
Training Language Models to Reason Efficiently

Daman Arora

Carnegie Mellon University
damana@andrew.cmu.edu

Andrea Zanette

Carnegie Mellon University
zanette@cmu.edu

Abstract

Scaling model size and training data has led to great advances in the performance of Large Language Models (LLMs). However, the diminishing returns of this approach necessitate alternative methods to improve model capabilities, particularly in tasks requiring advanced reasoning. Large reasoning models, which leverage long chain-of-thoughts, bring unprecedented breakthroughs in problem-solving capabilities but at a substantial deployment cost associated to longer generations. Reducing inference costs is crucial for the economic feasibility, user experience, and environmental sustainability of these models.

In this work, we propose to train large reasoning models to reason efficiently. Our method incentivizes models to minimize unnecessary computational overhead while largely maintaining accuracy, thereby achieving substantial deployment efficiency gains. It enables the derivation of a family of reasoning models with varying efficiency levels, controlled via a single hyperparameter. Experiments on two open-weight large reasoning models demonstrate significant reductions in inference cost while preserving most of the accuracy.

1 Introduction

Large language models (LLMs) have made significant advancements by pre-training larger models with extensive datasets (Kaplan et al., 2020), but this approach faces diminishing returns due to limited high-quality training data. An alternative to improve model capabilities, especially in domains involving careful reasoning, involves allowing models to “think” before answering, as seen in frontier reasoning models like OpenAI’s o1 (OpenAI et al., 2024), Gemini 2.0 Flash Thinking Experimental, and DeepSeek-R1 (Guo et al., 2025). These models produce intermediate tokens during inference, collectively referred to as *chain-of-thoughts* (Wei et al., 2022), to perform additional computations before returning an answer. The process of generating a long chain of thought before answering the user query is called *reasoning*. More precisely, *large reasoning models* with chain-of-thoughts capable of performing advanced reasoning emerge from RL (Sutton & Barto, 2018; Guo et al., 2025) on base models using ground-truth scoring functions (e.g., correctness on math problems).

These reasoning models use test-time compute in the form of very long chain-of-thoughts, an approach that incur a high inference cost due to the quadratic cost of the attention mechanism and linear growth of the KV cache for transformer-based architectures Vaswani (2017). However, effective deployment of LLMs demands models that are not only capable but also computationally efficient to serve. Even for resource-rich organizations such as large tech companies that have the resources to train reasoning models, excessive inference costs may mean operating at a loss rather than at a profit in order to match the competitor’s offering. Furthermore, reducing inference costs often reduces latency, improves responsiveness, and therefore improves user experience. Finally, lowering the inference computation positively benefits the environment by reducing carbon emissions.

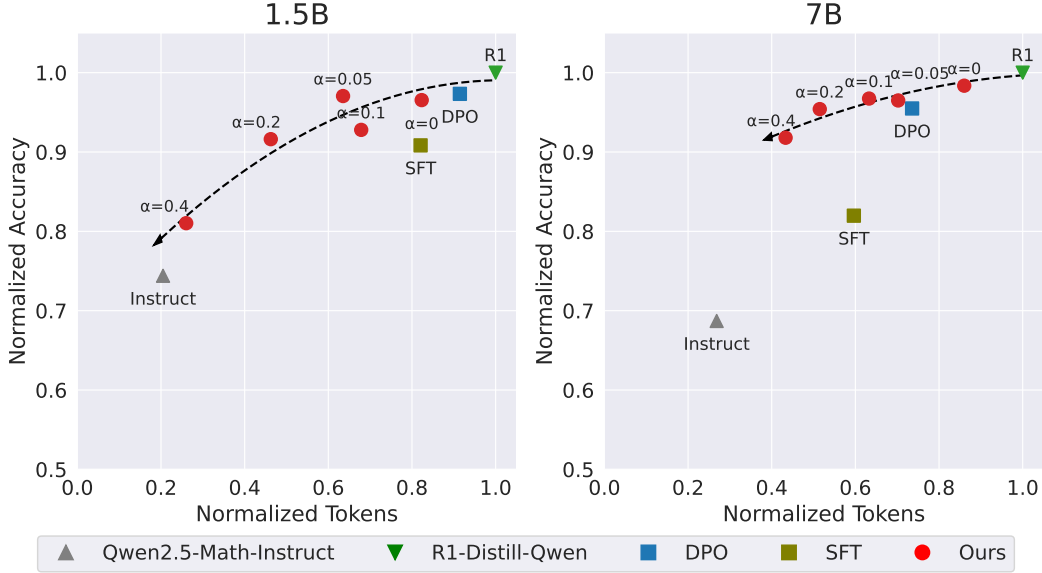


Figure 1: This figure describes the results of our training procedure aggregated on 5 evaluation datasets on both 1.5B and 7B scales. The x-axis represents response length normalized by the performance of the R1-Distilled model and the y-axis represents accuracy normalized by the accuracy of the R1-Distilled model. Each value of α has been aggregated over 3 seeds. Our training reward design allows us to easily and effectively navigate the token-accuracy tradeoff curve. Increasing the value of α generally results in models that have lower response lengths.

To achieve this goal, we use policy gradient methods (Sutton & Barto, 2018) to train the model to use lower number of tokens to reach the correct solution, thereby minimizing inference costs, ideally without compromising on accuracy. We use a modified reinforcement learning reward function which *encourages the model to produce correct answers with short chain-of-thoughts*. At the time of the initial publication of our pre-print on arXiv, we were among the first to consider training the model to be efficient at inference time, and we discuss concurrent as well as related literature in Section 2. Our method allows the user to *control* the reduction in inference compute by adjusting a scalar coefficient in an intuitive way. In other words, starting from a reasoning model, our procedure allows us to derive a *family* of models, each with increased generation efficiency (i.e., shorter chain-of-thoughts).

We perform numerical experiments on two recently released open-weight large reasoning models, DeepSeek-R1-Distill-Qwen-1.5B and DeepSeek-R1-Distill-Qwen-7B Guo et al. (2025) and derive models with a substantial reduction in reasoning cost while approximately maintaining accuracy, see Figure 1 for a summary of our results. We observe that our training procedure allows us to gracefully navigate the compute-performance tradeoff curve, that is, reduce compute significantly with minimal loss in performance. For instance, for the 7B model, we can achieve a reduction of 50% in tokens with less than 5% reduction in accuracy.

We also observe that the reduction in response length depends on the hardness of the problem. For the 7B model, our method can reduce token usage by 16% on the competition-level benchmark: American Invitational Mathematics Examination 2024 with a drop of 3.3% points in accuracy, up to 37% tokens on the MATH500 (Hendrycks et al., 2021) dataset with a drop of 2.2% in accuracy, and up to 65% tokens on the GSM8K (Cobbe et al., 2021b) dataset with a drop of 1.7% in accuracy.

Our work only requires a couple of lines of changes to any standard reinforcement learning implementation. Beyond its simplicity, an attractive property of our approach is its *computational efficiency*: although training reasoning models with large scale RL has a prohibitive cost Guo et al. (2025), our procedure shows that training them to reason efficiently is highly viable even with modest academic resources: our models are obtained with only 100 RL steps (~ 200 gradient updates). The fact that we achieve a performance comparable to that of the original reasoning model with a short training is surprising because in a few RL steps the model learns to optimize for shorter, more efficient reasoning patterns compared to the original model. All of our code and trained models would be made public.

2 Related Work

Improving model capabilities with test-time compute Several techniques such as Chain of Thoughts (Wei et al., 2022), Self Consistency (Wang et al., 2022), Best-of-N, Monte Carlo Tree Search (Silver et al., 2017), Tree-of-thoughts (Yao et al., 2024), Stream-of-Search (Gandhi et al., 2024), Graph-of-Thoughts (Besta et al., 2024), Process Reward Models (Lightman et al., 2024) and Self-Correction (Kumar et al., 2024; Welleck et al., 2023) have been proposed to improve performance by spending more inference time compute. It has also been shown that in some cases scaling inference compute can be more effective than scaling model size (Snell et al., 2025; Liu et al., 2025a). While the above mentioned techniques can be effective in specialized scenarios, modern large scale reasoning models are trained with large scale RL with verifiable rewards and perform autoregressive generation.

Large Reasoning Models Frontier reasoning such as OpenAI o1, Deepseek R1 and QwQ-Preview rely on long, monolithic chain-of-thoughts to perform advanced reasoning. They are trained with large scale reinforcement learning Guo et al. (2025), which leads them to develop abilities such as branching, verification and backtracking. Our approach aims to make these models more efficient.

Efficient Serving While we focus on developing reasoning models that can be served efficiently, our approach is orthogonal to existing methods from the literature of efficient LLMs; see Zhou et al. (2024) for a recent survey. For example, system-level techniques build a system to accelerate inference. Some examples include speculative decoding Leviathan et al. (2023) and batch engines like vLLM Kwon et al. (2023a); both can be directly combined with our method. Model-based techniques, on the other hand, act directly on the model to accelerate inference. Some examples include weight pruning Liu et al. (2018) and quantization Lin et al. (2024), which are all complementary and can be combined with our methodology.

Concurrent works Chen et al. (2024) investigate the overthinking phenomena and propose methods to mitigate it by using heuristics such as First-Correct Solutions (FCS) and Greedy Diverse Solutions (GDS) to generate preference data which is then used for offline policy optimization. However, this method doesn’t allow easily tuning the model to the user’s compute budget. At the time of our pre-print coming online on arXiv, several papers with similar objectives were released as well. For instance, the Kimi k1.5 (Team et al., 2025) also report a method to shorten the chain-of-thought using a length penalty in the reward function while doing online RL, a procedure similar in principle but not identical to ours. We note that their procedure does not appear to have a tunable parameter which allows to obtain a family of models—each with varying trade-offs—as we do. Another work in this direction is O1-Pruner (Luo et al., 2025) which proposes a slightly different offline RL objective to minimize tokens while maintaining accuracy while we use online RL. Subsequent to our pre-print, multiple works have tried to look into this research direction. For instance, Qu et al. (2025) pose the problem of minimizing compute usage as a meta-RL problem and use online RL to solve it. Aggarwal & Welleck (2025) train reasoning models to follow exact token constraints specified in prompts. Xia et al. (2025) prune tokens from Chain-of-Thoughts using semantic importance and perform SFT to obtain models with controllable compression.

Efficiency of Chain-of-Thought Jin et al. (2024) find that lengthening chain-of-thought has a correlation with improving performance. Conditional training as done by Kang et al. (2024) is also another approach to the problem of generating shorter chain-of-thoughts. Explicitly trying to control the number of tokens by prompt engineering has been explored by Nayab et al. (2025) and Han et al. (2024). However, none of these methods have explored models that generate a long CoT and don’t use RL to train models to be less verbose. Also, they hinge on the critical assumption that models can follow instructions of the format: ‘Please answer this query in less than X tokens’. However, we find that distilled reasoning models are not capable of following such instructions. In fact, there has been work on this after our pre-print was released which trains models to explicitly follow such instructions (Aggarwal & Welleck, 2025). See Appendix Section F for a detailed discussion.

3 Setup

Let p be a language model. When provided with a prompt x , the language model produces a response $y = (y^1, y^2, \dots, y^t)$, where y^i represents the i -th token in the response and t is the total number of tokens in the response sequence. More precisely, the generation is *auto-regressive*, meaning that given the prompt x and the tokens $y^{\leq k} = (y^1, y^2, \dots, y^k)$ generated so far, the next token y^{k+1} is

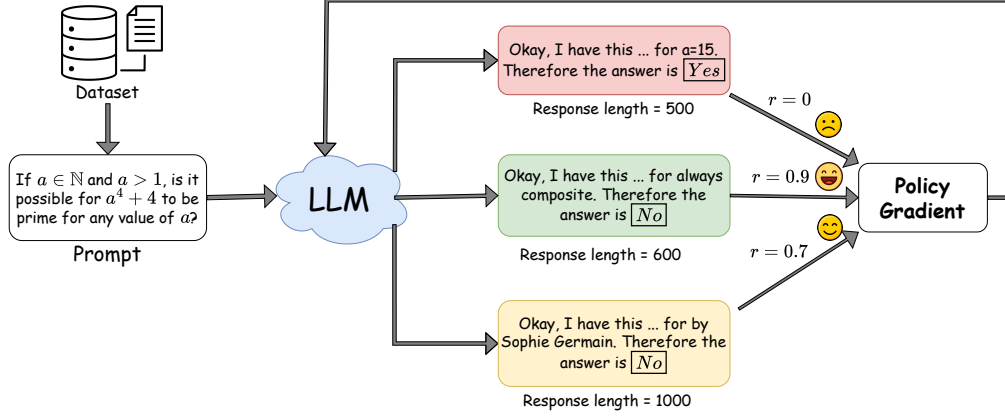


Figure 2: Pipeline depicting our method. For every prompt, multiple solutions are sampled and rewarded based on correctness and response length. The shortest correct answers are rewarded the highest and the language model is then updated using policy gradients.

generated from the conditional model

$$y^{k+1} \sim p(\cdot \mid x, y^{\leq k}). \quad (1)$$

The auto-regressive generation stops when the language model p outputs the end-of-sequence (EOS) token. Therefore, if $y = (y^1, y^2, \dots, y^t)$ is a full response, y^t is always the EOS token. With a little abuse of notation, we also let $y \sim p(\cdot \mid x)$ denote the process of sampling the full response $y = (y^1, y^2, \dots, y^t)$ from the model p via auto-regressive sampling according to Equation (1).

Chain-of-Thoughts Chain of thoughts, introduced by Wei et al. (2022), is a key framework to implement reasoning. Given a prompt x , the LLM is said to produce a “chain of thought” when it produces intermediate tokens that are not part of the output before generating the final answer in an autoregressive way. Typically, the final answer is not formally separated from the chain-of-thoughts, and so we let y denote the full output of the model $y \sim p(x)$.

Objective function and reinforcement learning We consider problems where the responses generated from an LLM can be evaluated by a scoring function $f(x, y) \mapsto \mathbb{R}$, often called *reward model* or *verifier*, that measures the suitability of the response. For math problems, such as those that we consider in this paper, the reward function establishes whether the solution to the problem is correct (Cobbe et al., 2021b; Hendrycks et al., 2021)

$$f(x, y) = 1\{y = y^*(x)\} \quad (2)$$

where $y^*(x)$ is the correct answer to the math problem x . Since y is the full output of the model, including the chain of thought, the relation $y = y^*(x)$ tests whether the final answer generated by the model coincides with the gold answer, rather than checking the equivalence between strings.

Large reasoning models Guo et al. (2025) are reportedly trained with reinforcement learning Sutton & Barto (2018). When a chain of thoughts is used, the objective function to maximize can be written as

$$\text{ACCURACY}(p) = \mathbb{E}_{x \sim \rho} \mathbb{E}_{y \sim p(x)} [1\{y = y^*(x)\}]. \quad (3)$$

where ρ is the prompt distribution. In the sequel, we simply write \mathbb{E} to denote the expectation. For math problems, maximizing Equation (3) directly maximizes the probability that the model correctly solves a random question from the prompt distribution.

4 Method

We aim to design a method that trains models to use lower amount of inference time compute to arrive at the correct answer. For simpler math problems, such as those in GSM8K Cobbe et al. (2021b),

the model should recognize when it has reached the correct solution within a few hundred tokens. In contrast, for competition-level problems like those in the American Invitational Mathematics Examination (AIME), the model should be capable of expending thousands of tokens if that is necessary to find a strategy that solves these exceptionally challenging questions.

One attractive option is to train the model on an objective function derived from Equation (3) that *encourages the model to produce correct solutions with the minimum amount of tokens*. In order to achieve the latter goal, we penalize the length of the correct responses

$$\mathbb{E}\left[1\{y = y^*(x)\}(1 - \alpha f(\text{LEN}(y)))\right] \quad (4)$$

using a monotonic function f of the input and a tunable parameter $\alpha \in [0, 1)$. The choice $\alpha = 0$ yields the reinforcement learning objective (3); increasing α increases the regularization towards shorter—but correct—responses.

In order to ensure that the length regularization is effective, we first normalize the length of the responses and then use the sigmoid function σ to soft-clip it, obtaining

$$f(\text{LEN}(y)) = \sigma\left(\frac{\text{LEN}(y) - \text{MEAN}(x)}{\text{STD}(x)}\right) \quad (5)$$

where

$$\text{MEAN}(x) = \mathbb{E}_{\substack{y \sim p(x), \\ \text{s.t. } 1\{y=y^*\}=1}} [\text{LEN}(y)], \quad \text{STD}(x) = \sqrt{\text{Var}_{\substack{y \sim p(x), \\ \text{s.t. } 1\{y=y^*\}=1}} [\text{LEN}(y)]}$$

are the *per-prompt* mean and standard deviation of the length, respectively. The per-prompt normalization ensures that longer chains of thought on hard problems are not disproportionately penalized compared to shorter ones on easier problems. When $\alpha \in [0, 1)$, the sigmoid ensures that the objective function is always bounded between $[0, 1]$ even for abnormally long or short generations, and that correct responses, even if long, are always preferred to incorrect ones. In practice, both the standard deviation and the mean are directly estimated from the rollouts during online training.

4.1 Optimizing the objective function with Reinforcement Learning

Since optimizing Equation (4) involves sampling from the model auto-regressively, the objective function is non-differentiable; however, it can be optimized with reinforcement learning, for instance with policy gradient methods (Sutton & Barto, 2018).

One popular option is proximal policy optimization (PPO) Schulman et al. (2017) which considers the (local) objective function

$$\min\{f_\theta^t(y, x)\mathcal{A}(y^{<t}, x), \text{clip}_{1-\epsilon}^{1+\epsilon}[f_\theta^t(y, x)]\mathcal{A}(y^{<t}, x)\}$$

defined using the density ratio

$$f_\theta^t(y, x) = \frac{\pi_\theta(y^t|x + y^{<t})}{\pi_{old}(y^t|x + y^{<t})}$$

and for a suitable choice for the advantage estimator $\mathcal{A}(y^{<t}, x)$. Traditionally, in deep reinforcement learning Schulman et al. (2017) the advantage estimator involves a neural network.

With language models, maintaining a separate value network to obtain a variance-reduced advantage estimator Schulman et al. (2017) may add significant computational and implementation complexity without necessarily increasing performance Kool et al. (2019); Ahmadian et al. (2024). One simple and effective alternative is to just estimate the advantage using Monte Carlo (MC) as proposed by Kool et al. (2019); Ahmadian et al. (2024). Such estimator is also called REINFORCE Leave One Out (RLOO) estimator. To be precise, the trajectory advantage can be estimated as

$$\mathcal{A}(y_i, x) = \mathcal{R}(y_i, x) - \frac{1}{n-1} \sum_{j \neq i} \mathcal{R}(y_j, x)$$

where \mathcal{R} is the trajectory return and y_i is the i^{th} generation for prompt x . We then simply use the sequence level advantage as the token level advantage, namely $\mathcal{A}(y^{<t}, x) = \mathcal{A}(y, x)$. In essence, we use PPO with the RLOO advantage estimator.

4.2 Population-level optimality guarantees

In this section we analyze the population-level maximizer of Equation (4) in a highly simplified setup and show how this can lead to the desired behavior of shortening the chain-of-thoughts without compromising accuracy. Additionally, we establish that the population-level maximizer is a model that generates the shortest correct response for any prompt in \mathcal{X} .

Consider the following simplified setup, where the language model p_θ conditioned on a prompt x is a multinomial distribution over N possible responses y_1, \dots, y_N . More precisely, given $|\mathcal{X}|$ multinomial distributions $p(\cdot | x)$ on the prompt space \mathcal{X} , there exists a value of the parameter θ that realizes such a choice. This is formalized by the following Assumptions.

Assumption 4.1 (Tabular Representation). For every choice of p such that

$$p(y_i | x) \in [0, 1], \quad \forall x \in \mathcal{X}, i \in [N] \quad (6)$$

$$\sum_i p(y_i | x) = 1, \quad \forall x \in \mathcal{X} \quad (7)$$

there exists a θ such that

$$p_\theta(y_i | x) = p(y_i | x), \quad \forall i \in [N], \forall x \in \mathcal{X}. \quad (8)$$

This assumption can be justified by the expressive power of the neural network. The following assumption ensures coverage, namely that for every prompt, there exists at least a correct response that the LLM can output for an appropriate value of θ . It encodes the fact that an LLM can learn the correct solution if given enough data.

Assumption 4.2 (Coverage). For all prompts $x \in \mathcal{X}$, $\exists y \in \{y_i\}_{i=1}^N$ such that $y = y^*(x)$.

Let p_{θ^*} denote the reasoning model that is the population level maximizer of the accuracy:

$$\theta^* = \arg \max_{\theta} \mathbb{E}_{x \sim \rho} \mathbb{E}_{y \sim p_\theta(x)} [1\{y = y^*(x)\}] \quad (9)$$

where ρ is the distribution over the prompts. From our simplified setup, in particular from Assumption 4.1 and Assumption 4.2 it is easy to see that

$$\text{ACCURACY}(p_{\theta^*}) = 1. \quad (10)$$

In other words, if the language model has enough expressive power that it can cover the correct solution for each of the prompts, maximization of the population level RL training objective (9) leads to a model that can output the correct solution over each prompt in the training dataset.

Let θ_{eff}^* denote the population-level parameters of the reasoning model obtained by maximizing Equation (4), i.e.,

$$\theta_{eff}^* = \arg \max_{\theta} \left\{ \mathbb{E}_{x \sim \rho} \mathbb{E}_{y \sim p_\theta(x)} [1\{y = y^*(x)\}(1 - \alpha f(\text{LEN}(y)))] \right\} \quad (11)$$

for a certain choice of a monotonically increasing function $f(\cdot) \in [0, 1]$ and scalar value $\alpha \in [0, 1]$.

We can prove that the population-level maximizer $p_{\theta_{eff}^*}$ is as accurate as the population-level maximizer p_{θ^*} and that it recovers the shortest solution for every prompt in \mathcal{X} . For brevity, the proofs are moved to Section E in the Appendix.

Proposition 4.3 (Accuracy is Preserved). *With the setup just described,*

$$\text{ACCURACY}(p_{\theta_{eff}^*}) = 1. \quad (12)$$

Proposition 4.4 (Training Objective yields the Shortest Correct Solution). *Under the assumptions listed, $\forall x \in \mathcal{X}$, $\forall y'$ such that $y' = y^*(x)$,*

$$\mathbb{E}_{y \sim p_{\theta_{eff}^*}(x)} [\text{LEN}(y)] \leq \text{LEN}(y') \quad (13)$$

Note that our claim is about the population-level maximizer; finite-sample guarantees can be obtained for parametric and nonparametric models using standard techniques from statistics Wainwright (2019).

Intuitively, the average length is reduced by virtue of our objective function (4), while accuracy is preserved in the idealized setting that we consider. Note that in our work, the function f is not a true function since f is computed using data generated in the sampling process, specifically $\text{MEAN}(x)$ and $\text{STD}(x)$.

5 Experiments

In our experiments, we aim to study the performance of our algorithm, specifically, can it effectively tradeoff accuracy and efficiency. Additionally, we compare with some simple baselines. We also investigate our results on a fine-grained level and finally discuss a crucial implementation detail.

5.1 Setup

Models and Datasets We evaluate all algorithms on the DeepSeek-R1-Distill-Qwen-1.5B and DeepSeek-R1-Distill-Qwen-7B Guo et al. (2025) models. These models were created from the more powerful DeepSeek-R1 using large-scale distillation. Along with a LLaMA-variant distilled by the same authors Guo et al. (2025), they are the only open-weight reasoning models of their size.

For post-training the model using our technique, we choose 3.2k prompts from the MATH, cn_k12, AIME, AoPS and the Olympiad subsets of the Numina Math dataset LI et al. (2024). The dataset includes problems that lack an objective answer, such as proof-based questions. We filter out such problems and ensure that the selected training problems have a numerical answer that can be parsed.

Evaluation We report the training logs and also evaluate the models on three test datasets namely: GSM8K Cobbe et al. (2021a), which contains grade-school-level math problems, MATH500 Hendrycks et al. (2021) which is a standard benchmark containing harder problems than GSM8K, and The American Invitational Mathematics Examination (AIME) 2024, a competition-level dataset of challenging mathematical problems. Additionally, to verify the robustness of our training methodology to datasets other than those based on mathematics, we evaluate models on CommonSenseQA and Logical Deduction from BIG-Bench (Srivastava et al., 2023).

For all models, we set the temperature to 0.6 as suggested in the model’s card¹ and set the token limit to 32K. We use vLLM Kwon et al. (2023b) for efficient batch inference. We use the parser created by the Qwen Team for the evaluation of their models² to measure correctness. We report the *average pass rate@k* for all models. Specifically, for each prompt, we sample k responses and compute the average accuracy per prompt, which is then averaged across the entire dataset. For GSM8K, we set $k = 1$ due to its large number of test samples. In contrast, for MATH500, we use $k = 3$, and for AIME2024, we set $k = 10$ given its limited set of only 30 questions. Implementation details along with computational requirements are given in Section A of the Appendix.

5.2 Baselines

Apart from the concurrent and related work discussed in Section 2, to our knowledge there are no prior studies in this setting. Alongside our method, we introduce and implement simple baseline approaches that help balance inference cost and accuracy.

1. **Generation Cutoff:** This simple baseline imposes a maximum token limit during the vLLM generation. If a response exceeds the token limit and remains incomplete, it is assigned a score of 0. We evaluate token cutoffs at 8,000, 16,000, 24,000, and 32,000.
2. **Rejection Sampling + SFT:** In this baseline, we generate 8 solutions per prompt using the distilled 1.5B and 7B models. From the generated solutions, we select the *shortest correct* responses and perform SFT on those responses. For a dataset of 3,200 prompts, this process yields approximately 2,200 and 2,500 valid responses for the 1.5B and 7B models, respectively. We experiment with three learning rates: 1×10^{-5} , 5×10^{-6} , and 2×10^{-6} . We find that 5×10^{-6} effectively reduces response length in a meaningful way.
3. **DPO:** Using the same dataset as above, we select response pairs consisting of the longest and shortest correct solutions and apply Direct Preference Optimization (DPO) Rafailov et al. (2023) on these preference pairs. While other preference optimization algorithms are applicable in this setting, we choose DPO for its popularity and ease of use. Similar to the SFT baseline, we experiment with three learning rates: 1×10^{-5} , 5×10^{-6} , and 2×10^{-6} . We observe that 1×10^{-5} effectively reduces response length.

¹<https://huggingface.co/deepseek-ai/DeepSeek-R1-Distill-Qwen-7B>

²<https://github.com/QwenLM/Qwen2.5-Math>

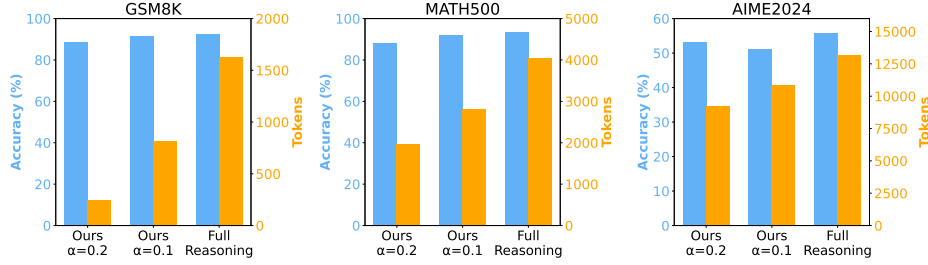


Figure 3: Our procedure trains models to be more token-efficient on easier problems, such as GSM8K, while preserving accuracy on harder problems, such as AIME2024. Full Reasoning refers to the reasoning model DeepSeek-R1-Distill-Qwen-7B. It is noteworthy that for the same value of α , there is higher reduction in tokens for an easier dataset (GSM8K) vs. a harder dataset (AIME2024)

5.3 Results

We train DeepSeek-R1-Distill-Qwen-1.5B and DeepSeek-R1-Distill-Qwen-7B models using different values of $\alpha \in [0, 0.05, 0.1, 0.2, 0.4]$ to illustrate the trade-offs between models with different lengths for the chain-of-thoughts. We report the aggregate results in Figure 1 and discuss the detailed results below. Note that the Generation Cutoff baseline has been excluded to make the figure more readable. Results including the Generation Cutoff baseline can be found in Figure 5.

5.3.1 Performance on the test sets

We provide results for MATH500, GSM8K and AIME2024 in Figure 3. Detailed results for all datasets are provided in Section C of the Appendix. As evident from Figure 1, our method enables smooth trade-offs of compute cost and accuracy, allowing models to be tailored to the specific requirements of downstream tasks or users based on different values of α . For instance, with $\alpha = 0.1$, the length of the chain-of-thought of the 7B model on the MATH dataset decreases by 36% (from ~ 4000 to ~ 2600 tokens) while the accuracy loss is only 2.2%. Similarly, in the AIME dataset, setting $\alpha = 0.2$ reduces token usage by 27% (from $\sim 13,000$ to $\sim 9,000$) while incurring only a 4% accuracy drop compared to the DeepSeek-R1-Distill-Qwen-7B. We offer several remarks:

Firstly, our training method allows for a *more performant* model compared to SFT and DPO given the same token usage as seen in Figure 1. This shows the effectiveness of the online RL training. An additional benefit of our methodology is the fact that the decrease in token usage is controllable using α , however its unclear how to do it using SFT or DPO.

Secondly, *increasing the value of α results in a greater decrease in response length* as seen in Figure 6. This is expected since the magnitude of the reward difference between shortest and longest correct responses increases as α increases. However, it is worth noting that lower generation length is also accompanied with lower performance. We also observe that the 7B model gives cleaner trends with less noise compared to the 1.5B model. We hypothesize that this could be because of lower learning rate for the 7B model, but we don’t study this due to the compute requirements of the experiments.

Thirdly, we also observe that the decrease in response length is *greater for easier problems* compared to harder problems. For instance, in Figure 3, we observe that in the 7B model, $\alpha = 0.2$ brings a token saving of 27% on AIME2024 and of 83% on GSM8K. We hypothesize that this is because there is a larger relative ‘spread’ in easier problems as compared to harder problems, that is, reasoning models are more wasteful on easier problems (see Appendix I for a more detailed study). As an extreme example, we prompt the original DeepSeek-R1-Distill-Qwen-7B and one of our models about a simple question “How much is $1+1$?”. While DeepSeek-R1-Distill-Qwen-7B reasoning model expends several tokens (more than a page in Appendix B) to arrive at the correct solution, the model trained with our method reaches the same conclusion within a few tokens.

Finally, even without any length penalty (i.e., $\alpha = 0$), we observe a reduction in response length on both the MATH and AIME datasets. We find that this is because of the recently discovered *biased nature* of RLOO’s implementation Liu et al. (2025b) in popular open-source libraries such as OpenRLHF which normalizes the loss by the length of the response. We find that removing the length bias from the implementation leads to no change in the response length of the model after

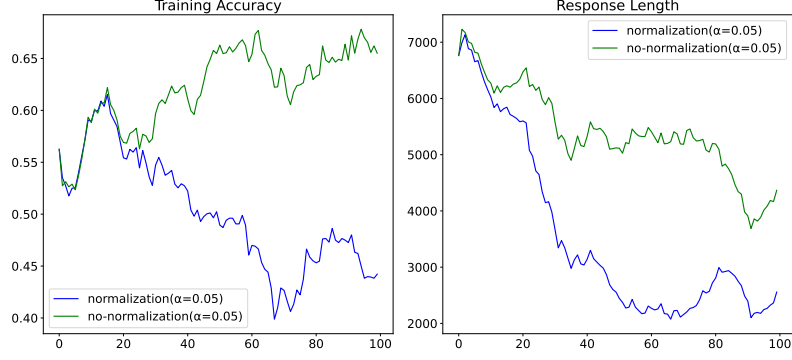


Figure 4: Advantage normalization rapidly decreases the response length alongside accuracy.

training. Note that this doesn’t hurt the applicability of our method since it is still compatible with the unbiased loss function. See Appendix Section G for more details.

5.4 Ablations

We perform an ablation to study a critical design component in the implementation of our method, namely the decision of not normalizing the advantage function in the RL training procedure.

In fact, it is a common practice (e.g., GRPO Shao et al. (2024)) to normalize the token-level advantage function and obtain $\hat{A}_{i,t} = \frac{r_i - r_{mean}}{r_{std}}$ where r_{mean} is the mean reward and r_{std} is the standard deviation of the rewards. While this choice is sensible in a more standard setting, it can have unintended consequences when the objective function contains the length penalty.

Consider the case where for a prompt x , all responses are correct. In that case, all rewards will be distributed within $[1 - \alpha, 1]$. Assume that the reward distribution is uniformly distributed in $[1 - \alpha, 1]$. In that case, the mean reward is $1 - \frac{\alpha}{2}$ and the standard deviation is $\frac{\alpha}{\sqrt{12}}$. The normalized advantage value for a correct response with maximum value $r = 1$ (i.e., the shortest correct response) becomes $\frac{1 - (1 - \alpha/2)}{\frac{\alpha}{\sqrt{12}}} = \sqrt{3}$ which is independent of α ! In other words, the advantage normalization, can bring a length decrease independent of α . Figure 4 shows that the resulting length decrease is generally too substantial for the model to absorb, leading to a sharp drop in accuracy.

6 Limitations

Our optimization procedure, while effective, is somewhat more involved than SFT or DPO-derived techniques because of the reinforcement learning setup. Furthermore, the choice of the penalty coefficient α affects the overall generation cost but does not precisely target a precise generation length, which may be required by some latency-constrained applications. This precise problem has been tackled by Aggarwal & Welleck (2025). Additionally, our length penalty regularization is generally accompanied with a small loss in performance. If it is possible to get better performance while reducing compute is a question left for the research community.

7 Conclusion

In this work, we introduced a novel and simple methodology that significantly reduces the inference cost for reasoning models while minimally affecting its accuracy. Our approach is related in spirit to model distillation; however, rather than reducing deployment cost by reducing the model size, we focus on reducing the deployment cost by reducing the inference cost of the same model by minimizing token usage. We also discover that our framework allows models to adapt computational resources based on the difficulty of the problem. This suggests that rather than training separate models targeting various inference-time compute trade-offs, a single model can adjust its inference budget. This property holds promise for applications requiring scalable, cost-effective AI solutions that are highly efficient without compromising on accuracy.

8 Acknowledgements

This work used DeltaAI at NCSA through allocation CIS240294: Advancing Large Language Model Alignment through Hybrid Feedback Integration and CIS250018: Training LLMs to Reasoning for Longer from the Advanced Cyberinfrastructure Coordination Ecosystem: Services & Support (ACCESS) program (Boerner et al., 2023), which is supported by U.S. National Science Foundation grants #2138259, #2138286, #2138307, #2137603, and #2138296.

References

- Aggarwal, P. and Welleck, S. L1: Controlling how long a reasoning model thinks with reinforcement learning, 2025. URL <https://arxiv.org/abs/2503.04697>.
- Ahmadian, A., Cremer, C., Gallé, M., Fadaee, M., Kreutzer, J., Pietquin, O., Üstün, A., and Hooker, S. Back to basics: Revisiting reinforce style optimization for learning from human feedback in llms, 2024. URL <https://arxiv.org/abs/2402.14740>.
- Besta, M., Blach, N., Kubicek, A., Gerstenberger, R., Podstawski, M., Gianinazzi, L., Gajda, J., Lehmann, T., Niewiadomski, H., Nyczyk, P., et al. Graph of thoughts: Solving elaborate problems with large language models. In *Proceedings of the AAAI Conference on Artificial Intelligence*, volume 38, pp. 17682–17690, 2024.
- Boerner, T. J., Deems, S., Furlani, T. R., Knuth, S. L., and Towns, J. Access: Advancing innovation: Nsf’s advanced cyberinfrastructure coordination ecosystem: Services & support. In *Practice and Experience in Advanced Research Computing 2023: Computing for the Common Good*, PEARC ’23, pp. 173–176, New York, NY, USA, 2023. Association for Computing Machinery. ISBN 9781450399852. doi: 10.1145/3569951.3597559. URL <https://doi.org/10.1145/3569951.3597559>.
- Chen, X., Xu, J., Liang, T., He, Z., Pang, J., Yu, D., Song, L., Liu, Q., Zhou, M., Zhang, Z., Wang, R., Tu, Z., Mi, H., and Yu, D. Do not think that much for $2+3=?$ on the overthinking of o1-like llms, 2024. URL <https://arxiv.org/abs/2412.21187>.
- Cobbe, K., Kosaraju, V., Bavarian, M., Chen, M., Jun, H., Kaiser, L., Plappert, M., Tworek, J., Hilton, J., Nakano, R., Hesse, C., and Schulman, J. Training verifiers to solve math word problems, 2021a. URL <https://arxiv.org/abs/2110.14168>.
- Cobbe, K., Kosaraju, V., Bavarian, M., Chen, M., Jun, H., Kaiser, L., Plappert, M., Tworek, J., Hilton, J., Nakano, R., et al. Training verifiers to solve math word problems. *arXiv preprint arXiv:2110.14168*, 2021b.
- Gandhi, K., Lee, D., Grand, G., Liu, M., Cheng, W., Sharma, A., and Goodman, N. D. Stream of search (sos): Learning to search in language. *arXiv preprint arXiv:2404.03683*, 2024.
- Guo, D., Yang, D., Zhang, H., Song, J., Zhang, R., Xu, R., Zhu, Q., Ma, S., Wang, P., Bi, X., et al. Deepseek-r1: Incentivizing reasoning capability in llms via reinforcement learning. *arXiv preprint arXiv:2501.12948*, 2025.
- Han, T., Wang, Z., Fang, C., Zhao, S., Ma, S., and Chen, Z. Token-budget-aware llm reasoning, 2024. URL <https://arxiv.org/abs/2412.18547>.
- Hendrycks, D., Burns, C., Kadavath, S., Arora, A., Basart, S., Tang, E., Song, D., and Steinhardt, J. Measuring mathematical problem solving with the math dataset. *arXiv preprint arXiv:2103.03874*, 2021.
- Hu, J., Wu, X., Zhu, Z., Xianyu, Wang, W., Zhang, D., and Cao, Y. Openrlhf: An easy-to-use, scalable and high-performance rlhf framework, 2024. URL <https://arxiv.org/abs/2405.11143>.
- Jin, M., Yu, Q., Shu, D., Zhao, H., Hua, W., Meng, Y., Zhang, Y., and Du, M. The impact of reasoning step length on large language models. In Ku, L.-W., Martins, A., and Srikumar, V. (eds.), *Findings of the Association for Computational Linguistics: ACL 2024*, pp. 1830–1842, Bangkok, Thailand, August 2024. Association for Computational Linguistics. doi: 10.18653/v1/2024.findings-acl.108. URL <https://aclanthology.org/2024.findings-acl.108/>.

- Kang, Y., Sun, X., Chen, L., and Zou, W. C3ot: Generating shorter chain-of-thought without compromising effectiveness, 2024. URL <https://arxiv.org/abs/2412.11664>.
- Kaplan, J., McCandlish, S., Henighan, T., Brown, T. B., Chess, B., Child, R., Gray, S., Radford, A., Wu, J., and Amodei, D. Scaling laws for neural language models. *arXiv preprint arXiv:2001.08361*, 2020.
- Kingma, D. P. and Ba, J. Adam: A method for stochastic optimization, 2017. URL <https://arxiv.org/abs/1412.6980>.
- Kool, W., van Hoof, H., and Welling, M. Buy 4 REINFORCE samples, get a baseline for free!, 2019. URL <https://openreview.net/forum?id=r1lgTGL5DE>.
- Kumar, A., Zhuang, V., Agarwal, R., Su, Y., Co-Reyes, J. D., Singh, A., Baumli, K., Iqbal, S., Bishop, C., Roelofs, R., et al. Training language models to self-correct via reinforcement learning. *arXiv preprint arXiv:2409.12917*, 2024.
- Kwon, W., Li, Z., Zhuang, S., Sheng, Y., Zheng, L., Yu, C. H., Gonzalez, J., Zhang, H., and Stoica, I. Efficient memory management for large language model serving with pagedattention. In *Proceedings of the 29th Symposium on Operating Systems Principles*, pp. 611–626, 2023a.
- Kwon, W., Li, Z., Zhuang, S., Sheng, Y., Zheng, L., Yu, C. H., Gonzalez, J. E., Zhang, H., and Stoica, I. Efficient memory management for large language model serving with pagedattention, 2023b. URL <https://arxiv.org/abs/2309.06180>.
- Leviathan, Y., Kalman, M., and Matias, Y. Fast inference from transformers via speculative decoding. In *International Conference on Machine Learning*, pp. 19274–19286. PMLR, 2023.
- LI, J., Beeching, E., Lewis Tunstall, B. L., Soletskyi, R., Huang, S. C., Kashif Rasul, L. Y., Albert Jiang, Z. S., Qin, Z., Dong, B., Zhou, L., Fleureau, Y., Lample, G., and Polu, S. Numinamath. [<https://github.com/project-numina/aimo-progress-prize>] (https://github.com/project-numina/aimo-progress-prize/blob/main/report/numina_dataset.pdf), 2024.
- Lightman, H., Kosaraju, V., Burda, Y., Edwards, H., Baker, B., Lee, T., Leike, J., Schulman, J., Sutskever, I., and Cobbe, K. Let’s verify step by step. In *The Twelfth International Conference on Learning Representations*, 2024. URL <https://openreview.net/forum?id=v8L0pN6E0i>.
- Lin, J., Tang, J., Tang, H., Yang, S., Chen, W.-M., Wang, W.-C., Xiao, G., Dang, X., Gan, C., and Han, S. Awq: Activation-aware weight quantization for on-device llm compression and acceleration. *Proceedings of Machine Learning and Systems*, 6:87–100, 2024.
- Liu, R., Gao, J., Zhao, J., Zhang, K., Li, X., Qi, B., Ouyang, W., and Zhou, B. Can 1b llm surpass 405b llm? rethinking compute-optimal test-time scaling, 2025a. URL <https://arxiv.org/abs/2502.06703>.
- Liu, Z., Sun, M., Zhou, T., Huang, G., and Darrell, T. Rethinking the value of network pruning. *arXiv preprint arXiv:1810.05270*, 2018.
- Liu, Z., Chen, C., Li, W., Qi, P., Pang, T., Du, C., Lee, W. S., and Lin, M. Understanding rl-zero-like training: A critical perspective, 2025b. URL <https://arxiv.org/abs/2503.20783>.
- Luo, H., Shen, L., He, H., Wang, Y., Liu, S., Li, W., Tan, N., Cao, X., and Tao, D. O1-pruner: Length-harmonizing fine-tuning for o1-like reasoning pruning, 2025. URL <https://arxiv.org/abs/2501.12570>.
- Nayab, S., Rossolini, G., Simoni, M., Saracino, A., Buttazzo, G., Manes, N., and Giacomelli, F. Concise thoughts: Impact of output length on llm reasoning and cost, 2025. URL <https://arxiv.org/abs/2407.19825>.
- OpenAI, :, Jaech, A., Kalai, A., Lerer, A., Richardson, A., El-Kishky, A., Low, A., Helyar, A., Madry, A., Beutel, A., Carney, A., Iftimie, A., Karpenko, A., Passos, A. T., Neitz, A., Prokofiev, A., Wei, A., Tam, A., Bennett, A., Kumar, A., Saraiva, A., Vallone, A., Duberstein, A., Kondrich, A., Mishchenko, A., Applebaum, A., Jiang, A., Nair, A., Zoph, B., Ghorbani, B., Rossen, B.,

- Sokolowsky, B., Barak, B., McGrew, B., Minaiev, B., Hao, B., Baker, B., Houghton, B., McKinzie, B., Eastman, B., Lugaresi, C., Bassin, C., Hudson, C., Li, C. M., de Bourcy, C., Voss, C., Shen, C., Zhang, C., Koch, C., Orsinger, C., Hesse, C., Fischer, C., Chan, C., Roberts, D., Kappler, D., Levy, D., Selsam, D., Dohan, D., Farhi, D., Mely, D., Robinson, D., Tsipras, D., Li, D., Oprica, D., Freeman, E., Zhang, E., Wong, E., Proehl, E., Cheung, E., Mitchell, E., Wallace, E., Ritter, E., Mays, E., Wang, F., Such, F. P., Raso, F., Leoni, F., Tsimpourlas, F., Song, F., von Lohmann, F., Sulit, F., Salmon, G., Parascandolo, G., Chabot, G., Zhao, G., Brockman, G., Leclerc, G., Salman, H., Bao, H., Sheng, H., Andrin, H., Bagherinezhad, H., Ren, H., Lightman, H., Chung, H. W., Kivlichan, I., O’Connell, I., Osband, I., Gilaberte, I. C., Akkaya, I., Kostrikov, I., Sutskever, I., Kofman, I., Pachocki, J., Lennon, J., Wei, J., Harb, J., Twore, J., Feng, J., Yu, J., Weng, J., Tang, J., Yu, J., Candela, J. Q., Palermo, J., Parish, J., Heidecke, J., Hallman, J., Rizzo, J., Gordon, J., Uesato, J., Ward, J., Huizinga, J., Wang, J., Chen, K., Xiao, K., Singhal, K., Nguyen, K., Cobbe, K., Shi, K., Wood, K., Rimbach, K., Gu-Lemberg, K., Liu, K., Lu, K., Stone, K., Yu, K., Ahmad, L., Yang, L., Liu, L., Maksin, L., Ho, L., Fedus, L., Weng, L., Li, L., McCallum, L., Held, L., Kuhn, L., Kondraciuk, L., Kaiser, L., Metz, L., Boyd, M., Trebacz, M., Joglekar, M., Chen, M., Tintor, M., Meyer, M., Jones, M., Kaufer, M., Schwarzer, M., Shah, M., Yatbaz, M., Guan, M. Y., Xu, M., Yan, M., Glaese, M., Chen, M., Lampe, M., Malek, M., Wang, M., Fradin, M., McClay, M., Pavlov, M., Wang, M., Wang, M., Murati, M., Bavarian, M., Rohaninejad, M., McAleese, N., Chowdhury, N., Chowdhury, N., Ryder, N., Tezak, N., Brown, N., Nachum, O., Boiko, O., Murk, O., Watkins, O., Chao, P., Ashbourne, P., Izmailov, P., Zhokhov, P., Dias, R., Arora, R., Lin, R., Lopes, R. G., Gaon, R., Miyara, R., Leike, R., Hwang, R., Garg, R., Brown, R., James, R., Shu, R., Cheu, R., Greene, R., Jain, S., Altman, S., Toizer, S., Toyer, S., Miserendino, S., Agarwal, S., Hernandez, S., Baker, S., McKinney, S., Yan, S., Zhao, S., Hu, S., Santurkar, S., Chaudhuri, S. R., Zhang, S., Fu, S., Papay, S., Lin, S., Balaji, S., Sanjeev, S., Sidor, S., Broda, T., Clark, A., Wang, T., Gordon, T., Sanders, T., Patwardhan, T., Sottiaux, T., Degry, T., Dimson, T., Zheng, T., Garipov, T., Stasi, T., Bansal, T., Creech, T., Peterson, T., Eloundou, T., Qi, V., Kosaraju, V., Monaco, V., Pong, V., Fomenko, V., Zheng, W., Zhou, W., McCabe, W., Zaremba, W., Dubois, Y., Lu, Y., Chen, Y., Cha, Y., Bai, Y., He, Y., Zhang, Y., Wang, Y., Shao, Z., and Li, Z. Openai o1 system card, 2024. URL <https://arxiv.org/abs/2412.16720>.
- Qu, Y., Yang, M. Y. R., Setlur, A., Tunstall, L., Beeching, E. E., Salakhutdinov, R., and Kumar, A. Optimizing test-time compute via meta reinforcement fine-tuning, 2025. URL <https://arxiv.org/abs/2503.07572>.
- Rafailov, R., Sharma, A., Mitchell, E., Manning, C. D., Ermon, S., and Finn, C. Direct preference optimization: Your language model is secretly a reward model. In *Thirty-seventh Conference on Neural Information Processing Systems*, 2023. URL <https://openreview.net/forum?id=HPuSIXJaa9>.
- Rajbhandari, S., Rasley, J., Ruwase, O., and He, Y. Zero: Memory optimizations toward training trillion parameter models, 2020. URL <https://arxiv.org/abs/1910.02054>.
- Schulman, J., Wolski, F., Dhariwal, P., Radford, A., and Klimov, O. Proximal policy optimization algorithms, 2017. URL <https://arxiv.org/abs/1707.06347>.
- Shao, Z., Wang, P., Zhu, Q., Xu, R., Song, J., Bi, X., Zhang, H., Zhang, M., Li, Y. K., Wu, Y., and Guo, D. Deepseekmath: Pushing the limits of mathematical reasoning in open language models, 2024. URL <https://arxiv.org/abs/2402.03300>.
- Silver, D., Schrittwieser, J., Simonyan, K., Antonoglou, I., Huang, A., Guez, A., Hubert, T., Baker, L., Lai, M., Bolton, A., et al. Mastering the game of go without human knowledge. *nature*, 550 (7676):354–359, 2017.
- Snell, C. V., Lee, J., Xu, K., and Kumar, A. Scaling LLM test-time compute optimally can be more effective than scaling parameters for reasoning. In *The Thirteenth International Conference on Learning Representations*, 2025. URL <https://openreview.net/forum?id=4FWAwZtd2n>.
- Srivastava, A., Rastogi, A., Rao, A., Shoeb, A. A. M., Abid, A., Fisch, A., Brown, A. R., Santoro, A., Gupta, A., Garriga-Alonso, A., Kluska, A., Lewkowycz, A., Agarwal, A., Power, A., Ray, A., Warstadt, A., Kocurek, A. W., Safaya, A., Tazarv, A., Xiang, A., Parrish, A., Nie, A., Hussain, A., Askell, A., Dsouza, A., Slone, A., Rahane, A., Iyer, A. S., Andreassen, A. J., Madotto, A., Santilli, A., Stuhlmüller, A., Dai, A. M., La, A., Lampinen, A. K., Zou, A., Jiang, A., Chen, A.,

Vuong, A., Gupta, A., Gottardi, A., Norelli, A., Venkatesh, A., Gholamidavoodi, A., Tabassum, A., Menezes, A., Kirubarajan, A., Mullokandov, A., Sabharwal, A., Herrick, A., Efrat, A., Erdem, A., Karakaş, A., Roberts, B. R., Loe, B. S., Zoph, B., Bojanowski, B., Özyurt, B., Hedayatnia, B., Neyshabur, B., Inden, B., Stein, B., Ekmekci, B., Lin, B. Y., Howald, B., Orinion, B., Diao, C., Dour, C., Stinson, C., Argueta, C., Ferri, C., Singh, C., Rathkopf, C., Meng, C., Baral, C., Wu, C., Callison-Burch, C., Waites, C., Voigt, C., Manning, C. D., Potts, C., Ramirez, C., Rivera, C. E., Siro, C., Raffel, C., Ashcraft, C., Garbacea, C., Sileo, D., Garrette, D., Hendrycks, D., Kilman, D., Roth, D., Freeman, C. D., Khashabi, D., Levy, D., González, D. M., Perszyk, D., Hernandez, D., Chen, D., Ippolito, D., Gilboa, D., Dohan, D., Drakard, D., Jurgens, D., Datta, D., Ganguli, D., Emelin, D., Kleyko, D., Yuret, D., Chen, D., Tam, D., Hupkes, D., Misra, D., Buzan, D., Mollo, D. C., Yang, D., Lee, D.-H., Schrader, D., Shutova, E., Cubuk, E. D., Segal, E., Hagerman, E., Barnes, E., Donoway, E., Pavlick, E., Rodolà, E., Lam, E., Chu, E., Tang, E., Erdem, E., Chang, E., Chi, E. A., Dyer, E., Jerzak, E., Kim, E., Manyasi, E. E., Zheltonozhskii, E., Xia, F., Siar, F., Martínez-Plumed, F., Happé, F., Chollet, F., Rong, F., Mishra, G., Winata, G. I., de Melo, G., Kruszewski, G., Parascandolo, G., Mariani, G., Wang, G. X., Jaimovitch-Lopez, G., Betz, G., Gur-Ari, G., Galijasevic, H., Kim, H., Rashkin, H., Hajishirzi, H., Mehta, H., Bogar, H., Shevlin, H. F. A., Schuetze, H., Yakura, H., Zhang, H., Wong, H. M., Ng, I., Noble, I., Jumelet, J., Geissinger, J., Kernion, J., Hilton, J., Lee, J., Fisac, J. F., Simon, J. B., Koppel, J., Zheng, J., Zou, J., Kocon, J., Thompson, J., Wingfield, J., Kaplan, J., Radom, J., Sohl-Dickstein, J., Phang, J., Wei, J., Yosinski, J., Novikova, J., Bosscher, J., Marsh, J., Kim, J., Taal, J., Engel, J., Alabi, J., Xu, J., Song, J., Tang, J., Waweru, J., Burden, J., Miller, J., Balis, J. U., Batchelder, J., Berant, J., Froberg, J., Rozen, J., Hernandez-Orallo, J., Boudeman, J., Guerr, J., Jones, J., Tenenbaum, J. B., Rule, J. S., Chua, J., Kanclerz, K., Livescu, K., Krauth, K., Gopalakrishnan, K., Ignatyeva, K., Markert, K., Dhole, K., Gimpel, K., Omondi, K., Mathewson, K. W., Chiafullo, K., Shkaruta, K., Shridhar, K., McDonell, K., Richardson, K., Reynolds, L., Gao, L., Zhang, L., Dugan, L., Qin, L., Contreras-Ochando, L., Morency, L.-P., Moschella, L., Lam, L., Noble, L., Schmidt, L., He, L., Oliveros-Colón, L., Metz, L., Senel, L. K., Bosma, M., Sap, M., Hoeve, M. T., Farooqi, M., Faruqui, M., Mazeika, M., Baturan, M., Marelli, M., Maru, M., Ramirez-Quintana, M. J., Tolkiehn, M., Giulianelli, M., Lewis, M., Potthast, M., Leavitt, M. L., Hagen, M., Schubert, M., Baitemirova, M. O., Arnaud, M., McElrath, M., Yee, M. A., Cohen, M., Gu, M., Ivanitskiy, M., Starritt, M., Strube, M., Śwędrowski, M., Bevilacqua, M., Yasunaga, M., Kale, M., Cain, M., Xu, M., Suzgun, M., Walker, M., Tiwari, M., Bansal, M., Aminnaseri, M., Geva, M., Gheini, M., T. M. V., Peng, N., Chi, N. A., Lee, N., Krakover, N. G.-A., Cameron, N., Roberts, N., Doiron, N., Martinez, N., Nangia, N., Deckers, N., Muennighoff, N., Keskar, N. S., Iyer, N. S., Constant, N., Fiedel, N., Wen, N., Zhang, O., Agha, O., Elbaghdadi, O., Levy, O., Evans, O., Casares, P. A. M., Doshi, P., Fung, P., Liang, P. P., Vicol, P., Alipoormolabashi, P., Liao, P., Liang, P., Chang, P. W., Eckersley, P., Htut, P. M., Hwang, P., Miłkowski, P., Patil, P., Pezeshkpour, P., Oli, P., Mei, Q., Lyu, Q., Chen, Q., Banjade, R., Rudolph, R. E., Gabriel, R., Habacker, R., Risco, R., Millièrre, R., Garg, R., Barnes, R., Saurous, R. A., Arakawa, R., Raymaekers, R., Frank, R., Sikand, R., Novak, R., Sitelew, R., Bras, R. L., Liu, R., Jacobs, R., Zhang, R., Salakhutdinov, R., Chi, R. A., Lee, S. R., Stovall, R., Teehan, R., Yang, R., Singh, S., Mohammad, S. M., Anand, S., Dillavou, S., Shleifer, S., Wiseman, S., Gruetter, S., Bowman, S. R., Schoenholz, S. S., Han, S., Kwatra, S., Rous, S. A., Ghazarian, S., Ghosh, S., Casey, S., Bischoff, S., Gehrmann, S., Schuster, S., Sadeghi, S., Hamdan, S., Zhou, S., Srivastava, S., Shi, S., Singh, S., Asaadi, S., Gu, S. S., Pachchigar, S., Toshniwal, S., Upadhyay, S., Debnath, S. S., Shakeri, S., Thormeyer, S., Melzi, S., Reddy, S., Makini, S. P., Lee, S.-H., Torene, S., Hatwar, S., Dehaene, S., Divic, S., Ermon, S., Biderman, S., Lin, S., Prasad, S., Piantadosi, S., Shieber, S., Misherghi, S., Kiritchenko, S., Mishra, S., Linzen, T., Schuster, T., Li, T., Yu, T., Ali, T., Hashimoto, T., Wu, T.-L., Desbordes, T., Rothschild, T., Phan, T., Wang, T., Nkinyili, T., Schick, T., Kornev, T., Tunduny, T., Gerstenberg, T., Chang, T., Neeraj, T., Khot, T., Shultz, T., Shaham, U., Misra, V., Demberg, V., Nyamai, V., Raunak, V., Ramasesh, V. V., vinay uday prabhu, Padmakumar, V., Srikumar, V., Fedus, W., Saunders, W., Zhang, W., Vossen, W., Ren, X., Tong, X., Zhao, X., Wu, X., Shen, X., Yaghoobzadeh, Y., Lakretz, Y., Song, Y., Bahri, Y., Choi, Y., Yang, Y., Hao, S., Chen, Y., Belinkov, Y., Hou, Y., Hou, Y., Bai, Y., Seid, Z., Zhao, Z., Wang, Z., Wang, Z. J., Wang, Z., and Wu, Z. Beyond the imitation game: Quantifying and extrapolating the capabilities of language models. *Transactions on Machine Learning Research*, 2023. ISSN 2835-8856. URL <https://openreview.net/forum?id=uyTL5Bvosj>. Featured Certification.

- Team, K., Du, A., Gao, B., Xing, B., Jiang, C., Chen, C., Li, C., Xiao, C., Du, C., Liao, C., Tang, C., Wang, C., Zhang, D., Yuan, E., Lu, E., Tang, F., Sung, F., Wei, G., Lai, G., Guo, H., Zhu, H., Ding, H., Hu, H., Yang, H., Zhang, H., Yao, H., Zhao, H., Lu, H., Li, H., Yu, H., Gao, H., Zheng, H., Yuan, H., Chen, J., Guo, J., Su, J., Wang, J., Zhao, J., Zhang, J., Liu, J., Yan, J., Wu, J., Shi, L., Ye, L., Yu, L., Dong, M., Zhang, N., Ma, N., Pan, Q., Gong, Q., Liu, S., Ma, S., Wei, S., Cao, S., Huang, S., Jiang, T., Gao, W., Xiong, W., He, W., Huang, W., Wu, W., He, W., Wei, X., Jia, X., Wu, X., Xu, X., Zu, X., Zhou, X., Pan, X., Charles, Y., Li, Y., Hu, Y., Liu, Y., Chen, Y., Wang, Y., Liu, Y., Qin, Y., Liu, Y., Yang, Y., Bao, Y., Du, Y., Wu, Y., Wang, Y., Zhou, Z., Wang, Z., Li, Z., Zhu, Z., Zhang, Z., Wang, Z., Yang, Z., Huang, Z., Huang, Z., Xu, Z., and Yang, Z. Kimi k1.5: Scaling reinforcement learning with llms, 2025. URL <https://arxiv.org/abs/2501.12599>.
- Vaswani, A. Attention is all you need. *Advances in Neural Information Processing Systems*, 2017.
- Wainwright, M. J. *High-dimensional statistics: A non-asymptotic viewpoint*, volume 48. Cambridge university press, 2019.
- Wang, X., Wei, J., Schuurmans, D., Le, Q., Chi, E., Narang, S., Chowdhery, A., and Zhou, D. Self-consistency improves chain of thought reasoning in language models. *arXiv preprint arXiv:2203.11171*, 2022.
- Wei, J., Wang, X., Schuurmans, D., Bosma, M., Xia, F., Chi, E., Le, Q. V., Zhou, D., et al. Chain-of-thought prompting elicits reasoning in large language models. *Advances in neural information processing systems*, 35:24824–24837, 2022.
- Welleck, S., Lu, X., West, P., Brahman, F., Shen, T., Khashabi, D., and Choi, Y. Generating sequences by learning to self-correct. In *The Eleventh International Conference on Learning Representations*, 2023. URL <https://openreview.net/forum?id=hH36JeQZDa0>.
- Xia, H., Li, Y., Leong, C. T., Wang, W., and Li, W. Tokenskip: Controllable chain-of-thought compression in llms, 2025. URL <https://arxiv.org/abs/2502.12067>.
- Yao, S., Yu, D., Zhao, J., Shafran, I., Griffiths, T., Cao, Y., and Narasimhan, K. Tree of thoughts: Deliberate problem solving with large language models. *Advances in Neural Information Processing Systems*, 36, 2024.
- Zhou, Z., Ning, X., Hong, K., Fu, T., Xu, J., Li, S., Lou, Y., Wang, L., Yuan, Z., Li, X., Yan, S., Dai, G., Zhang, X.-P., Dong, Y., and Wang, Y. A survey on efficient inference for large language models, 2024. URL <https://arxiv.org/abs/2404.14294>.

A Implementation Details

We build on the OpenRLHF codebase Hu et al. (2024). For the 1.5B model, we use 4 GH200 GPUs on one low-density node and for the 7B model, we use 8 GH200 GPUs distributed across two low-density nodes (4 GPUs per node). We set vLLM to the maximum context length (32K) during generation and set the generation temperature to 1. For training the 1.5B, ZeRO Stage 2 Rajbhandari et al. (2020) is used and for the 7B, ZeRO Stage 3 with activation checkpointing is required to prevent out of memory errors. The training precision is set to bfloat16. We generate 8 responses for each prompt. For every iteration, 32 prompts are selected from the dataset and the global batch size is set to 128 which leads to 2 gradient steps per RL iteration. For the 1.5B, the learning rate is set to $5 \cdot 10^{-6}$ and for the 7B, it is set to $2 \cdot 10^{-6}$. For all experiments, Adam Kingma & Ba (2017) is used as the standard optimizer. We experiment with 4 values of α in the following range: 0.05, 0.1, 0.2 and 0.4. For all RL experiments, the value of the KL coefficient is set to $1 \cdot 10^{-3}$. The experiments on both model take approximately 20 hours. We use the same prompt template for all models which can be found in Appendix D.

B Qualitative examples of the model before and after training

We prompt the Distilled 7B model and our model trained with $\alpha = 0.4$ on a simple prompt.

Prompt: Please reason step by step, and put your final answer within \boxed{ }. Question: What is 1+1?

Model (before training):

```
<think> Okay, so I need to figure out what 1 plus 1 is. Hmm, let me think about this. I remember in school, they taught us that adding numbers is like combining groups of things. So, if I have one apple and someone gives me another apple, how many apples do I have? I guess that's two apples. Wait, let me try to visualize it. Imagine I have one block here and another block there. If I put them together, how many blocks do I have in total? That's right, two blocks. So, 1 plus 1 equals 2. But hold on, is this always the case? Like, in different number systems or something? I mean, I know in binary, 1 plus 1 is 10, but I think that's a different topic. The question is about regular addition, right? So, in standard arithmetic, 1 plus 1 is definitely 2. I can also think about the number line. If I start at 1 and move one step forward, where do I land? That's 2. So, that makes sense too. Maybe I can use my fingers to check. If I hold up one finger and then hold up another finger, how many fingers am I showing? Two fingers. Yep, that confirms it. Is there a more mathematical way to look at this? I remember something about addition being the simplest arithmetic operation. It's the basis for more complex calculations. So, 1 plus 1 is the first step in understanding addition. I also recall that in set theory, the sum of two sets with one element each results in a set with two elements. So, if set A has one element and set B has one element, then the union of A and B has two elements. That's another way to see that 1 plus 1 equals 2. Wait, but in some programming languages, like Python, what happens if I write 1 + 1? Let me think. If I open a Python interpreter, and type 1 + 1, it should output 2. Yeah, that's correct. So, even in code, it's 2. Is there any context where 1 plus 1 doesn't equal 2? Like in different mathematical structures or something? Hmm, in Boolean algebra, sometimes 1 plus 1 is 1, but that's not standard arithmetic anymore. It's more like a logical OR operation. But the original question is about regular addition, so that's not applicable here. Also, in some contexts like in a group where the operation is defined differently, 1 + 1 might not be 2, but again, that's beyond the scope of basic arithmetic. I guess another way to think about it is using the Peano axioms, which are the foundation of arithmetic. According to these axioms, 1 is a natural number, and the successor of 1 is 2. So, adding 1 to 1 gives you the successor of 1, which is 2. That's
```

pretty solid. Maybe I can also use real-life examples. If I have a glass with one liter of water and I pour another liter into it, how much water do I have? Two liters. So, that's another practical demonstration that $1 + 1 = 2$. I wonder if there's a geometric interpretation. If I have a line segment of length 1 and another line segment of length 1, and I place them end to end, the total length would be 2. Yep, that works too. So, after thinking through all these different angles-counting, number lines, sets, programming, Boolean algebra, axioms, real-life examples, and geometry-it all consistently points to the fact that $1 + 1 = 2$. There doesn't seem to be any situation within the realm of standard arithmetic where this isn't true. I guess another way to reinforce this is by using objects I have at hand. Like, if I have one book and someone gives me another book, how many books do I have? Two books. Or if I have one toy car and find another toy car, now I have two toy cars. This practical application really helps solidify the concept. I also remember learning about the commutative property of addition, which states that the order of numbers being added doesn't change the result. So, $1 + 1$ is the same as $1 + 1$, which is still 2. And the associative property, which allows us to add more than two numbers together without changing the result. Although, in this case, we're only adding two numbers, so that might not apply directly, but it's good to know for more complex problems. Additionally, thinking about money could help. If I have one dollar and I receive another dollar, how much money do I have? Two dollars. That's another real-world example supporting the idea that $1 + 1 = 2$. I suppose another way to view it is through tally marks. If I make one tally mark and then make another tally mark, how many marks do I have? Two tally marks. So, that's another visual method to confirm it. I also recall that in some early number systems, like Roman numerals, 1 is represented by 'I', and adding another 'I' gives you 'II', which stands for two. So, even in different numeral systems, the concept holds. Thinking about the history of mathematics, the concept of $1 + 1 = 2$ was established long ago. Euclid defined numbers and their operations in his "Elements," laying the foundation for arithmetic as we know it today. So, this basic principle has been consistent throughout history. In conclusion, after considering all these various perspectives-counting, number lines, sets, programming, Boolean algebra, axioms, real-life examples, geometry, and even the history of math-it's clear that $1 + 1 = 2$ in standard arithmetic. There's no ambiguity or alternative result in these different contexts. So, I can confidently say that the answer is 2. $1 + 1 = 2$. This is consistently true across various mathematical contexts, from basic arithmetic to more complex structures. 2

Model (after training, $\alpha = 0.4$): $\langle \text{think} \rangle$ To solve $(1 + 1)$, I start by identifying the two numbers involved, which are both 1. Next, I add these two numbers together. Finally, the sum of 1 and 1 is 2. $\langle \text{think} \rangle$ To solve $1 + 1$, follow these simple steps: 1. ****Identify the numbers to add:**** Both numbers are 1. 2. ****Perform the addition:**** $1 + 1 = 2$ Therefore, the final answer is 2.

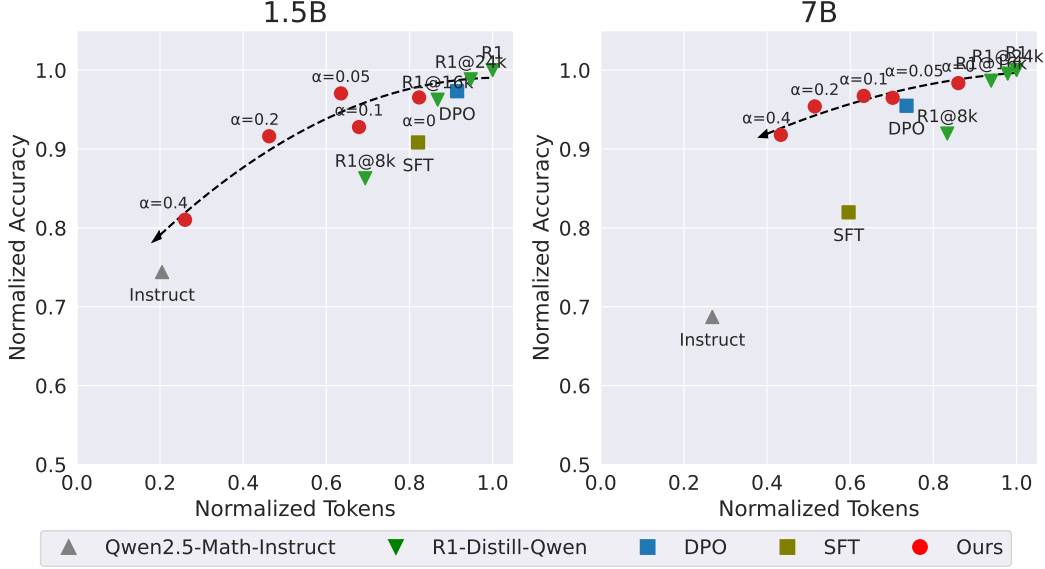


Figure 5: This figures describes the results of our training on 5 evaluation benchmarks. The green triangles represent the vLLM-cutoff baseline at different values of the token cutoff limit. As we can observe, the vLLM-cutoff baseline performs poorly, because we can get higher performing models with lower compute requirements.

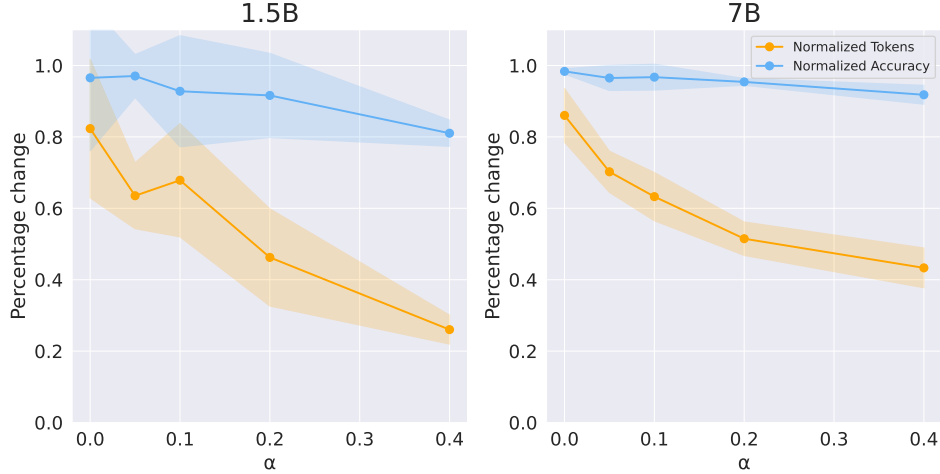


Figure 6: This plot demonstrates the change in normalized accuracy and normalized token usage as a function of the α parameter for the 1.5B and 7B models after training. We use three seeds to create the error bars and aggregate over 5 datasets. As evident, the token usage significantly reduces after training with little drop in accuracy.

C Extended Results

Here, we provide detailed results in tabular and graphical format for all the models and baselines. Firstly, we present the vLLM-cutoff baseline along with the other baselines in Figure 5. Secondly, Figure 6 contains a sensitivity plot of the model’s accuracy and token usage as α increases. The numbers are present in the numerical format for ease of future usage by researchers in Tables 1 to 10.

Next, we provide results on all datasets in Figures 7 to 11.

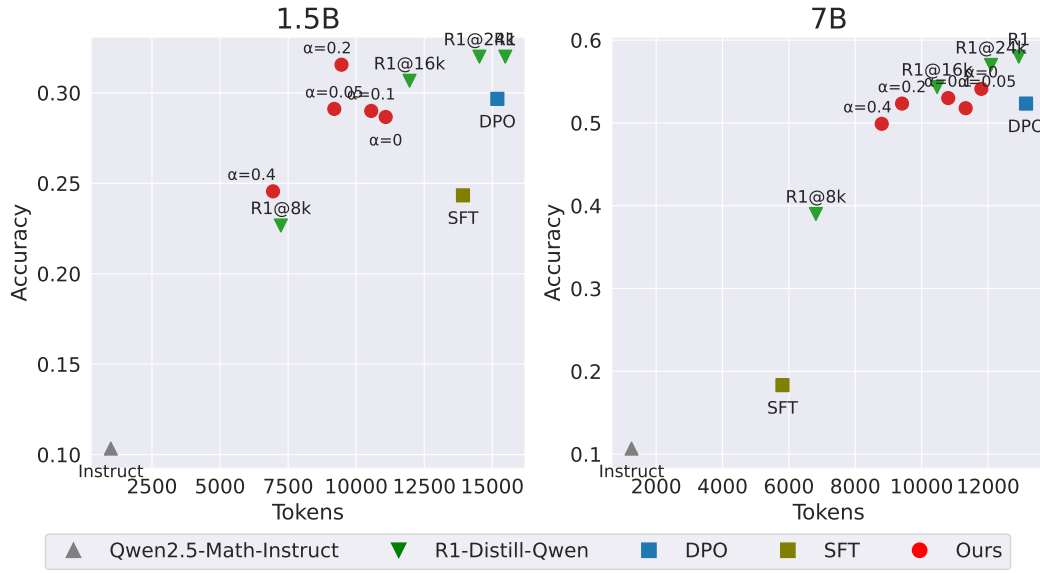


Figure 7: Detailed results on the AIME 2024 dataset.

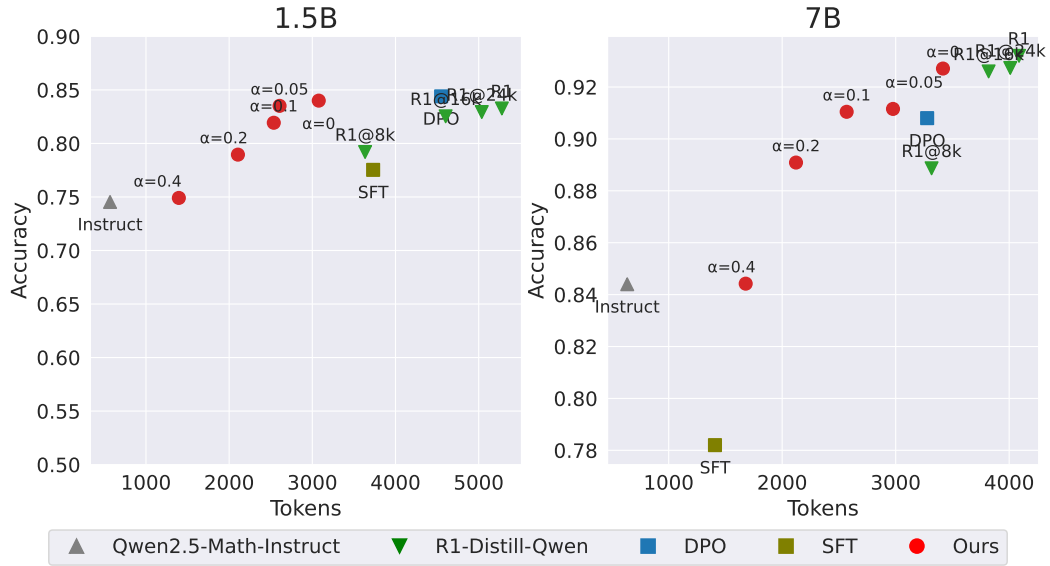


Figure 8: Detailed results on the MATH500 dataset.

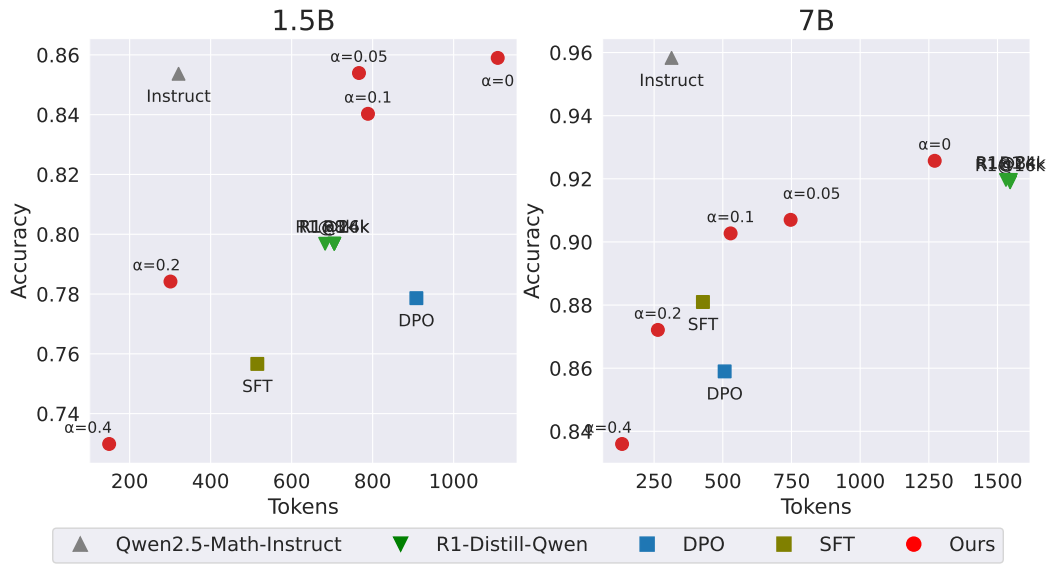


Figure 9: Detailed results on the GSM8K dataset.

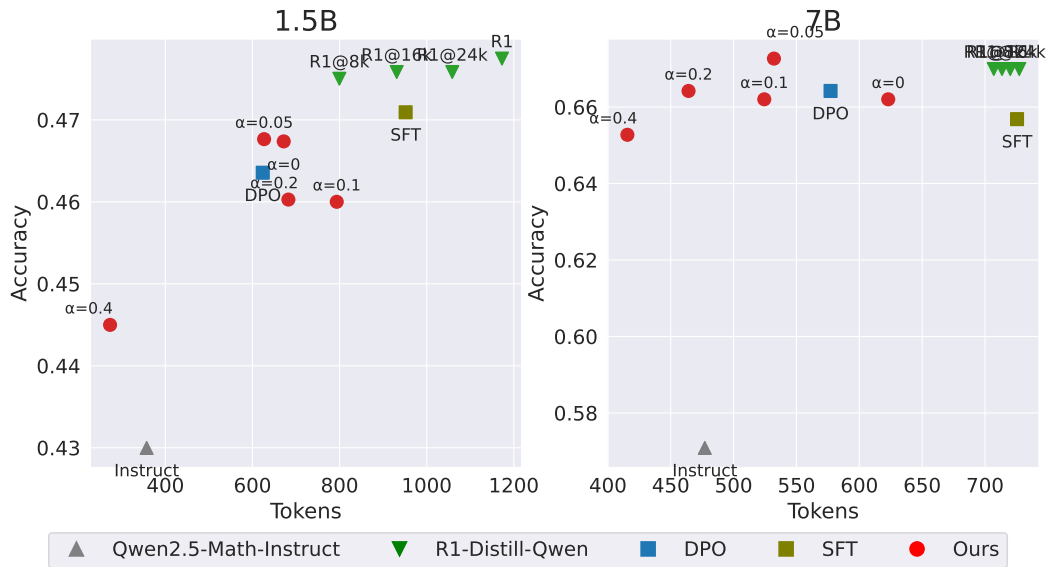


Figure 10: Detailed results on the CommonSenseQA dataset.

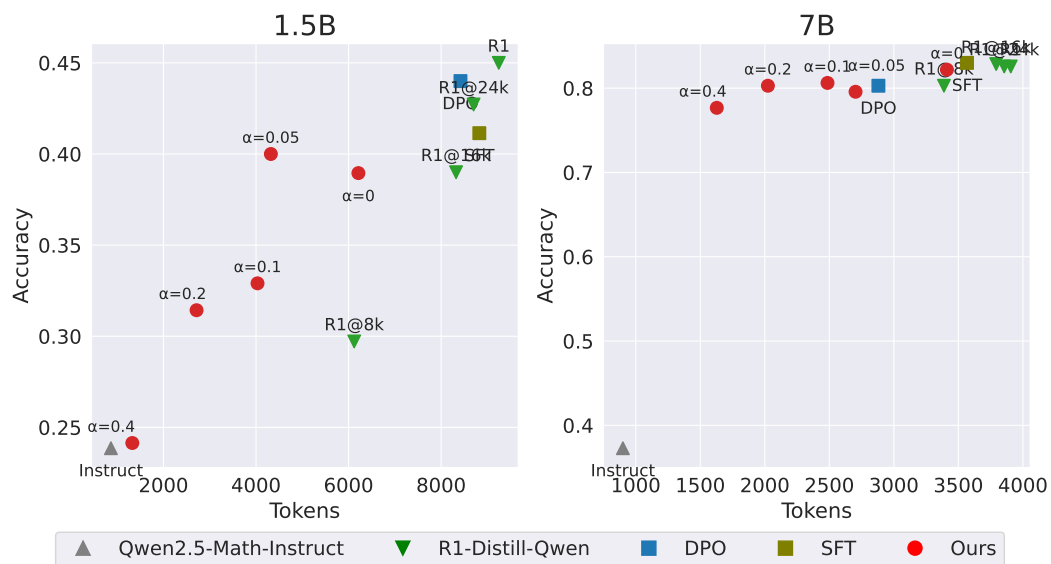


Figure 11: Detailed results on the Logical Deduction dataset.

Model Type	α	Tokens	Avg. Pass Rate	Ctx.
DPO	-	15184	0.297	32k
Distill	-	15476	0.320	32k
Distill	-	14525	0.320	24k
Distill	-	11960	0.307	16k
Distill	-	7232	0.227	8k
RL	0	11083 ± 2022.255	0.287 ± 0.007	32k
RL	0.05	9196 ± 897.023	0.291 ± 0.010	32k
RL	0.1	10552 ± 1219.273	0.290 ± 0.000	32k
RL	0.2	9459 ± 208.410	0.316 ± 0.037	32k
RL	0.4	6945 ± 1661.037	0.246 ± 0.006	32k
SFT	-	13922	0.243	32k

Table 1: AIME2024 results for 1.5B model

Model Type	α	Tokens	Avg. Pass Rate	Ctx.
DPO	-	4548	0.844	32k
Distill	-	5278	0.833	32k
Distill	-	5037	0.829	24k
Distill	-	4602	0.825	16k
Distill	-	3634	0.792	8k
RL	0	3077 ± 478.990	0.840 ± 0.004	32k
RL	0.05	2606 ± 189.452	0.835 ± 0.009	32k
RL	0.1	2536 ± 162.608	0.819 ± 0.010	32k
RL	0.2	2105 ± 203.024	0.790 ± 0.025	32k
RL	0.4	1395 ± 263.378	0.749 ± 0.009	32k
SFT	-	3731	0.775	32k

Table 2: MATH500 results for 1.5B model

Model Type	α	Tokens	Avg. Pass Rate	Ctx.
DPO	-	907	0.779	32k
Distill	-	705	0.797	32k
Distill	-	705	0.797	24k
Distill	-	705	0.797	16k
Distill	-	682	0.797	8k
RL	0	1108 ± 114.065	0.859 ± 0.008	32k
RL	0.05	766 ± 145.017	0.854 ± 0.012	32k
RL	0.1	788 ± 130.661	0.840 ± 0.010	32k
RL	0.2	300 ± 80.443	0.784 ± 0.023	32k
RL	0.4	149 ± 14.366	0.730 ± 0.036	32k
SFT	-	515	0.757	32k

Table 3: GSM8K results for 1.5B model

Model Type	α	Tokens	Avg. Pass Rate	Ctx.
DPO	-	623	0.464	32k
Distill	-	1172	0.477	32k
Distill	-	1058	0.476	24k
Distill	-	931	0.476	16k
Distill	-	799	0.475	8k
RL	0	671 ± 67.573	0.467 ± 0.012	32k
RL	0.05	626 ± 109.267	0.468 ± 0.013	32k
RL	0.1	793 ± 199.640	0.460 ± 0.009	32k
RL	0.2	682 ± 211.393	0.460 ± 0.007	32k
RL	0.4	273 ± 77.614	0.445 ± 0.015	32k
SFT	-	951	0.471	32k

Table 4: CommonSenseQA results for 1.5B model

Model Type	α	Tokens	Avg. Pass Rate	Ctx.
DPO	-	8418	0.440	32k
Distill	-	9246	0.450	32k
Distill	-	8701	0.427	24k
Distill	-	8322	0.390	16k
Distill	-	6121	0.297	8k
RL	0	6212 ± 1793.041	0.390 ± 0.092	32k
RL	0.05	4320 ± 851.996	0.400 ± 0.027	32k
RL	0.1	4032 ± 1464.899	0.329 ± 0.070	32k
RL	0.2	2715 ± 1257.025	0.314 ± 0.053	32k
RL	0.4	1326 ± 373.231	0.241 ± 0.016	32k
SFT	-	8825	0.411	32k

Table 5: Logical Deduction results for 1.5B model

Model Type	α	Tokens	Avg. Pass Rate	Ctx.
DPO	-	13147	0.523	32k
Distill	-	12837	0.563	32k
Distill	-	12091	0.570	24k
Distill	-	10467	0.543	16k
Distill	-	6814	0.390	8k
RL	0	11800 ± 593.175	0.541 ± 0.026	32k
RL	0.05	11328 ± 578.224	0.518 ± 0.014	32k
RL	0.1	10802 ± 371.321	0.530 ± 0.019	32k
RL	0.2	9410 ± 84.672	0.523 ± 0.010	32k
RL	0.4	8798 ± 811.499	0.499 ± 0.032	32k
SFT	-	5806	0.183	32k

Table 6: AIME2024 results for 7B model

Model Type	α	Tokens	Avg. Pass Rate	Ctx.
DPO	-	3277	0.908	32k
Distill	-	4086	0.932	32k
Distill	-	4009	0.927	24k
Distill	-	3818	0.926	16k
Distill	-	3316	0.889	8k
RL	0	3416 ± 192.578	0.927 ± 0.005	32k
RL	0.05	2976 ± 143.409	0.912 ± 0.004	32k
RL	0.1	2568 ± 206.046	0.910 ± 0.010	32k
RL	0.2	2121 ± 121.491	0.891 ± 0.020	32k
RL	0.4	1678 ± 152.337	0.844 ± 0.020	32k
SFT	-	1407	0.782	32k

Table 7: MATH500 results for 7B model

Model Type	α	Tokens	Avg. Pass Rate	Ctx.
DPO	-	506	0.859	32k
Distill	-	1547	0.920	32k
Distill	-	1547	0.920	24k
Distill	-	1545	0.919	16k
Distill	-	1530	0.920	8k
RL	0	1271 ± 72.028	0.926 ± 0.002	32k
RL	0.05	746 ± 224.053	0.907 ± 0.011	32k
RL	0.1	528 ± 177.573	0.903 ± 0.008	32k
RL	0.2	263 ± 45.282	0.872 ± 0.018	32k
RL	0.4	133 ± 30.929	0.836 ± 0.007	32k
SFT	-	427	0.881	32k

Table 8: GSM8k results for 7B model

Model Type	α	Tokens	Avg. Pass Rate	Ctx.
DPO	-	577	0.664	32k
Distill	-	727	0.670	32k
Distill	-	720	0.670	24k
Distill	-	713	0.670	16k
Distill	-	706	0.670	8k
RL	0	623 ± 17.082	0.662 ± 0.012	32k
RL	0.05	532 ± 17.618	0.673 ± 0.014	32k
RL	0.1	524 ± 44.455	0.662 ± 0.011	32k
RL	0.2	464 ± 9.552	0.664 ± 0.004	32k
RL	0.4	415 ± 21.115	0.653 ± 0.007	32k
SFT	-	725	0.657	32k

Table 9: CommonSenseQA results for 7B model

Model Type	α	Tokens	Avg. Pass Rate	Ctx.
DPO	-	2880	0.803	32k
Distill	-	3905	0.826	32k
Distill	-	3855	0.826	24k
Distill	-	3792	0.829	16k
Distill	-	3388	0.803	8k
RL	0	3411 ± 292.207	0.822 ± 0.006	32k
RL	0.05	2703 ± 224.043	0.796 ± 0.029	32k
RL	0.1	2485 ± 262.447	0.806 ± 0.030	32k
RL	0.2	2023 ± 182.110	0.803 ± 0.008	32k
RL	0.4	1627 ± 216.580	0.777 ± 0.022	32k
SFT	-	3566	0.830	32k

Table 10: Logical Deduction results for 7B model

D Prompt template for training

For all training purposes, we use the following prompt template:

Please reason step by step, and put your final answer within `\boxed{\}`.
Question: \$QUESTION

E Omitted short proofs

E.1 Proof of Proposition 4.3

Proof. Notice that the objective function Equation (4) can be written as

$$\frac{1}{|\mathcal{X}|} \sum_{x \in \mathcal{X}} \frac{1}{N} \sum_{i \in [N]} p_{\theta}(y_i | x) g(y_i) \quad (14)$$

for a positive function $g(\cdot) > 0$. Consider the following lemma.

Lemma E.1. *For a given prompt x , if there exist a correct answer $y' = y^*$, then the the population maximizer p_{θ} of Equation (14) places no mass on the incorrect answers for that prompt, i.e.,*

$$p_{\theta}(y | x) = 0, \text{ if } y \neq y^*. \quad (15)$$

Proof. Suppose the above claim did not hold; in other words, suppose that for some incorrect answer $y \neq y^*$, we have that $p_{\theta}(y | x) > 0$ and that p_{θ} maximizes Equation (14). Then consider the distribution $p_{\theta'}$ defined as

$$p_{\theta'}(y | x) = 0, \text{ if } y \neq y^* \quad (16)$$

$$p_{\theta'}(y | x) \propto p_{\theta}(y | x), \text{ if } y = y^*. \quad (17)$$

It can be verified that such distribution increases the value of the objective function (14) because it places more mass on the positive terms, contradicting the optimality of p_{θ} . \square

Lemma E.1 can be applied to establish the following: if for prompt x there exists a correct answer $y = y^*$, then

$$p_{\theta_{eff}^*}(y | x) = 0, \text{ if } y \neq y^* \quad (18)$$

which implies that both $p_{\theta_{eff}^*}$ has its support on the correct answers only, proving the claim. \square

E.2 Proof of Proposition 4.4

Proof. Assume the contrapositive, that is $\exists x \in \mathcal{X}, \exists y'$ such that $y' = y^*(x)$

$$\mathbb{E}_{y \sim p_{\theta_{eff}^*}(x)} [\text{LEN}(y)] > \text{LEN}(y')$$

Then consider the modified distribution $p_{\theta'}$ which places all mass on the response y' , that is,

$$p_{\theta'}(y|x) = 1\{y = y'\}$$

then

$$\mathbb{E}_{y \sim p_{\theta_{eff}^*}(x)} [\text{LEN}(y)] > \text{LEN}(y') = \mathbb{E}_{y \sim p_{\theta'}(x)} [\text{LEN}(y)]$$

Since f is a monotonically increasing function bounded in $[0, 1]$ and $\alpha \in [0, 1]$

$$\mathbb{E}_{y \sim p_{\theta_{eff}^*}(x)} [1 - \alpha \cdot f(\text{LEN}(y))] < \mathbb{E}_{y \sim p_{\theta'}(x)} [1 - \alpha \cdot f(\text{LEN}(y))]$$

Since $\text{ACCURACY}(p_{\theta_{eff}^*}) = \text{ACCURACY}(p_{\theta'}) = 1$, we have that

$$\Rightarrow \mathbb{E}_{y \sim p_{\theta_{eff}^*}(x)} [1\{y = y^*(x)\}(1 - \alpha \cdot f(\text{LEN}(y)))] < \mathbb{E}_{y \sim p_{\theta'}(x)} [1\{y = y^*(x)\}(1 - \alpha \cdot f(\text{LEN}(y)))]$$

This would imply that θ_{eff}^* is not the population-level maximizer which is a contradiction. \square

F Do Distilled Models follow Length Constraints?

One simple way to improve efficiency of reasoning models would be to simply prompt with an instruction saying: ‘Respond in less than X tokens’. This would prevent very verbose responses from the model. However, we perform an experiment to check whether the distilled models even have such a capability. For the MATH500 test set, we prompt the model in the following manner:

Please think step by step and answer in less than X tokens.

Question: {question}

Answer:

We vary X in the range 256, 512, 1024, 2048, 4096 and measure the number of output tokens. The results for the 1.5B and 7B Distilled models are in Tables 11 and 12. The results clearly demonstrate that there is no correlation between the requested number of tokens and the number of tokens that are actually generated by the model. This leads us to believe that the Distilled models are not capable of following length constraints out-of-the-box.

Token Limit	Tokens Generated
256	4609.34
512	4915.71
768	5228.85
1024	4913.84
1280	5306.68
2048	5064.06
4096	5245.11

Table 11: This table lists the average number of tokens generated by Distilled-R1-Qwen-1.5B for varying token limits as mentioned in the prompt for the MATH500 test set.

Token Limit	Tokens Generated
256	3434.56
512	3587.05
768	3518.34
1024	3716.17
1280	3524.46
2048	3688.01
4096	3815.11

Table 12: This table lists the average number of tokens generated by Distilled-R1-Qwen-7B for varying token limits as mentioned in the prompt for the MATH500 test set.

G Length Reduction without Length Penalty?

Reduction in length when $\alpha = 0$ is an intriguing observation. Recent work by Liu et al. (2025b) pointed out a bias in the GRPO loss function: it averages per-token loss across entire sequences, which unintentionally favors shorter correct sequences over longer correct ones, and longer incorrect sequences over shorter incorrect ones. Another inadvertent issue that this creates is that it increases loss weightage on problems which are easier since they generally have shorter Chain-of-Thoughts. This may explain the unexpected reduction in reasoning length, even when $\alpha = 0$. We tested the fix proposed by Liu et al. (2025b) and observed that the length reduction disappears when the fix is applied. Table 13 below shows changes in normalized accuracy and token usage (relative to a baseline 7B distilled model). The average results highlight that the fix mitigates the unintended length bias.

Dataset	RLOO+Fix (Δ NT)	RLOO+Fix (Δ NA)	RLOO (Δ NT)	RLOO (Δ NA)	Baseline (NT)	Baseline (NA)
MATH500	2.3	-0.4	-17.4	-0.6	100	100
AIME2024	8	-3	-10.9	-3.6	100	100
GSM8k	-12.2	-3.37	-17.2	1.08	100	100
Average	-0.64	-2.25	-15.16	-1.04	100	100

Table 13: Table showing the effects of fixes proposed by Liu et al. (2025b) Δ NT refers to change in normalized tokens. Δ NA refers to change in normalized accuracy. All numbers are normalized based on the Baseline scores. All experiments have been conducted on the 7B Distilled model.

H Training dynamics

We present the performance on the training dataset in Figure 12. Notably, setting $\alpha = 0$ corresponds to applying RL without any length penalty. Increasing α results in a significant reduction in token usage—up to 50% compared to the initial model—while maintaining the same level of accuracy as at the beginning of RL training. Lower values of α improve performance while still reducing the number of tokens.

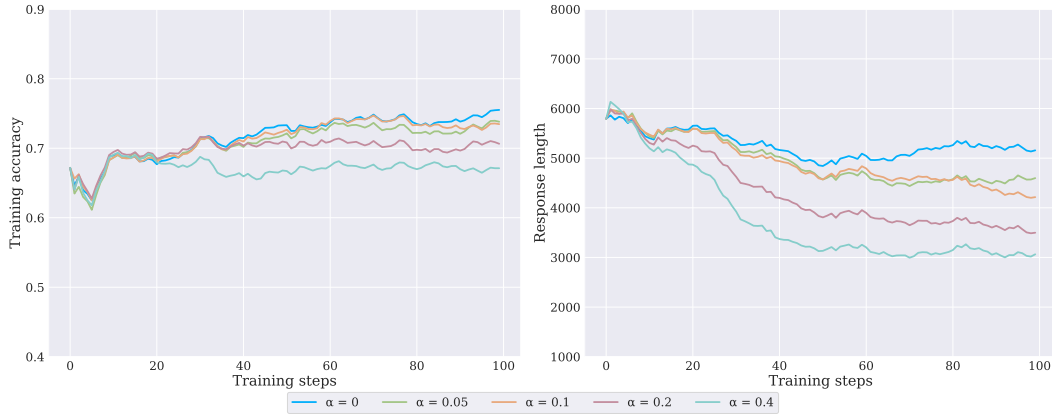


Figure 12: The figure shows the dynamics of the training accuracy and the corresponding generation lengths with varying values of α for the 7B model. The training accuracy and response length have been smoothed out using running averages over 25 training iterations.

I Difficulty based analysis

As discussed in the Results 5.3 earlier, we observe a bigger reduction in easier datasets like GSM8K as compared to harder datasets such as AIME2024. We hypothesize that this is because for easier datasets, reasoning models have a larger relative ‘spread’. To study this rigorously, we compute the quantities: $\frac{\text{STD}(x)}{\text{MEAN}(x)}$ and $\frac{\text{MEAN}(x) - \text{MIN}(x)}{\text{MEAN}(x)}$ for GSM8K, MATH500 and AIME2024 datasets for the Distilled 7B model. The values we get are reported in Table I. As we can see that the relative spread is much larger for the easier dataset (GSM8K) compared to the hardest dataset (AIME2024). This clearly demonstrates that reasoning models waste relatively more tokens thinking about easier problems compared to hard problems.

Dataset	$\frac{\text{STD}(x)}{\text{MEAN}(x)}$	$\frac{\text{MEAN}(x) - \text{MIN}(x)}{\text{MEAN}(x)}$
GSM8K	0.357	0.437
MATH500	0.271	0.323
AIME	0.254	0.264

Table 14: Relative spread for different datasets for the Distilled 7B model. We only aggregate over prompts where there is atleast 1 correct response out of 8 generations.