$: The graph of $g(x)$ is the reflection |
| | and b is the number of points where the graphs of $y = f(x)$ | of the graph of $f(x)$ across the x-axis. Therefore: - The roots of $g(x)$ are also at $x = 1$ and $x = 3$ The vertex of the |
| | and $y = h(x)$ intersect, what is | parabola for $g(x)$ is at $x = 2$. |
| | $\begin{array}{c c} and \ y = h(x) \text{ intersect, what is} \\ 10a + b? \end{array}$ | 3. Graph of $h(x) = f(-x)$:** The graph of $h(x)$ is the reflec- |
| | Answer: 21 | tion of the graph of $f(x)$ across the y-axis. Therefore: - The |
| | | roots of $h(x)$ are at $x = -1$ and $x = -3$ The vertex of the |
| | | parabola for $h(x)$ is at $x = -2$. |
| | | Next, we determine the points of intersection: |
| | | 4. Intersection of $f(x)$ and $g(x)$: Since $g(x) = -f(x)$, |
| | | the graphs of $f(x)$ and $g(x)$ intersect at the points where |
| | | f(x) = -f(x). This simplifies to: |
| | | f(x) = 0 |
| | | The roots of $f(x)$ are at $x = 1$ and $x = 3$. Therefore, the |
| | | graphs of $f(x)$ and $g(x)$ intersect at two points: |
| | | a=2 |
| | | 5. Intersection of $f(x)$ and $h(x)$: Since $h(x) = f(-x)$, the graphs of $f(x)$ and $h(x)$ intersect at the points where |
| | | f(x) = f(-x). This implies that $f(x)$ is an even function, |
| | | and its graph is symmetric about the y-axis. The roots of $f(x)$ |
| | | are at $x = 1$ and $x = 3$, but these are not symmetric about |
| | | the y-axis. Therefore, the graphs of $f(x)$ and $h(x)$ do not |
| | | intersect. |
| | | Thus: |
| | | b = 0 |
| | | Finally, we calculate $10a + b$: |
| | | $10a + b = 10 \cdot 2 + 0 = 20$ |
| | | Thus, the final answer is: |
| | | 20 |
| Table 10: High N | THR (highlight in red) exam | nple from gwen2.5math-1.5B. Core steps in correct |

Table 10: High NTHR (highlight in red) example from qwen2.5math-1.5B. Core steps in correct parts are with high NTHR. The erroneous words of step 5 are with low NTHR values.

| Model & Data | Problem & Answer | Incorrect Response |
|---------------|--|---|
| Qwen2.5-0.5B- | Question : What is the greatest | To solve the inequality $n^2 - 11n + 24 \le 0$, we need to |
| Ins | integer n such that $n^2 - 11n +$ | find the values of n that satisfy the inequality. Since the |
| Math | $24 \le 0$? | inequality is an inequality, we should first find the roots of the |
| | Answer: 8 | corresponding equation $n^2 - 11n + 24 = 0$. Let's solve the |
| | | equation $n^2 - 11n + 24 = 0$: 1. Factor the quadratic equation: |
| | | $n^2 - 11n + 24 = (n-3)(n-8) = 0$. 2. Set each factor equal |
| | | to zero and solve for n: $n-3=0$ or $n-8=0$. 3. Solve |
| | | for n: $n = 3$ or $n = 8$. Therefore, the roots of the equation |
| | | $n^2-11n+24=0$ are $n=3$ and $n=8$. Now, we need to find |
| | | the values of n that satisfy the inequality $n^2 - 11n + 24 \le 0$. |
| | | Since the inequality is an inequality, we should check the |
| | | intervals between the roots. The roots $n=3$ and $n=8$ |
| | | divide the number line into three intervals: $(-\infty, 3]$, $[3, 8]$, |
| | | and $[8,\infty)$. Let 's test a point in each interval to determine |
| | | where the inequality holds true. 1. For n in $(-\infty, 3)$, let's test |
| | | $n = 0$: $0^2 - 11(0) + 24 = 24 > 0$. 2. For n in (3, 8), let 's test |
| | | $n = 5$: $5^2 - 11(5) + 24 = 25 - 55 + 24 = -4 < 0.3$. For n |
| | | in $(8, \infty)$, let 's test $n = 9$: $9^2 - 11(9) + 24 = 81 - 99 + 24 = 91 + 11(9) + 124 = 11(9) + 124$ |
| | | $6 > 0$. Therefore, the inequality $n^2 - 11n + 24 < 0$ holds |
| | | true for $3 \le n \le 8$. Hence, the greatest integer n that satisfies |
| | | the inequality is $n = 8$. In boxed format, the answer is 8. |
| | 1 | are mequanty is to the content format, the this wer is of |

Table 11: High NTHR (highlight in red) example from qwen2.5-0.5B-Ins. This example have correct answer but in wrong format. Core steps and correct answer 8 are with high NTHR.

NeurIPS Paper Checklist

1. Claims

Question: Do the main claims made in the abstract and introduction accurately reflect the paper's contributions and scope?

Answer: [Yes]

Justification: We give clear motivation and summarized the contribution of our work in introduction.

Guidelines:

- The answer NA means that the abstract and introduction do not include the claims made in the paper.
- The abstract and/or introduction should clearly state the claims made, including the contributions made in the paper and important assumptions and limitations. A No or NA answer to this question will not be perceived well by the reviewers.
- The claims made should match theoretical and experimental results, and reflect how much the results can be expected to generalize to other settings.
- It is fine to include aspirational goals as motivation as long as it is clear that these goals are not attained by the paper.

2. Limitations

Question: Does the paper discuss the limitations of the work performed by the authors?

Answer: [Yes]

Justification: In conclusion, we discussed the limitaions.

Guidelines:

- The answer NA means that the paper has no limitation while the answer No means that the paper has limitations, but those are not discussed in the paper.
- The authors are encouraged to create a separate "Limitations" section in their paper.
- The paper should point out any strong assumptions and how robust the results are to violations of these assumptions (e.g., independence assumptions, noiseless settings, model well-specification, asymptotic approximations only holding locally). The authors should reflect on how these assumptions might be violated in practice and what the implications would be.
- The authors should reflect on the scope of the claims made, e.g., if the approach was only tested on a few datasets or with a few runs. In general, empirical results often depend on implicit assumptions, which should be articulated.
- The authors should reflect on the factors that influence the performance of the approach. For example, a facial recognition algorithm may perform poorly when image resolution is low or images are taken in low lighting. Or a speech-to-text system might not be used reliably to provide closed captions for online lectures because it fails to handle technical jargon.
- The authors should discuss the computational efficiency of the proposed algorithms and how they scale with dataset size.
- If applicable, the authors should discuss possible limitations of their approach to address problems of privacy and fairness.
- While the authors might fear that complete honesty about limitations might be used by reviewers as grounds for rejection, a worse outcome might be that reviewers discover limitations that aren't acknowledged in the paper. The authors should use their best judgment and recognize that individual actions in favor of transparency play an important role in developing norms that preserve the integrity of the community. Reviewers will be specifically instructed to not penalize honesty concerning limitations.

3. Theory assumptions and proofs

Question: For each theoretical result, does the paper provide the full set of assumptions and a complete (and correct) proof?

Answer: [Yes]

Justification: We provide the main theory in main paper and proof and assumptions in appendix.

Guidelines:

- The answer NA means that the paper does not include theoretical results.
- All the theorems, formulas, and proofs in the paper should be numbered and crossreferenced.
- All assumptions should be clearly stated or referenced in the statement of any theorems.
- The proofs can either appear in the main paper or the supplemental material, but if they appear in the supplemental material, the authors are encouraged to provide a short proof sketch to provide intuition.
- Inversely, any informal proof provided in the core of the paper should be complemented by formal proofs provided in appendix or supplemental material.
- Theorems and Lemmas that the proof relies upon should be properly referenced.

4. Experimental result reproducibility

Question: Does the paper fully disclose all the information needed to reproduce the main experimental results of the paper to the extent that it affects the main claims and/or conclusions of the paper (regardless of whether the code and data are provided or not)?

Answer: [Yes]

Justification: We provide algorithms and details in appendix to ensure reproducibility.

Guidelines:

- The answer NA means that the paper does not include experiments.
- If the paper includes experiments, a No answer to this question will not be perceived well by the reviewers: Making the paper reproducible is important, regardless of whether the code and data are provided or not.
- If the contribution is a dataset and/or model, the authors should describe the steps taken to make their results reproducible or verifiable.
- Depending on the contribution, reproducibility can be accomplished in various ways. For example, if the contribution is a novel architecture, describing the architecture fully might suffice, or if the contribution is a specific model and empirical evaluation, it may be necessary to either make it possible for others to replicate the model with the same dataset, or provide access to the model. In general, releasing code and data is often one good way to accomplish this, but reproducibility can also be provided via detailed instructions for how to replicate the results, access to a hosted model (e.g., in the case of a large language model), releasing of a model checkpoint, or other means that are appropriate to the research performed.
- While NeurIPS does not require releasing code, the conference does require all submissions to provide some reasonable avenue for reproducibility, which may depend on the nature of the contribution. For example
- (a) If the contribution is primarily a new algorithm, the paper should make it clear how to reproduce that algorithm.
- (b) If the contribution is primarily a new model architecture, the paper should describe the architecture clearly and fully.
- (c) If the contribution is a new model (e.g., a large language model), then there should either be a way to access this model for reproducing the results or a way to reproduce the model (e.g., with an open-source dataset or instructions for how to construct the dataset).
- (d) We recognize that reproducibility may be tricky in some cases, in which case authors are welcome to describe the particular way they provide for reproducibility. In the case of closed-source models, it may be that access to the model is limited in some way (e.g., to registered users), but it should be possible for other researchers to have some path to reproducing or verifying the results.

5. Open access to data and code

Question: Does the paper provide open access to the data and code, with sufficient instructions to faithfully reproduce the main experimental results, as described in supplemental material?

Answer: [Yes]

Justification: The data are all opensource data, we plan to open-source all the code after the paper is accepted.

Guidelines:

- The answer NA means that paper does not include experiments requiring code.
- Please see the NeurIPS code and data submission guidelines (https://nips.cc/public/guides/CodeSubmissionPolicy) for more details.
- While we encourage the release of code and data, we understand that this might not be possible, so "No" is an acceptable answer. Papers cannot be rejected simply for not including code, unless this is central to the contribution (e.g., for a new open-source benchmark).
- The instructions should contain the exact command and environment needed to run to reproduce the results. See the NeurIPS code and data submission guidelines (https://nips.cc/public/guides/CodeSubmissionPolicy) for more details.
- The authors should provide instructions on data access and preparation, including how
 to access the raw data, preprocessed data, intermediate data, and generated data, etc.
- The authors should provide scripts to reproduce all experimental results for the new
 proposed method and baselines. If only a subset of experiments are reproducible, they
 should state which ones are omitted from the script and why.
- At submission time, to preserve anonymity, the authors should release anonymized versions (if applicable).
- Providing as much information as possible in supplemental material (appended to the paper) is recommended, but including URLs to data and code is permitted.

6. Experimental setting/details

Question: Does the paper specify all the training and test details (e.g., data splits, hyper-parameters, how they were chosen, type of optimizer, etc.) necessary to understand the results?

Answer: [Yes]

Justification: We report the details in appendix

Guidelines:

- The answer NA means that the paper does not include experiments.
- The experimental setting should be presented in the core of the paper to a level of detail that is necessary to appreciate the results and make sense of them.
- The full details can be provided either with the code, in appendix, or as supplemental material.

7. Experiment statistical significance

Question: Does the paper report error bars suitably and correctly defined or other appropriate information about the statistical significance of the experiments?

Answer: [Yes]

Justification: We run on different models, some model we report the average performance.

- The answer NA means that the paper does not include experiments.
- The authors should answer "Yes" if the results are accompanied by error bars, confidence intervals, or statistical significance tests, at least for the experiments that support the main claims of the paper.
- The factors of variability that the error bars are capturing should be clearly stated (for example, train/test split, initialization, random drawing of some parameter, or overall run with given experimental conditions).
- The method for calculating the error bars should be explained (closed form formula, call to a library function, bootstrap, etc.)
- The assumptions made should be given (e.g., Normally distributed errors).

- It should be clear whether the error bar is the standard deviation or the standard error
 of the mean.
- It is OK to report 1-sigma error bars, but one should state it. The authors should preferably report a 2-sigma error bar than state that they have a 96% CI, if the hypothesis of Normality of errors is not verified.
- For asymmetric distributions, the authors should be careful not to show in tables or figures symmetric error bars that would yield results that are out of range (e.g. negative error rates).
- If error bars are reported in tables or plots, The authors should explain in the text how they were calculated and reference the corresponding figures or tables in the text.

8. Experiments compute resources

Question: For each experiment, does the paper provide sufficient information on the computer resources (type of compute workers, memory, time of execution) needed to reproduce the experiments?

Answer: [Yes]

Justification: Yes, the paper provides sufficient information on the computer resources needed to reproduce the experiments in appendix

Guidelines:

- The answer NA means that the paper does not include experiments.
- The paper should indicate the type of compute workers CPU or GPU, internal cluster, or cloud provider, including relevant memory and storage.
- The paper should provide the amount of compute required for each of the individual experimental runs as well as estimate the total compute.
- The paper should disclose whether the full research project required more compute than the experiments reported in the paper (e.g., preliminary or failed experiments that didn't make it into the paper).

9. Code of ethics

Question: Does the research conducted in the paper conform, in every respect, with the NeurIPS Code of Ethics https://neurips.cc/public/EthicsGuidelines?

Answer: [Yes]

Justification: Yes, the research conducted in the paper conforms, in every respect, with the NeurIPS Code of Ethics.

Guidelines:

- The answer NA means that the authors have not reviewed the NeurIPS Code of Ethics.
- If the authors answer No, they should explain the special circumstances that require a
 deviation from the Code of Ethics.
- The authors should make sure to preserve anonymity (e.g., if there is a special consideration due to laws or regulations in their jurisdiction).

10. Broader impacts

Question: Does the paper discuss both potential positive societal impacts and negative societal impacts of the work performed?

Answer: [Yes]

Justification: this can help achieve better LLM reasoning.

- The answer NA means that there is no societal impact of the work performed.
- If the authors answer NA or No, they should explain why their work has no societal impact or why the paper does not address societal impact.
- Examples of negative societal impacts include potential malicious or unintended uses (e.g., disinformation, generating fake profiles, surveillance), fairness considerations (e.g., deployment of technologies that could make decisions that unfairly impact specific groups), privacy considerations, and security considerations.

- The conference expects that many papers will be foundational research and not tied to particular applications, let alone deployments. However, if there is a direct path to any negative applications, the authors should point it out. For example, it is legitimate to point out that an improvement in the quality of generative models could be used to generate deepfakes for disinformation. On the other hand, it is not needed to point out that a generic algorithm for optimizing neural networks could enable people to train models that generate Deepfakes faster.
- The authors should consider possible harms that could arise when the technology is being used as intended and functioning correctly, harms that could arise when the technology is being used as intended but gives incorrect results, and harms following from (intentional or unintentional) misuse of the technology.
- If there are negative societal impacts, the authors could also discuss possible mitigation strategies (e.g., gated release of models, providing defenses in addition to attacks, mechanisms for monitoring misuse, mechanisms to monitor how a system learns from feedback over time, improving the efficiency and accessibility of ML).

11. Safeguards

Question: Does the paper describe safeguards that have been put in place for responsible release of data or models that have a high risk for misuse (e.g., pretrained language models, image generators, or scraped datasets)?

Answer: [NA]
Justification:
Guidelines:

- The answer NA means that the paper poses no such risks.
- Released models that have a high risk for misuse or dual-use should be released with necessary safeguards to allow for controlled use of the model, for example by requiring that users adhere to usage guidelines or restrictions to access the model or implementing safety filters.
- Datasets that have been scraped from the Internet could pose safety risks. The authors should describe how they avoided releasing unsafe images.
- We recognize that providing effective safeguards is challenging, and many papers do not require this, but we encourage authors to take this into account and make a best faith effort.

12. Licenses for existing assets

Question: Are the creators or original owners of assets (e.g., code, data, models), used in the paper, properly credited and are the license and terms of use explicitly mentioned and properly respected?

Answer: [Yes]

Justification: All foundation models and dataets used are properly cited.

- The answer NA means that the paper does not use existing assets.
- The authors should cite the original paper that produced the code package or dataset.
- The authors should state which version of the asset is used and, if possible, include a URL.
- The name of the license (e.g., CC-BY 4.0) should be included for each asset.
- For scraped data from a particular source (e.g., website), the copyright and terms of service of that source should be provided.
- If assets are released, the license, copyright information, and terms of use in the package should be provided. For popular datasets, paperswithcode.com/datasets has curated licenses for some datasets. Their licensing guide can help determine the license of a dataset.
- For existing datasets that are re-packaged, both the original license and the license of the derived asset (if it has changed) should be provided.

• If this information is not available online, the authors are encouraged to reach out to the asset's creators.

13. New assets

Question: Are new assets introduced in the paper well documented and is the documentation provided alongside the assets?

Answer: [NA]

Justification: The paper does not introduce any new assets such as code, data, or models. Therefore, there is no documentation provided alongside new assets.

Guidelines:

- The answer NA means that the paper does not release new assets.
- Researchers should communicate the details of the dataset/code/model as part of their submissions via structured templates. This includes details about training, license, limitations, etc.
- The paper should discuss whether and how consent was obtained from people whose asset is used.
- At submission time, remember to anonymize your assets (if applicable). You can either create an anonymized URL or include an anonymized zip file.

14. Crowdsourcing and research with human subjects

Question: For crowdsourcing experiments and research with human subjects, does the paper include the full text of instructions given to participants and screenshots, if applicable, as well as details about compensation (if any)?

Answer: [NA].

Justification: The paper does not involve crowdsourcing experiments or research with human subjects. Therefore, it does not include instructions given to participants, screenshots, or details about compensation.

Guidelines:

- The answer NA means that the paper does not involve crowdsourcing nor research with human subjects.
- Including this information in the supplemental material is fine, but if the main contribution of the paper involves human subjects, then as much detail as possible should be included in the main paper.
- According to the NeurIPS Code of Ethics, workers involved in data collection, curation, or other labor should be paid at least the minimum wage in the country of the data collector.

15. Institutional review board (IRB) approvals or equivalent for research with human subjects

Question: Does the paper describe potential risks incurred by study participants, whether such risks were disclosed to the subjects, and whether Institutional Review Board (IRB) approvals (or an equivalent approval/review based on the requirements of your country or institution) were obtained?

Answer: [NA]

Justification: The paper does not involve any study participants, crowdsourcing experiments, or research with human subjects. Therefore, it does not describe potential risks incurred by study participants, disclose such risks to subjects, or obtain Institutional Review Board (IRB) approvals.

- The answer NA means that the paper does not involve crowdsourcing nor research with human subjects.
- Depending on the country in which research is conducted, IRB approval (or equivalent) may be required for any human subjects research. If you obtained IRB approval, you should clearly state this in the paper.

- We recognize that the procedures for this may vary significantly between institutions and locations, and we expect authors to adhere to the NeurIPS Code of Ethics and the guidelines for their institution.
- For initial submissions, do not include any information that would break anonymity (if applicable), such as the institution conducting the review.

16. Declaration of LLM usage

Question: Does the paper describe the usage of LLMs if it is an important, original, or non-standard component of the core methods in this research? Note that if the LLM is used only for writing, editing, or formatting purposes and does not impact the core methodology, scientific rigorousness, or originality of the research, declaration is not required.

Answer: [NA]

Justification: we don't rely on LLM in developing our method.

- The answer NA means that the core method development in this research does not involve LLMs as any important, original, or non-standard components.
- Please refer to our LLM policy (https://neurips.cc/Conferences/2025/LLM) for what should or should not be described.