

# Logit Arithmetic Elicits Long Reasoning Capabilities Without Training

Anonymous ACL submission

## Abstract

Large reasoning models (LRMs) can do complex reasoning via long chain-of-thought (CoT) involving cognitive strategies such as backtracking and self-correction. Recent studies suggest that some models inherently possess these long reasoning abilities, which may be unlocked via extra training. Our work investigates whether we can elicit such behavior without *any* training. To that goal, we propose a decoding-time approach, **THINKLOGIT**, which utilizes logits arithmetic (Liu et al., 2024) to tune a target large LM for long reasoning using a substantially smaller model as the guider. We then show that we can further boost its performance by training the guider model with preference optimization over correct/incorrect reasoning pairs sampled from both the target and guider model—a setup we refer to as **THINKLOGIT-DPO**. Our experiments demonstrate that THINKLOGIT and THINKLOGIT-DPO achieve a relative improvement in pass@1 by 24.5% and 29.1%, respectively, over five mathematical and scientific reasoning datasets using the Qwen2.5-32B when guided by R1-Distill-Qwen-1.5B—a model 21x smaller. Ablation studies confirm that THINKLOGIT-DPO succeeds only when it couples a preference-learning objective with training pairs drawn from both the target and guider models. Our work presents a computationally-efficient method to elicit long reasoning in large models with minimal or no additional training.

## 1 Introduction

Large reasoning models (LRMs), such as DeepSeek-R1 (DeepSeek-AI et al., 2025), OpenAI o1 (OpenAI, 2024), and Qwen3 (Qwen Team, 2025), have significantly advanced reasoning by leveraging inference-time compute (Snell et al., 2024; Brown et al., 2024). These models generate very long chain-of-thought (CoT) traces involving planning, reflection, and self-correction (Gandhi

et al., 2025). It is widely believed that such behavior requires training, either through reinforcement learning with verifiable rewards (DeepSeek-AI et al., 2025; Lambert et al., 2024; Shao et al., 2024) or supervised distillation (Muennighoff et al., 2025; Li et al., 2025b). However, this training is costly due to the length of reasoning traces and extensive sampling. While such costs are prohibitive for large models, small models can be trained with modest compute (Dang and Ngo, 2025; Luo et al., 2025). This observation motivates our central research question: *Can a small reasoning model elicit long CoT behavior in a larger model at inference time, without training the larger model?*

We address this question with a decoding-time technique, enabling a small reasoning model to guide a target model, mainly by manipulating its logits. Specifically, we use logit arithmetic (Liu et al., 2021, 2024; Mitchell et al., 2024; Fan et al., 2024a) to combine the output distributions of both models, allowing the target model to benefit from the guider model’s long-chain-of-thought capabilities, without any additional training. We call this base approach **THINKLOGIT**. Furthermore, as the output distribution of both models may substantially differ, we align them by further training the small guider model. This training process uses Direct Preference Optimization (DPO; Rafailov et al., 2023) on preference pairs sampled from the guider and target models, thereby making THINKLOGIT more *on-policy*, and then applies logit arithmetic using the fine-tuned guider. We refer to this approach as **THINKLOGIT-DPO** and show that such training can further boost performance compared to THINKLOGIT.

We evaluate our methods on five challenging benchmarks covering mathematical and scientific reasoning. Our results show that fusing the logits of a small reasoning model (DeepSeek-R1-Distill-Qwen-1.5B) with those of a large target (Qwen2.5-32B) improves pass@1 by 24.5% with



a short-CoT model into a long-CoT one. Intuitively, adding this delta to  $L$  induces analogous long reasoning behavior without altering its weights.

**Warm-up for Stable Decoding.** We empirically observe that directly applying logit arithmetic at each decoding step would cause many repetitive generations in the long-CoT scenario. To stabilize generations, we defer guidance until a prefix of length  $T$ :

$$\tilde{\ell}_{t+1} = \begin{cases} \ell_{t+1}^{(L)}, & t+1 \leq T, \\ \ell_{t+1}^{(L)} + \alpha(\ell_{t+1}^{(S^*)} - \ell_{t+1}^{(S)}), & t+1 > T, \end{cases} \quad (1)$$

We set  $T=100$  tokens in all experiments unless otherwise specified.

## 2.2 THINKLOGIT-DPO

The effectiveness of THINKLOGIT can be limited by mismatches between the output distributions of the guider and target models. We therefore construct preference pairs that capture complementary strengths:

**Type-1:**  $(x, y^{L\checkmark}, y^{S\times})$  — The *large* model’s correct (short) CoT is preferred over the *small* model’s incorrect (long) one. This encourages the guider to preserve the correctness of the target model and avoid introducing new errors.

**Type-2:**  $(x, y^{S\checkmark}, y^{L\times})$  — The *small* model’s correct (long) CoT is preferred over the *large* model’s incorrect (short) one, teaching the guider to be more confident at fixing the large model’s reasoning errors.

We gather these pairs from training queries  $x$  by independently sampling CoTs from  $L$  and  $S^*$  and labeling correctness based on the final answer. Let  $\theta$  denote the parameters of the preference-optimized guider, initialized from  $S^*$ . We train  $\theta$  with a DPO objective function that mixes the two pair types:

$$\mathcal{L}_{\text{DPO}}(\theta) = \lambda \mathbb{E}_{(x, y^{L\checkmark}, y^{S\times}) \sim \mathcal{D}_1} \ell_{\theta}(x; y^{L\checkmark}, y^{S\times}) + (1 - \lambda) \mathbb{E}_{(x, y^{S\checkmark}, y^{L\times}) \sim \mathcal{D}_2} \ell_{\theta}(x; y^{S\checkmark}, y^{L\times}), \quad (2)$$

where  $\ell_{\theta}(x; y^+, y^-) = \log \sigma(r_{\theta}(x, y^+) - r_{\theta}(x, y^-))$ ,  $\sigma$  is the sigmoid function,  $r_{\theta}(x, y) = \beta[\log \pi_{\theta}(y | x) - \log \pi_{\text{ref}}(y | x)]$  is the implicit reward of trajectory  $y$ , and  $\lambda \in [0, 1]$  balances the two datasets  $\mathcal{D}_1$  (Type-1) and  $\mathcal{D}_2$  (Type-2). We use  $\lambda = \frac{|\mathcal{D}_1|}{|\mathcal{D}_1| + |\mathcal{D}_2|}$  by default, directly concatenating

two datasets as DPO training data without further rebalancing. After fine-tuning, we replace  $S^*$  in THINKLOGIT with the optimized guider to obtain THINKLOGIT-DPO.

## 3 Experiments and Results

### 3.1 Experimental Setup

**Benchmarks.** We evaluate models on five widely used reasoning benchmarks. Four of them are competition math problems sources from AIME2024 (30 problems), AIME2025 (30 problems), AMC23 (40 problems), and a subset of 134 hard problems (level 5) from MATH500 (Lightman et al., 2024). We also evaluate on another scientific reasoning dataset GPQA Diamond (Rein et al., 2023), consisting of 198 PhD-level science questions in Biology, Chemistry, and Physics. For each dataset, we independently sample 8 completions and compute pass@k (Chen et al., 2021). A problem is marked as solved if any of the  $k$  sampled outputs is correct, so pass@k helps reveal a model’s potential to solve a problem. We primarily use pass@1 unless otherwise specified.

**Models.** Our major target model is **Qwen2.5-32B** (Yang et al., 2024a). We utilize a long CoT post-trained 1.5B models as the guider which is **R1-Distill-Qwen-1.5B** (DeepSeek-AI et al., 2025), a version based on Qwen2.5-Math-1.5B (Yang et al., 2024b) that has been supervised fine-tuned on 800K long-CoT examples distilled from DeepSeek-R1. Because all three models use the identical tokenizer, their output logits are directly comparable and can be combined arithmetically.

**Preference Data Construction.** We use the level 4–5 subset of the MATH training set (Hendrycks et al., 2021) and independently sample 5 completions from both the guider model ( $S^*$ ) and the target model ( $L$ ). Each completion is checked for final-answer correctness against the gold label.<sup>1</sup>

The target model  $L$  yields 12,412 correct completions ( $y^{L\checkmark}$ ) and 16,448 incorrect ones ( $y^{L\times}$ ), whereas the guider  $S^*$  produces 18,651 correct ( $y^{S\checkmark}$ ) and 10,209 incorrect ( $y^{S\times}$ ) completions. Forming the Cartesian product for each question gives 11,974 Type-1 preference pairs ( $y^{L\checkmark}, y^{S\times}$ ) and 43,209 Type-2 pairs ( $y^{S\checkmark}, y^{L\times}$ ), for a total of

<sup>1</sup>We extract answers from `\boxed{\}` and compute exact match with ground-truths based on this script <https://github.com/openai/prm800k/blob/main/prm800k/grading/grader.py> by (Lightman et al., 2024).

Model	# Training Examples	# Trainable Params	AIME 2024	AIME 2025	AMC 23	MATH Level 5	GPQA Diamond	Average
Qwen2.5-32B ( <i>Target</i> )	-	-	14.6	8.3	<b>57.2</b>	44.7	<b>36.9</b>	32.3
R1-Distill-1.5B ( <i>Guider</i> )	-	-	<b>16.2</b>	<b>18.8</b>	51.2	<b>47.5</b>	28.9	<b>32.5</b>
<i>No Fine-tuning of the Target</i>								
Target + THINKLOGIT	0	0	<b>22.5</b>	19.2	62.2	55.3	41.8	40.2
Target + THINKLOGIT-DPO	10K	78M	22.1	<b>21.7</b>	<b>63.7</b>	<b>58.5</b>	<b>42.4</b>	<b>41.7</b>
<i>Full Fine-tuning of the Target</i>								
s1.1-32B	1K	32B	32.9	25.4	70.0	72.2	51.9	44.5
R1-Distill-32B	800K	32B	<b>45.8</b>	<b>35.0</b>	<b>76.9</b>	<b>72.7</b>	<b>55.6</b>	<b>57.2</b>

Table 1: Comparison of **pass@1** performance across five reasoning benchmarks. The best results in each section are marked in **bold**. Key takeaways include: (1) fusing target and guider logits (THINKLOGIT) yields substantial accuracy gains on top of both models; (2) DPO-trained guider (THINKLOGIT-DPO) adds further improvement; (3) math-only guidance alignment transfers effectively to out-of-domain scientific reasoning (GPQA Diamond); (4) THINKLOGIT-DPO partially recovers benefits of full fine-tuning with fewer trainable parameters and less training data.

55,183 pairs. We then randomly select 10K preference pairs from the total 55K pairs for DPO fine-tuning. We applied LoRA (Hu et al., 2022) with a rank size of 64 for parameter-efficient fine-tuning of the guider model. For all models, decoding is performed with a temperature of 0.6, a maximum length of 8192 tokens, and guidance strength of  $\alpha = 1$ . More training details are in Appendix A.1.

### 3.2 Main Results

Table 1 presents the pass@1 scores for all systems. We highlight three key observations. *First*, THINKLOGIT boosts reasoning accuracy upon both target and guider model, and THINKLOGIT-DPO raises it further. Combining the logits of the 32B target with those of the 1.5B guider (THINKLOGIT) raises the average pass@1 by 24.5% relative to the frozen target and by 23.7% relative to the guider. Replacing the vanilla guider with the DPO-trained guider (THINKLOGIT-DPO) brings the relative improvement to 29.1% over the target model, without any extra inference cost. The performance gains are consistent across five tested datasets.

*Second*, a guider trained only on mathematics problems maintains effectiveness on out-of-domain scientific reasoning. Although the DPO alignment phase trains the 1.5B guider solely to mathematics problems, it still maintains and even slightly improves performance of vanilla logit arithmetic (THINKLOGIT) from 41.8 to 42.4 on the out-of-domain GPQA Diamond benchmark, which spans biology, chemistry, and physics. This

indicates that while DPO shrinks the distribution gap between guider and target outputs on maths data, it also makes the guider’s token-level probabilities easier for the target to follow regardless of subject matter. A recent study by Tang et al. (2025) also reports a similar domain-general nature of long CoT based on model representation analysis, reinforcing the potential for enhanced guidance from a math-trained guider to transfer broadly across disciplines.

*Third*, our approach recovers most of the benefit of full-parameter fine-tuning while touching only a small subset of weights and using far less data. With LoRA, we adjust just 78M adapter parameters and train on 10K preference pairs, yet THINKLOGIT-DPO closes 77% of the pass@1 gap between the frozen 32B target and the fully fine-tuned s1.1-32B (Muennighoff et al., 2025), which updates all 32B parameters using 1K carefully selected long-CoT examples from a corpus of 59K examples. Our pipeline does rely on a 1.5B guider (R1-distill-1.5B) that was already fine-tuned on 800K distilled examples from DeepSeek-R1, but training this smaller model is far cheaper than R1-distill-32B, and once trained, the same guider can be reused for many other larger targets, especially those in the same model family (e.g., Qwen2.5-72B, with results shown in Figure 4), at no extra cost. Consequently, the cumulative data and compute requirements of our pipeline remain well below those of fully fine-tuning large models, while still delivering substantial accuracy gains.



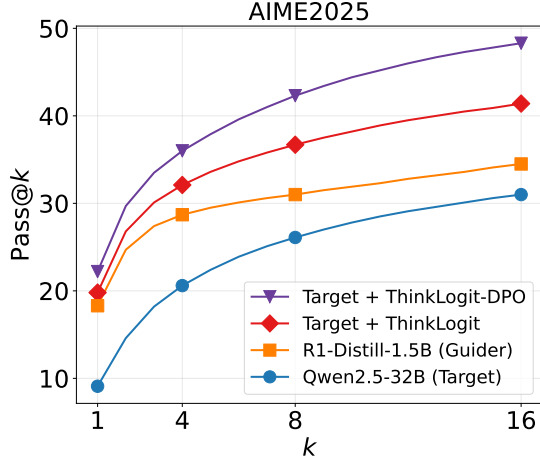


Figure 2: Test-time scaling on AIME2025.  $\text{Pass}@k$  for  $k = 1-16$  comparing the target, guider, their direct logit fusion (THINKLOGIT), and the DPO-aligned fusion (THINKLOGIT-DPO). Our methods not only increase sample efficiency but also broaden the reasoning boundary of the target model.

### 3.3 Test-Time Scaling Properties

Figure 2 plots  $\text{pass}@k$  for  $k = 1-16$  on AIME2025, the dataset where the 32B target performs worst and scaling effects are therefore most visible. Both THINKLOGIT and THINKLOGIT-DPO surpass the target’s  $\text{pass}@16$  performance with only four generations, achieving a four-fold improvement in *sample efficiency*. The advantage widens as  $k$  grows: at  $k = 16$  our DPO-aligned guider leads the target by roughly 17 points. Unlike the baseline’s early plateau, our curve keeps rising, implying that logits guidance *broaden the reasoning boundary* (Yue et al., 2025b) rather than merely re-ranking similar completions.

### 3.4 Comparison with Other Training-Free Methods for Long CoT Elicitation

Figure 3 contrasts our approach against two training-free baselines for long chain-of-thought elicitation. First, the budget-forcing heuristic introduced by Muennighoff et al. (2025) replaces end-of-sentence tokens with a placeholder string like “Wait” to artificially increase output length. While this does produce longer completions, it consistently hurts performance, showing that verbosity alone does not lead to deeper reasoning.<sup>2</sup> Second,

<sup>2</sup>We note that while Muennighoff et al. (2025) demonstrate the effectiveness of budget-forcing on a Qwen2.5-32B-Instruct model fine-tuned on 1K long CoTs, they do not evaluate this technique directly on the untuned model.



Figure 3: Comparison of our training-free long chain-of-thought elicitation method (THINKLOGIT) against two baselines: budget-forcing and one-shot long CoT in-context learning (ICL). The left panel shows  $\text{pass}@1$  on AIME2025 and AMC23; the right panel shows the average chain-of-thought length in tokens. While budget-forcing and long CoT ICL increase verbosity, they degrade accuracy, whereas THINKLOGIT produces genuinely extended reasoning that boosts performance.

inserting a single long CoT example in the prompt (sampled from the s1.1-1K dataset (Muennighoff et al., 2025)) for in-context learning (ICL; Brown et al., 2020; Min et al., 2022; Dong et al., 2024) also degrades performance despite of longer outputs from the target model. In contrast, THINKLOGIT-DPO uses logit-level guidance from a small reasoning model to steer the decoding towards genuine long chain-of-thoughts, which translates into a clear uplift in downstream accuracy. This shows that our improvements stem from *the quality of the guidance* being applied, rather than *the quantity of tokens* generated.

### 3.5 Cross-Model Transferability of THINKLOGIT-DPO

We evaluate whether the THINKLOGIT-DPO guider, optimized for the reasoning preference of Qwen2.5-32B, can be applied off-the-shelf to a larger model (Qwen2.5-72B) in the same family. At inference time, we fuse the R1-distill-1.5B guider logits with the 72B target (THINKLOGIT) and then swap in the DPO-trained guider (THINKLOGIT-DPO), without any additional fine-tuning on the outputs from the 72B model. As highlighted in Figure 4, THINKLOGIT-DPO consistently improves upon THINKLOGIT, confirming that preference signals learned via DPO on a 32B model transfer effectively to larger scales, offering a plug-and-play mechanism to boost long-chain reasoning in even more capable LLMs within the same family.

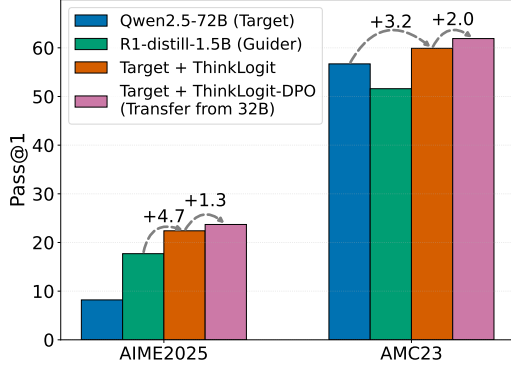


Figure 4: Pass@1 for eliciting long CoT in a 72B target model with logits arithmetic. THINKLOGIT-DPO delivers larger performance improvements on AIME2025 and AMC23 compared to THINKLOGIT, demonstrating that preference signals learned on a 32B model transfer effectively to a larger 72B model.

### 3.6 Ablation Study of THINKLOGIT-DPO

To further investigate the design choices in THINKLOGIT-DPO, we ablate both our mixed-pair data construction and preference-based learning objective (DPO) against single-source or supervised fine-tuning alternatives. Results in Table 2 answer the following research questions.

**Are preference pairs sourced from both the target and the guider necessary to maximize performance?** We construct the same amount of 10K preference pairs using only the guider’s correct vs. incorrect outputs, i.e.,  $(x, y^{S\checkmark}, y^{S\times})$ . DPO on this data underperforms markedly on AMC23 (58.8 vs. 63.7), confirming that mixing pairs which highlight *both* the target’s and guider’s strengths is crucial for maximal gains.

**Is training on both types of pairs necessary for the effectiveness of THINKLOGIT-DPO?** We next ablate by training on only one type of preference pairs at a time: using only Type-2 pairs  $(x, y^{S\checkmark}, y^{L\times})$  (i.e.,  $\lambda = 0$  in Equation 2) yields a pass@1 of 57.2, while using only Type-1 pairs  $(x, y^{L\checkmark}, y^{S\times})$  (i.e.,  $\lambda = 1$  in Equation 2) drops further to 51.9. Both are substantially below the 63.7 achieved by the full mixture, indicating that *both* Type-2 pairs (which teach the guider to correct target errors) and Type-1 pairs (which enforce preservation of correct target outputs) provide complementary signals necessary for optimal alignment.

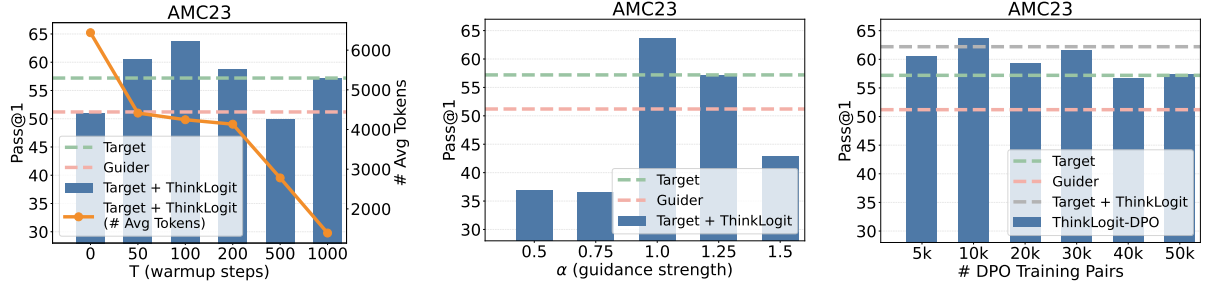
**Can supervised fine-tuning replace preference-based alignment of the guider?** We evaluate standard supervised fine-tuning (SFT) against

Model	Guider’s Training Data	Pass@1
THINKLOGIT-DPO (ours)	$(x, y^{L\checkmark}, y^{S\times}), (x, y^{S\checkmark}, y^{L\times})$	<b>63.7</b>
THINKLOGIT-DPO		
w/o dual sources	$(x, y^{S\checkmark}, y^{S\times})$	58.8
w/o Type-1 pairs	$(x, y^{S\checkmark}, y^{L\times})$	57.2
w/o Type-2 pairs	$(x, y^{L\checkmark}, y^{S\times})$	51.9
THINKLOGIT-SFT		
learning from target	$(x, y^{L\checkmark})$	44.7
self-learning	$(x, y^{S\checkmark})$	55.6
learning from teacher	$(x, y^{R1\checkmark})$	60.9

Table 2: Pass@1 on AMC23 under ablations of guider’s training data and objectives in THINKLOGIT-DPO. We compare the full DPO regime with mixed Type-1 and Type-2 pairs against single-source DPO (only guider outputs, only Type-1, only Type-2) and supervised fine-tuning variants. The dual-source, mixed-pair DPO yields the highest accuracy, demonstrating the necessity of complementary preference signals and preference-based alignment.

DPO by training the guider on three equally sized sets of high-quality completions: (1) the target model’s correct outputs  $y^{L\checkmark}$ , (2) its own correct outputs  $y^{S\checkmark}$  (also known as rejection-sampling fine-tuning (Yuan et al., 2023)), and (3) R1-distilled completions  $y^{R1\checkmark}$ . Although SFT on (1) and (2) turns the guider into a better *standalone* reasoner, none of these variants rival the performance of the DPO-aligned guider. This gap demonstrates that *optimizing with pairwise preference comparisons yields a better guider than optimizing solely for correctness*. While SFT can adapt the guider toward the target’s short-CoT reasoning style in (3) and thus reduce the distributional gap, it tends to overwrite the guider’s native strengths. In contrast, DPO *maintains the guider’s intrinsic reasoning capabilities*—preserving its long reasoning capability—while *selectively aligning it to the target’s preferences* through pairwise comparisons.

Overall, the best performance arises when the guider is aligned with the target via DPO and trained on a mixture of Type-1 and Type-2 preference pairs sourced from both models. Naive SFT—whether on the guider itself, the target’s outputs, or a stronger teacher—fails to match these gains, underscoring key factors behind the effectiveness of THINKLOGIT-DPO.



(a) Varying the warm-up steps  $T$  reveals that applying guidance too early (i.e.,  $T=0$ ) leads to repeated generations and lower accuracy, while a moderate  $T$  improves pass@1 and stabilize decoding.

(b) Sweeping the guidance strength  $\alpha$  shows that  $\alpha = 1.0$  yields the best trade-off between guider influence and target model priors.

(c) Increasing DPO training pairs from 5K to 10K improves pass@1, but further scaling leads to negative returns, likely due to redundancy and overfitting in preference data.

Figure 5: Impact of hyperparameter tuning on THINKLOGIT (Figure 5a and 5b) and training data size on THINKLOGIT-DPO (Figure 5c).

### 3.7 Impact of Hyperparameter Tuning on THINKLOGIT

In THINKLOGIT, two hyperparameters play a critical role in balancing stability, accuracy, and generation efficiency: the warm-up length  $T$  and the guidance strength  $\alpha$  (Eq. 1). We evaluate their effects on the AMC23 benchmark, which presents a suitable mix of problem difficulties and clearly exhibits both stability and guidance effects.

To assess warm-up, we vary  $T$  over  $\{0, 50, 100, 200, 500, 1000\}$  with  $\alpha$  fixed to 1. For each possible values of  $T$ , we sample 8 completions per question and compute pass@1 alongside the average number of generated tokens (Figure 5a). When  $T = 0$ , guidance is applied immediately from the beginning, causing repetitive loops in early decoding and yielding the lowest accuracy. Allowing 50–200 tokens of unguided generation stabilizes the chain-of-thought, improving pass@1 over both target and guider models and reducing generation length. Increasing  $T$  beyond 200 causes the model to revert to the shorter CoTs typically produced by the target model, leading to an accuracy drop and shorter outputs.

With  $T$  fixed at 100, we sweep  $\alpha$  over  $\{0.5, 0.75, 1.0, 1.25, 1.5\}$  to control how strongly the guider’s delta-logits modify the target’s distribution (Figure 5b). At  $\alpha = 1.0$ , we observe the highest pass@1 together with moderate generation length, indicating an optimal trade-off between the guider’s corrective signal and the target model’s own priors. Therefore, we set  $\alpha = 1.0$  as the default guidance strength. Future work might explore adaptive, context-aware schedules for  $T$  and  $\alpha$  (Fan et al., 2024a) to optimize this trade-off further.

### 3.8 Impact of Training Data Size on THINKLOGIT-DPO

To determine the optimal number of DPO preference pairs, we randomly sampled subsets of 5K, 10K, 20K, 30K, 40K, and 50K from our full pool of 55K Type-1 and Type-2 pairs. Figure 5c plots pass@1 against training dataset size. With 5K pairs, THINKLOGIT-DPO’s pass@1 remains lower than that of the vanilla THINKLOGIT; increasing to 10K pairs raises pass@1 above THINKLOGIT while keeping training cost moderate. Beyond 10K pairs, adding more data leads to a decline in pass@1. Since our preference set is constructed via the Cartesian product of correct and incorrect generations, we hypothesize that pairs beyond 10K primarily recombine existing chains-of-thought rather than introduce new solution patterns, resulting in redundant examples and a higher risk of overfitting.

## 4 Related Work

### 4.1 Long Chain-of-Thought (CoT) Reasoning

Large reasoning models, such as OpenAI’s o1 and o3 (OpenAI, 2024, 2025), DeepSeek-R1 (DeepSeek-AI et al., 2025), and QwQ (Team, 2025), achieve state-of-the-art results on mathematical and coding benchmarks by generating CoT traces that often extend to thousands of tokens, enabling systematic backtracking, verification, and self-reflection before a final answer is produced (Gandhi et al., 2025). One way to elicit such long-form reasoning is through **reinforcement learning with verifiable rewards** (Lambert et al., 2024). Pioneered by Group-Relative Policy Optimization (GRPO) (Shao et al., 2024) and refined by more stable and token-efficient variants such as DAPO (Yu

et al., 2025) and Dr. GRPO (Liu et al., 2025), this approach optimizes outcome-based rewards for correctness; nevertheless, mounting evidence shows that it mainly re-weights reasoning patterns already latent in the base model (Liu et al., 2025; Yue et al., 2025a). A complementary line of work demonstrates that the same capability can be acquired with **data-efficient supervised fine-tuning**. Distilled long CoTs from stronger teacher models allows a student to extend its reasoning length and thus improve accuracy using only about one thousand examples (Muennighoff et al., 2025; Xu et al., 2025; Ye et al., 2025; Li et al., 2025b). Finally, **training-free** methods exploit the fact that pretrained LLMs already exhibit long-CoT behaviours (Liu et al., 2025; Gandhi et al., 2025). (Tang et al., 2025) inject contrastive long- versus short-CoT representations into hidden states via representation engineering (Zou et al., 2023), whereas (Zhao et al., 2025) amplify a handful of key neurons at inference. Both techniques, however, require domain-specific long/short traces and white-box access, limiting their applicability in out-of-domain or black-box settings. THINKLOGIT sidesteps these constraints entirely. It keeps the target LLM frozen and, at inference time, fuses its logits with those of a lightweight “guider” model trained for long reasoning. This logit-fusion strategy recovers long-CoT behaviour induced by training-based methods while introducing no additional training cost or curated long-CoT examples.

## 4.2 Decoding Algorithms for LLM Reasoning

Decoding-time interventions offer an attractive alternative to full model fine-tuning: they can improve the reasoning capabilities of an off-the-shelf LLM with only a marginal increase in training or inference cost. The earliest line of work is Chain-of-Thought (CoT) prompting (Wei et al., 2022), which simply asks the model to “think aloud.” Subsequent self-consistency decoding (Wang et al., 2023) samples a set of diverse CoTs and majority-votes over their answers, while later work shows that short CoTs can even be elicited *without* any prompting (Wang and Zhou, 2024). Crucially, these traces are usually brief: they march directly to the answer without back-tracking or verification, and therefore do not unlock the *long-form* reasoning studied in our work. A second family, **guided decoding**, biases generation toward correctness using either self-evaluation signals from the model itself (Xie et al., 2023) or an external discriminator (Khalifa

et al., 2023). Accuracy is further improved by best-of- $n$  reranking with discriminative reward models that score either the final answer or the reasoning process (Cobbe et al., 2021; Lightman et al., 2024; Wang et al., 2024). Generative reward objectives extend this idea and generalise better across tasks (Hosseini et al., 2024; Zhang et al., 2025; Wang et al., 2025; Khalifa et al., 2025). However, all of these methods depend on sampling many complete reasoning traces and scoring them *after* they are generated, which both raises costs and keeps them in the short-CoT regime. Auxiliary-model approaches modify the output of a frozen *target* model on the fly. Contrastive decoding subtracts logits from an “amateur” model or layer to suppress low-quality outputs (Li et al., 2023; Chuang et al., 2024), while speculative decoding speeds inference by letting a small draft model propose tokens that the expert later accepts or rejects (Leviathan et al., 2023; Yang et al., 2025; Liao et al., 2025). A closely related strand, **logits arithmetic**, blends the output distributions of three models token-by-token (Liu et al., 2021; Ormazabal et al., 2023; Shi et al., 2024), successfully emulating task-specific fine-tuning (Liu et al., 2024; Fan et al., 2024b), scaling laws (Mitchell et al., 2024), unlearning (Huang et al., 2025) and even overriding safety filters (Zhao et al., 2024). THINKLOGIT follows this lightweight pathway by using a compact guider model to unlock long-form reasoning in a frozen large model, while THINKLOGIT-DPO additionally aligns the guider’s distribution with the target model’s, delivering further gains.

## 5 Conclusion and Future Work

We introduce THINKLOGIT and THINKLOGIT-DPO, two decoding-time techniques that unlock long chain-of-thought (CoT) reasoning in frozen, non-reasoning LLMs. THINKLOGIT injects logits from a small, long-CoT guider, boosting accuracy by 24.5 % on five reasoning benchmarks for only a 1.1 $\times$  increase in inference-time parameters, while THINKLOGIT-DPO aligns the guider with target distribution via Direct Preference Optimization for even higher gains. Together they offer a compute-efficient route to deploy long-CoT LLMs. Future work will combine heterogeneous model families, and develop context-aware guidance (e.g., adaptive strength  $\alpha$  as in Fan et al. (2024a)) to mitigate the over-thinking problem in long reasoning (Chen et al., 2024).



## Limitations

**Inference-Time Overhead.** Deploying THINKLOGIT requires hosting the large target model along with two smaller models—the base model  $S$  and the DPO-aligned guider  $S^*$ . In our primary experimental setup (guiding a 32B target with a 1.5B guider), the total parameter count increases by approximately  $1.1\times$  compared to using the frozen target alone. Instead of sequentially querying each model at every inference step, we implement asynchronous decoding to concurrently obtain logits from all three models. Profiling on NVIDIA L40S GPUs indicates a moderate inference slowdown (approximately 25% fewer tokens per second) compared to running only the target model. Since THINKLOGIT-DPO simply replaces the original guider model used in THINKLOGIT with a preference-optimized model of the same size, THINKLOGIT-DPO incurs no additional inference overhead beyond THINKLOGIT itself.

**Same-Family Constraint.** THINKLOGIT and THINKLOGIT-DPO computes *token-wise* differences between the guider’s and a base model’s logits and then adds that delta to the target model. Because the three models must share an identical vocabulary, we currently restrict all three to the same model family (e.g., Qwen2.5). Although Section 3.5 shows that a guider aligned on a 32B target transfers to a 72B target in the *same* family, we have not yet verified that the method generalizes to other model families. Extending the approach to heterogeneous families such as Mistral (Ras-togi et al., 2025), Llama (Dubey et al., 2024), or Gemma (Kamath et al., 2025) will require a robust tokenizer alignment algorithm (Fu et al., 2023; Li et al., 2025a) to ensure delta logits remain semantically meaningful across models. We leave the design and empirical validation of such cross-family fusion to future work.

**Limited Domains of Evaluation.** Our experiments focus on math- and science-oriented reasoning tasks. A broader evaluation suite, including coding (Jimenez et al., 2023; Jain et al., 2025), planning (Zheng et al., 2024a; Xie et al., 2024), and tool-use (Huang et al., 2024; Patil et al., 2025), is needed to understand failure modes that may emerge in less structurally similar settings.

**Offline alignment.** The guider is aligned with the target via Direct Preference Optimisation (DPO)

on a fixed set of preference pairs. This *offline* formulation cannot adapt once deployment uncovers new error patterns or distribution drift. Incorporating *online* reinforcement learning (Schulman et al., 2017; Shao et al., 2024) that updates the guider from streamed on-policy samples could, in principle, reduce this brittleness. However, on-policy RL introduces training efficiency and stability challenges that remain open research problems.

## References

- Bradley C. A. Brown, Jordan Juravsky, Ryan Ehrlich, Ronald Clark, Quoc V. Le, Christopher Ré, and Azalia Mirhoseini. 2024. [Large language monkeys: Scaling inference compute with repeated sampling](#). *CoRR*, abs/2407.21787.
- Tom B. Brown, Benjamin Mann, Nick Ryder, Melanie Subbiah, Jared Kaplan, Prafulla Dhariwal, Arvind Neelakantan, Pranav Shyam, Girish Sastry, Amanda Askell, Sandhini Agarwal, Ariel Herbert-Voss, Gretchen Krueger, Tom Henighan, Rewon Child, Aditya Ramesh, Daniel M. Ziegler, Jeffrey Wu, Clemens Winter, and 12 others. 2020. [Language models are few-shot learners](#). In *Advances in Neural Information Processing Systems 33: Annual Conference on Neural Information Processing Systems 2020, NeurIPS 2020, December 6-12, 2020, virtual*.
- Mark Chen, Jerry Tworek, Heewoo Jun, Qiming Yuan, Henrique Pondé de Oliveira Pinto, Jared Kaplan, Harri Edwards, Yuri Burda, Nicholas Joseph, Greg Brockman, Alex Ray, Raul Puri, Gretchen Krueger, Michael Petrov, Heidy Khlaaf, Girish Sastry, Pamela Mishkin, Brooke Chan, Scott Gray, and 39 others. 2021. [Evaluating large language models trained on code](#). *CoRR*, abs/2107.03374.
- Xingyu Chen, Jiahao Xu, Tian Liang, Zhiwei He, Jianhui Pang, Dian Yu, Linfeng Song, Qiuzhi Liu, Mengfei Zhou, Zhuosheng Zhang, Rui Wang, Zhaopeng Tu, Haitao Mi, and Dong Yu. 2024. [Do NOT think that much for 2+3=? on the overthinking of o1-like llms](#). *CoRR*, abs/2412.21187.
- Yung-Sung Chuang, Yujia Xie, Hongyin Luo, Yoon Kim, James R. Glass, and Pengcheng He. 2024. [Dola: Decoding by contrasting layers improves factuality in large language models](#). In *The Twelfth International Conference on Learning Representations, ICLR 2024, Vienna, Austria, May 7-11, 2024*. OpenReview.net.
- Karl Cobbe, Vineet Kosaraju, Mohammad Bavarian, Mark Chen, Heewoo Jun, Lukasz Kaiser, Matthias Plappert, Jerry Tworek, Jacob Hilton, Reiichiro Nakano, Christopher Hesse, and John Schulman. 2021. [Training verifiers to solve math word problems](#). *CoRR*, abs/2110.14168.
- Quy-Anh Dang and Chris Ngo. 2025. Reinforcement learning for reasoning in small llms: What works and what doesn’t. *arXiv preprint arXiv:2503.16219*.

683	DeepSeek-AI, Daya Guo, Dejian Yang, Haowei Zhang,	Arian Hosseini, Xingdi Yuan, Nikolay Malkin, Aaron C.	740
684	Junxiao Song, Ruoyu Zhang, Runxin Xu, Qihao Zhu,	Courville, Alessandro Sordoni, and Rishabh Agar-	741
685	Shirong Ma, Peiyi Wang, Xiao Bi, Xiaokang Zhang,	wal. 2024. <a href="#">V-star: Training verifiers for self-taught</a>	742
686	Xingkai Yu, Yu Wu, Z. F. Wu, Zhibin Gou, Zhi-	<a href="#">reasoners</a> . <i>CoRR</i> , abs/2402.06457.	743
687	hong Shao, Zhuoshu Li, Ziyi Gao, and 81 others.		
688	2025. <a href="#">Deepseek-r1: Incentivizing reasoning capa-</a>	Edward J. Hu, Yelong Shen, Phillip Wallis, Zeyuan	744
689	<a href="#">bility in llms via reinforcement learning</a> . <i>CoRR</i> ,	Allen-Zhu, Yuanzhi Li, Shean Wang, Lu Wang, and	745
690	abs/2501.12948.	Weizhu Chen. 2022. <a href="#">Lora: Low-rank adaptation of</a>	746
		<a href="#">large language models</a> . In <i>The Tenth International</i>	747
691	Qingxiu Dong, Lei Li, Damai Dai, Ce Zheng, Jingyuan	<i>Conference on Learning Representations, ICLR 2022,</i>	748
692	Ma, Rui Li, Heming Xia, Jingjing Xu, Zhiyong Wu,	<i>Virtual Event, April 25-29, 2022</i> . OpenReview.net.	749
693	Baobao Chang, Xu Sun, Lei Li, and Zhifang Sui.		
694	2024. <a href="#">A survey on in-context learning</a> . In <i>Proceed-</i>	James Y. Huang, Wenxuan Zhou, Fei Wang, Fred	750
695	<i>ings of the 2024 Conference on Empirical Methods in</i>	Morstatter, Sheng Zhang, Hoifung Poon, and Muhao	751
696	<i>Natural Language Processing, EMNLP 2024, Miami,</i>	Chen. 2025. <a href="#">Offset unlearning for large language</a>	752
697	<i>FL, USA, November 12-16, 2024</i> , pages 1107–1128.	<a href="#">models</a> . <i>Trans. Mach. Learn. Res.</i> , 2025.	753
698	Association for Computational Linguistics.		
		Shijue Huang, Wanjun Zhong, Jianqiao Lu, Qi Zhu, Ji-	754
699	Abhimanyu Dubey, Abhinav Jauhri, Abhinav Pandey,	ahui Gao, Weiwen Liu, Yutai Hou, Kingshan Zeng,	755
700	Abhishek Kadian, Ahmad Al-Dahle, Aiesha Letman,	Yasheng Wang, Lifeng Shang, Xin Jiang, Ruifeng	756
701	Akhil Mathur, Alan Schelten, Amy Yang, Angela	Xu, and Qun Liu. 2024. <a href="#">Planning, creation, usage:</a>	757
702	Fan, Anirudh Goyal, Anthony Hartshorn, Aobo Yang,	<a href="#">Benchmarking llms for comprehensive tool utiliza-</a>	758
703	Archi Mitra, Archie Sravankumar, Artem Korenev,	<a href="#">tion in real-world complex scenarios</a> . In <i>Findings of</i>	759
704	Arthur Hinsvark, Arun Rao, Aston Zhang, and 82	<i>the Association for Computational Linguistics, ACL</i>	760
705	others. 2024. <a href="#">The llama 3 herd of models</a> . <i>CoRR</i> ,	<i>2024, Bangkok, Thailand and virtual meeting, Au-</i>	761
706	abs/2407.21783.	<i>gust 11-16, 2024</i> , pages 4363–4400. Association for	762
		Computational Linguistics.	763
707	Chenghao Fan, Zhenyi Lu, Wei Wei, Jie Tian, Xiaoye		
708	Qu, Danyang Chen, and Yu Cheng. 2024a. On gi-	Naman Jain, King Han, Alex Gu, Wen-Ding Li, Fanjia	764
709	ant’s shoulders: Effortless weak to strong by dynamic	Yan, Tianjun Zhang, Sida Wang, Armando Solar-	765
710	logits fusion. In <i>The Thirty-eighth Annual Confer-</i>	Lezama, Koushik Sen, and Ion Stoica. 2025. <a href="#">Live-</a>	766
711	<i>ence on Neural Information Processing Systems</i> .	<a href="#">codebench: Holistic and contamination free evalua-</a>	767
		<a href="#">tion of large language models for code</a> . In <i>The Thir-</i>	768
712	Chenghao Fan, Zhenyi Lu, Wei Wei, Jie Tian, Xiaoye	<i>teenth International Conference on Learning Repre-</i>	769
713	Qu, Danyang Chen, and Yu Cheng. 2024b. <a href="#">On gi-</a>	<i>sentsations, ICLR 2025, Singapore, April 24-28, 2025</i> .	770
714	<a href="#">ant’s shoulders: Effortless weak to strong by dynamic</a>	OpenReview.net.	771
715	<a href="#">logits fusion</a> . In <i>Advances in Neural Information Pro-</i>		
716	<i>cessing Systems 38: Annual Conference on Neural</i>	Carlos E. Jimenez, John Yang, Alexander Wettig,	772
717	<i>Information Processing Systems 2024, NeurIPS 2024,</i>	Shunyu Yao, Kexin Pei, Ofir Press, and Karthik	773
718	<i>Vancouver, BC, Canada, December 10 - 15, 2024</i> .	Narasimhan. 2023. <a href="#">Swe-bench: Can language</a>	774
		<a href="#">models resolve real-world github issues?</a> <i>CoRR</i> ,	775
719	Yao Fu, Hao Peng, Litu Ou, Ashish Sabharwal, and	abs/2310.06770.	776
720	Tushar Khot. 2023. <a href="#">Specializing smaller language</a>		
721	<a href="#">models towards multi-step reasoning</a> . In <i>Interna-</i>	Aishwarya Kamath, Johan Ferret, Shreya Pathak, Nino	777
722	<i>tional Conference on Machine Learning, ICML 2023,</i>	Vieillard, Ramona Merhej, Sarah Perrin, Tatiana	778
723	<i>23-29 July 2023, Honolulu, Hawaii, USA</i> , volume	Matejovicova, Alexandre Ramé, Morgane Rivière,	779
724	202 of <i>Proceedings of Machine Learning Research</i> ,	Louis Rouillard, Thomas Mesnard, Geoffrey Cideron,	780
725	pages 10421–10430. PMLR.	Jean-Bastien Grill, Sabela Ramos, Edouard Yvinec,	781
		Michelle Casbon, Etienne Pot, Ivo Penchev, Gaël	782
726	Kanishk Gandhi, Ayush Chakravarthy, Anikait Singh,	Liu, and 79 others. 2025. <a href="#">Gemma 3 technical report</a> .	783
727	Nathan Lile, and Noah D. Goodman. 2025. <a href="#">Cog-</a>	<i>CoRR</i> , abs/2503.19786.	784
728	<a href="#">nitive behaviors that enable self-improving reason-</a>		
729	<a href="#">ers, or, four habits of highly effective stars</a> . <i>CoRR</i> ,	Muhammad Khalifa, Rishabh Agarwal, Lajanugen Lo-	785
730	abs/2503.01307.	geswaran, Jaekyeom Kim, Hao Peng, Moontae Lee,	786
		Honglak Lee, and Lu Wang. 2025. <a href="#">Process reward</a>	787
731	Shousheng Jia Haosheng Zou, Xiaowei Lv and Xi-	<a href="#">models that think</a> . <i>CoRR</i> , abs/2504.16828.	788
732	angzheng Zhang. 2024. <a href="#">360-llama-factory</a> .		
		Muhammad Khalifa, Lajanugen Logeswaran, Moontae	789
733	Dan Hendrycks, Collin Burns, Saurav Kadavath, Akul	Lee, Honglak Lee, and Lu Wang. 2023. <a href="#">GRACE:</a>	790
734	Arora, Steven Basart, Eric Tang, Dawn Song, and Ja-	<a href="#">discriminator-guided chain-of-thought reasoning</a> . In	791
735	cob Steinhardt. 2021. <a href="#">Measuring mathematical prob-</a>	<i>Findings of the Association for Computational Lin-</i>	792
736	<a href="#">lem solving with the MATH dataset</a> . In <i>Proceedings</i>	<i>guistics: EMNLP 2023, Singapore, December 6-10,</i>	793
737	<i>of the Neural Information Processing Systems Track</i>	<i>2023</i> , pages 15299–15328. Association for Computa-	794
738	<i>on Datasets and Benchmarks 1, NeurIPS Datasets</i>	<i>tional Linguistics</i> .	795
739	<i>and Benchmarks 2021, December 2021, virtual</i> .		

Nathan Lambert, Jacob Morrison, Valentina Pyatkin, Shengyi Huang, Hamish Ivison, Faeze Brahman, Lester James V. Miranda, Alisa Liu, Nouha Dziri, Shane Lyu, Yuling Gu, Saumya Malik, Victoria Graf, Jena D. Hwang, Jiangjiang Yang, Ronan Le Bras, Oyvind Tafjord, Chris Wilhelm, Luca Soldaini, and 4 others. 2024. <a href="#">Tülu 3: Pushing frontiers in open language model post-training</a> . <i>CoRR</i> , abs/2411.15124.	853
Yaniv Leviathan, Matan Kalman, and Yossi Matias. 2023. <a href="#">Fast inference from transformers via speculative decoding</a> . In <i>International Conference on Machine Learning, ICML 2023, 23-29 July 2023, Honolulu, Hawaii, USA</i> , volume 202 of <i>Proceedings of Machine Learning Research</i> , pages 19274–19286. PMLR.	854
Chong Li, Jiajun Zhang, and Chengqing Zong. 2025a. <a href="#">Tokalign: Efficient vocabulary adaptation via token alignment</a> . In <i>Proceedings of the 63rd Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers), ACL 2025, Vienna, Austria, July 27 - August 1, 2025</i> , pages 4109–4126. Association for Computational Linguistics.	855
Dacheng Li, Shiyi Cao, Tyler Griggs, Shu Liu, Xiangxi Mo, Eric Tang, Sumanth Hegde, Kourosh Hakhamaneshi, Shishir G. Patil, Matei Zaharia, Joseph E. Gonzalez, and Ion Stoica. 2025b. <a href="#">LLMs can easily learn to reason from demonstrations structure, not content, is what matters!</a> <i>CoRR</i> , abs/2502.07374.	856
Xiang Lisa Li, Ari Holtzman, Daniel Fried, Percy Liang, Jason Eisner, Tatsunori Hashimoto, Luke Zettlemoyer, and Mike Lewis. 2023. <a href="#">Contrastive decoding: Open-ended text generation as optimization</a> . In <i>Proceedings of the 61st Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers), ACL 2023, Toronto, Canada, July 9-14, 2023</i> , pages 12286–12312. Association for Computational Linguistics.	857
Baohao Liao, Yuhui Xu, Hanze Dong, Junnan Li, Christof Monz, Silvio Savarese, Doyen Sahoo, and Caiming Xiong. 2025. <a href="#">Reward-guided speculative decoding for efficient LLM reasoning</a> . <i>CoRR</i> , abs/2501.19324.	858
Hunter Lightman, Vineet Kosaraju, Yuri Burda, Harrison Edwards, Bowen Baker, Teddy Lee, Jan Leike, John Schulman, Ilya Sutskever, and Karl Cobbe. 2024. <a href="#">Let’s verify step by step</a> . In <i>The Twelfth International Conference on Learning Representations, ICLR 2024, Vienna, Austria, May 7-11, 2024</i> . OpenReview.net.	859
Alisa Liu, Xiaochuang Han, Yizhong Wang, Yulia Tsvetkov, Yejin Choi, and Noah A. Smith. 2024. <a href="#">Tuning language models by proxy</a> . volume abs/2401.08565.	860
Alisa Liu, Maarten Sap, Ximing Lu, Swabha Swayamdipta, Chandra Bhagavatula, Noah A. Smith, and Yejin Choi. 2021. <a href="#">Dexperts: Decoding-time controlled text generation with experts and anti-experts</a> . In <i>Proceedings of the 59th Annual Meeting of the Association for Computational Linguistics and the 11th International Joint Conference on Natural Language Processing, ACL/IJCNLP 2021, (Volume 1: Long Papers), Virtual Event, August 1-6, 2021</i> , pages 6691–6706. Association for Computational Linguistics.	861
Zichen Liu, Changyu Chen, Wenjun Li, Penghui Qi, Tianyu Pang, Chao Du, Wee Sun Lee, and Min Lin. 2025. Understanding r1-zero-like training: A critical perspective. <i>arXiv preprint arXiv:2503.20783</i> .	862
Michael Luo, Sijun Tan, Justin Wong, Xiaoxiang Shi, William Y. Tang, Manan Roongta, Colin Cai, Jeffrey Luo, Li Erran Li, Raluca Ada Popa, and Ion Stoica. 2025. Deepscaler: Surpassing o1-preview with a 1.5b model by scaling rl. <a href="#">Notion Blog</a> . Accessed 2025.	863
Sewon Min, Xinxu Lyu, Ari Holtzman, Mikel Artetxe, Mike Lewis, Hannaneh Hajishirzi, and Luke Zettlemoyer. 2022. <a href="#">Rethinking the role of demonstrations: What makes in-context learning work?</a> In <i>Proceedings of the 2022 Conference on Empirical Methods in Natural Language Processing, EMNLP 2022, Abu Dhabi, United Arab Emirates, December 7-11, 2022</i> , pages 11048–11064. Association for Computational Linguistics.	864
Eric Mitchell, Rafael Rafailov, Archit Sharma, Chelsea Finn, and Christopher D. Manning. 2024. <a href="#">An emulator for fine-tuning large language models using small language models</a> . In <i>The Twelfth International Conference on Learning Representations, ICLR 2024, Vienna, Austria, May 7-11, 2024</i> . OpenReview.net.	865
Niklas Muennighoff, Zitong Yang, Weijia Shi, Xiang Lisa Li, Li Fei-Fei, Hannaneh Hajishirzi, Luke Zettlemoyer, Percy Liang, Emmanuel J. Candès, and Tatsunori Hashimoto. 2025. <a href="#">s1: Simple test-time scaling</a> . <i>CoRR</i> , abs/2501.19393.	866
OpenAI. 2024. <a href="#">Learning to reason with LLMs</a> .	867
OpenAI. 2025. OpenAI o3 and o4-mini System Card. <a href="https://cdn.openai.com/pdf/2221c875-02dc-4789-800b-e7758f3722c1/o3-and-o4-mini-system-card.pdf">https://cdn.openai.com/pdf/2221c875-02dc-4789-800b-e7758f3722c1/o3-and-o4-mini-system-card.pdf</a> . Accessed: 2025-07-05.	868
Aitor Ormazabal, Mikel Artetxe, and Eneko Agirre. 2023. <a href="#">Comblm: Adapting black-box language models through small fine-tuned models</a> . In <i>Proceedings of the 2023 Conference on Empirical Methods in Natural Language Processing, EMNLP 2023, Singapore, December 6-10, 2023</i> , pages 2961–2974. Association for Computational Linguistics.	869
Shishir G. Patil, Huanzhi Mao, Charlie Cheng-Jie Ji, Fanjia Yan, Vishnu Suresh, Ion Stoica, and Joseph E. Gonzalez. 2025. The berkeley function calling leaderboard (bfcl): From tool use to agentic evaluation of large language models. In <i>Forty-second International Conference on Machine Learning</i> .	870



909	Qwen Team. 2025. <a href="#">Qwen3: Think Deeper, Act Faster   Qwen</a> .	
910		
911	Rafael Rafailov, Archit Sharma, Eric Mitchell, Christo-	
912	pher D. Manning, Stefano Ermon, and Chelsea Finn.	
913	2023. <a href="#">Direct preference optimization: Your language</a>	
914	<a href="#">model is secretly a reward model</a> . In <i>Advances in</i>	
915	<i>Neural Information Processing Systems 36: Annual</i>	
916	<i>Conference on Neural Information Processing Sys-</i>	
917	<i>tems 2023, NeurIPS 2023, New Orleans, LA, USA,</i>	
918	<i>December 10 - 16, 2023</i> .	
919	Abhinav Rastogi, Albert Q. Jiang, Andy Lo, Gabrielle	
920	Berrada, Guillaume Lample, Jason Rute, Joep Bar-	
921	mentlo, Karmesh Yadav, Kartik Khandelwal, Khy-	
922	athi Raghavi Chandu, Léonard Blier, Lucile Saulnier,	
923	Matthieu Dinot, Maxime Darrin, Neha Gupta, Ro-	
924	man Soletskyi, Sagar Vaze, Teven Le Scao, Yihan	
925	Wang, and 80 others. 2025. <a href="#">Magistral</a> . <i>CoRR</i> ,	
926	abs/2506.10910.	
927	David Rein, Betty Li Hou, Asa Cooper Stickland,	
928	Jackson Petty, Richard Yuanzhe Pang, Julien Di-	
929	rani, Julian Michael, and Samuel R. Bowman. 2023.	
930	<a href="#">GPQA: A graduate-level google-proof q&amp;a bench-</a>	
931	<a href="#">mark</a> . <i>CoRR</i> , abs/2311.12022.	
932	John Schulman, Filip Wolski, Prafulla Dhariwal, Alec	
933	Radford, and Oleg Klimov. 2017. <a href="#">Proximal policy</a>	
934	<a href="#">optimization algorithms</a> . <i>CoRR</i> , abs/1707.06347.	
935	Zhihong Shao, Peiyi Wang, Qihao Zhu, Runxin Xu,	
936	Junxiao Song, Mingchuan Zhang, Y. K. Li, Y. Wu,	
937	and Daya Guo. 2024. <a href="#">Deepseekmath: Pushing the</a>	
938	<a href="#">limits of mathematical reasoning in open language</a>	
939	<a href="#">models</a> . <i>CoRR</i> , abs/2402.03300.	
940	Weijia Shi, Xiaochuang Han, Mike Lewis, Yulia	
941	Tsvetkov, Luke Zettlemoyer, and Wen-tau Yih. 2024.	
942	<a href="#">Trusting your evidence: Hallucinate less with context-</a>	
943	<a href="#">aware decoding</a> . In <i>Proceedings of the 2024 Confer-</i>	
944	<i>ence of the North American Chapter of the Associ-</i>	
945	<i>ation for Computational Linguistics: Human Lan-</i>	
946	<i>guage Technologies: Short Papers, NAACL 2024,</i>	
947	<i>Mexico City, Mexico, June 16-21, 2024</i> , pages 783–	
948	791. Association for Computational Linguistics.	
949	Charlie Snell, Jaehoon Lee, Kelvin Xu, and Aviral Ku-	
950	mar. 2024. <a href="#">Scaling LLM test-time compute optimally</a>	
951	<a href="#">can be more effective than scaling model parameters</a> .	
952	<i>CoRR</i> , abs/2408.03314.	
953	Xinyu Tang, Xiaolei Wang, Zhihao Lv, Yingqian Min,	
954	Wayne Xin Zhao, Binbin Hu, Ziqi Liu, and Zhiqiang	
955	Zhang. 2025. <a href="#">Unlocking general long chain-of-</a>	
956	<a href="#">thought reasoning capabilities of large language mod-</a>	
957	<a href="#">els via representation engineering</a> . <i>arXiv preprint</i>	
958	<i>arXiv:2503.11314</i> .	
959	Qwen Team. 2025. <a href="#">Qwq-32b: Embracing the power of</a>	
960	<a href="#">reinforcement learning</a> .	
961	Chenglong Wang, Yang Gan, Yifu Huo, Yongyu Mu,	
962	Qiaozhi He, Murun Yang, Bei Li, Tong Xiao, Chun-	
963	liang Zhang, Tongran Liu, and 1 others. 2025. <a href="#">Gram:</a>	
964	<a href="#">A generative foundation reward model for reward</a>	
965	<a href="#">generalization</a> . <i>arXiv preprint arXiv:2506.14175</i> .	
	Peiyi Wang, Lei Li, Zhihong Shao, Runxin Xu, Damai	966
	Dai, Yifei Li, Deli Chen, Yu Wu, and Zhifang Sui.	967
	2024. <a href="#">Math-shepherd: Verify and reinforce llms step-</a>	968
	<a href="#">by-step without human annotations</a> . In <i>Proceedings</i>	969
	<i>of the 62nd Annual Meeting of the Association for</i>	970
	<i>Computational Linguistics (Volume 1: Long Papers),</i>	971
	<i>ACL 2024, Bangkok, Thailand, August 11-16, 2024,</i>	972
	pages 9426–9439. Association for Computational	973
	Linguistics.	974
	Xuezhi Wang, Jason Wei, Dale Schuurmans, Quoc V.	975
	Le, Ed H. Chi, Sharan Narang, Aakanksha Chowd-	976
	hery, and Denny Zhou. 2023. <a href="#">Self-consistency</a>	977
	<a href="#">improves chain of thought reasoning in language</a>	978
	<a href="#">models</a> . In <i>The Eleventh International Conference</i>	979
	<i>on Learning Representations, ICLR 2023, Kigali,</i>	980
	<i>Rwanda, May 1-5, 2023</i> . OpenReview.net.	981
	Xuezhi Wang and Denny Zhou. 2024. <a href="#">Chain-of-thought</a>	982
	<a href="#">reasoning without prompting</a> . In <i>Advances in Neural</i>	983
	<i>Information Processing Systems 38: Annual Confer-</i>	984
	<i>ence on Neural Information Processing Systems 2024,</i>	985
	<i>NeurIPS 2024, Vancouver, BC, Canada, December</i>	986
	<i>10 - 15, 2024</i> .	987
	Jason Wei, Xuezhi Wang, Dale Schuurmans, Maarten	988
	Bosma, Brian Ichter, Fei Xia, Ed H. Chi, Quoc V. Le,	989
	and Denny Zhou. 2022. <a href="#">Chain-of-thought prompting</a>	990
	<a href="#">elicits reasoning in large language models</a> . In <i>Ad-</i>	991
	<i>vances in Neural Information Processing Systems 35:</i>	992
	<i>Annual Conference on Neural Information Process-</i>	993
	<i>ing Systems 2022, NeurIPS 2022, New Orleans, LA,</i>	994
	<i>USA, November 28 - December 9, 2022</i> .	995
	Jian Xie, Kai Zhang, Jiangjie Chen, Tinghui Zhu, Renze	996
	Lou, Yuandong Tian, Yanghua Xiao, and Yu Su. 2024.	997
	<a href="#">Travelplanner: A benchmark for real-world planning</a>	998
	<a href="#">with language agents</a> . In <i>Forty-first International</i>	999
	<i>Conference on Machine Learning, ICML 2024, Vi-</i>	1000
	<i>enna, Austria, July 21-27, 2024</i> . OpenReview.net.	1001
	Yuxi Xie, Kenji Kawaguchi, Yiran Zhao, James Xu	1002
	Zhao, Min-Yen Kan, Junxian He, and Michael Qizhe	1003
	Xie. 2023. <a href="#">Self-evaluation guided beam search for</a>	1004
	<a href="#">reasoning</a> . In <i>Advances in Neural Information Pro-</i>	1005
	<i>cessing Systems 36: Annual Conference on Neural</i>	1006
	<i>Information Processing Systems 2023, NeurIPS 2023,</i>	1007
	<i>New Orleans, LA, USA, December 10 - 16, 2023</i> .	1008
	Haotian Xu, Xing Wu, Weinong Wang, Zhongzhi Li,	1009
	Da Zheng, Boyuan Chen, Yi Hu, Shijia Kang, Ji-	1010
	aming Ji, Yingying Zhang, Zhijiang Guo, Yaodong	1011
	Yang, Muhan Zhang, and Debing Zhang. 2025. <a href="#">Red-</a>	1012
	<a href="#">star: Does scaling long-cot data unlock better slow-</a>	1013
	<a href="#">reasoning systems?</a> <i>CoRR</i> , abs/2501.11284.	1014
	An Yang, Baosong Yang, Beichen Zhang, Binyuan	1015
	Hui, Bo Zheng, Bowen Yu, Chengyuan Li, Dayi-	1016
	heng Liu, Fei Huang, Haoran Wei, Huan Lin, Jian	1017
	Yang, Jianhong Tu, Jianwei Zhang, Jianxin Yang, Ji-	1018
	axi Yang, Jingren Zhou, Junyang Lin, Kai Dang, and	1019
	22 others. 2024a. <a href="#">Qwen2.5 technical report</a> . <i>CoRR</i> ,	1020
	abs/2412.15115.	1021



1022	An Yang, Beichen Zhang, Binyuan Hui, Bofei Gao,	Yaowei Zheng, Richong Zhang, Junhao Zhang, Yanhan	1078
1023	Bowen Yu, Chengpeng Li, Dayiheng Liu, Jian-	Ye, Zheyang Luo, Zhangchi Feng, and Yongqiang Ma.	1079
1024	hong Tu, Jingren Zhou, Junyang Lin, Keming Lu,	2024b. <a href="#">Llamafactory: Unified efficient fine-tuning</a>	1080
1025	Mingfeng Xue, Runji Lin, Tianyu Liu, Xingzhang	of 100+ language models. In <i>Proceedings of the</i>	1081
1026	Ren, and Zhenru Zhang. 2024b. <a href="#">Qwen2.5-math tech-</a>	62nd Annual Meeting of the Association for Computa-	1082
1027	anical report: Toward mathematical expert model via	tional Linguistics (Volume 3: System Demonstra-	1083
1028	self-improvement. <i>CoRR</i> , abs/2409.12122.	tions), Bangkok, Thailand. Association for Computa-	1084
1029	Wang Yang, Xiang Yue, Vipin Chaudhary, and Xiao-	tional Linguistics.	1085
1030	tian Han. 2025. <a href="#">Speculative thinking: Enhancing</a>		
1031	<a href="#">small-model reasoning with large model guidance at</a>	Andy Zou, Long Phan, Sarah Li Chen, James Campbell,	1086
1032	<a href="#">inference time</a> . <i>CoRR</i> , abs/2504.12329.	Phillip Guo, Richard Ren, Alexander Pan, Xuwang	1087
1033	Yixin Ye, Zhen Huang, Yang Xiao, Ethan Chern, Shijie	Yin, Mantas Mazeika, Ann-Kathrin Dombrowski,	1088
1034	Xia, and Pengfei Liu. 2025. <a href="#">LIMO: less is more for</a>	Shashwat Goel, Nathaniel Li, Michael J. Byun, Zifan	1089
1035	<a href="#">reasoning</a> . <i>CoRR</i> , abs/2502.03387.	Wang, Alex Mallen, Steven Basart, Sanmi Koyejo,	1090
1036	Qiying Yu, Zheng Zhang, Ruofei Zhu, Yufeng Yuan,	Dawn Song, Matt Fredrikson, and 2 others. 2023.	1091
1037	Xiaochen Zuo, Yu Yue, Tiantian Fan, Gaohong Liu,	<a href="#">Representation engineering: A top-down approach</a>	1092
1038	Lingjun Liu, Xin Liu, Haibin Lin, Zhiqi Lin, Bole	to AI transparency. <i>CoRR</i> , abs/2310.01405.	1093
1039	Ma, Guangming Sheng, Yuxuan Tong, Chi Zhang,		
1040	Mofan Zhang, Wang Zhang, Hang Zhu, and 16 others.	<b>A Technical Details</b>	1094
1041	2025. <a href="#">DAPO: an open-source LLM reinforcement</a>		
1042	<a href="#">learning system at scale</a> . <i>CoRR</i> , abs/2503.14476.	<b>A.1 Training Details</b>	1095
1043	Zheng Yuan, Hongyi Yuan, Chengpeng Li, Guanting	<b>Environment.</b> All experiments were conducted	1096
1044	Dong, Chuanqi Tan, and Chang Zhou. 2023. <a href="#">Scaling</a>	using NVIDIA A40/L40S GPUs with 48GB mem-	1097
1045	<a href="#">relationship on learning mathematical reasoning with</a>	ory. The software environment was configured as	1098
1046	<a href="#">large language models</a> . <i>CoRR</i> , abs/2308.01825.	follows:	1099
1047	Yang Yue, Zhiqi Chen, Rui Lu, Andrew Zhao, Zhaokai	• 360-LLaMA-Factory ( <a href="#">Haosheng Zou and</a>	1100
1048	Wang, Shiji Song, and Gao Huang. 2025a. Does	<a href="#">Zhang, 2024</a> ) (A long-CoT adapted version of	1101
1049	reinforcement learning really incentivize reasoning	LLaMA-Factory 0.9.1 ( <a href="#">Zheng et al., 2024b</a> ))	1102
1050	capacity in llms beyond the base model? <i>arXiv</i>		
1051	<i>preprint arXiv:2504.13837</i> .	• torch 2.7.0	1103
1052	Yang Yue, Zhiqi Chen, Rui Lu, Andrew Zhao, Zhaokai		
1053	Wang, Yang Yue, Shiji Song, and Gao Huang. 2025b.	• transformers 4.51.3	1104
1054	<a href="#">Does reinforcement learning really incentivize rea-</a>		
1055	<a href="#">soning capacity in llms beyond the base model?</a>	• accelerate 1.0.1	1105
1056	<i>CoRR</i> , abs/2504.13837.		
1057	Lunjun Zhang, Arian Hosseini, Hritik Bansal, Mehran	• datasets 3.1.0	1106
1058	Kazemi, Aviral Kumar, and Rishabh Agarwal. 2025.		
1059	<a href="#">Generative verifiers: Reward modeling as next-token</a>	• trl 0.9.6	1107
1060	<a href="#">prediction</a> . In <i>The Thirteenth International Confer-</i>		
1061	<i>ence on Learning Representations, ICLR 2025, Sin-</i>	• peft 0.12.0	1108
1062	<i>gapore, April 24-28, 2025</i> . OpenReview.net.		
1063	Xuandong Zhao, Xianjun Yang, Tianyu Pang, Chao Du,	• deepspeed 0.14.4	1109
1064	Lei Li, Yu-Xiang Wang, and William Yang Wang.		
1065	2024. <a href="#">Weak-to-strong jailbreaking on large language</a>	<b>LoRA Configuration.</b> We applied LoRA ( <a href="#">Hu</a>	1110
1066	<a href="#">models</a> . <i>CoRR</i> , abs/2401.17256.	<a href="#">et al., 2022</a> ) for parameter-efficient fine-tuning of	1111
1067	Zekai Zhao, Qi Liu, Kun Zhou, Zihan Liu, Yifei	the guider model:	1112
1068	Shao, Zhiting Hu, and Biwei Huang. 2025. Acti-		
1069	vation control for efficiently eliciting long chain-of-	• Rank: 64	1113
1070	thought ability of language models. <i>arXiv preprint</i>		
1071	<i>arXiv:2505.17697</i> .	• $\alpha_{\text{LoRA}}$ : 128	1114
1072	Huaixiu Steven Zheng, Swaroop Mishra, Hugh Zhang,		
1073	Xinyun Chen, Minmin Chen, Azade Nova, Le Hou,	• Target modules: q_proj, k_proj, v_proj,	1115
1074	Heng-Tze Cheng, Quoc V. Le, Ed H. Chi, and	o_proj	1116
1075	Denny Zhou. 2024a. <a href="#">NATURAL PLAN: bench-</a>		
1076	<a href="#">marking llms on natural language planning</a> . <i>CoRR</i> ,	• Bias: None	1117
1077	abs/2406.04520.		

**DPO Training.** For preference optimization with DPO, we used the following settings:

- Batch size: 32 (4 GPUs \* 8 Gradient Accumulation)
- Epoch: 1
- Learning rate: 5e-6
- Optimizer: AdamW
- Learning rate scheduler: cosine with warmup
- Warmup ratio: 0.1
- $\beta$  (reward scaling): 0.1
- Cutoff length: 8192