007 008

009

Training Language Models to Reason Efficiently

Anonymous Authors¹

Abstract

010Scaling model size and training data has led to011great advances in the performance of Large Lan-013guage Models. However, the diminishing returns014of this approach necessitate alternative methods to015improve model capabilities, particularly in tasks016requiring advanced reasoning. Large reasoning017models, which leverage long chain-of-thoughts,018bring unprecedented breakthroughs in problem-019solving capabilities but at a substantial deploy-020ment cost associated to longer generations. Re-021ducing inference costs is crucial for the economic022feasibility, user experience, and environmental023sustainability of these models.

In this work, we propose to train large reasoning models to reason efficiently. More precisely, we use reinforcement learning (RL) to train reasoning models to dynamically allocate inferencetime compute based on task complexity. Our method incentivizes models to minimize unnecessary computational overhead while maintaining accuracy, thereby achieving substantial 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, 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 reinforcement learning (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 commands 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 accurate 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 increases user experience. Finally, lowering the inference computation has a direct impact in reducing carbon emissions, with a positive benefit to both the environment and the society.

We aim to develop a procedure to *train* the model to use the appropriate amount of inference time compute to solve the problem at hand with reasoning. For straightforward problems, the resulting model would deliver efficient, direct solutions, while for more demanding tasks, it would invest additional computational effort to perform advanced reasoning. Such an adaptable model, that invests the minimum amount of compute to arrive at the correct solution, would represent a significant leap forward in terms of operational cost.

We use reinforcement learning policy gradient methods (Sutton & Barto, 2018) to train the model to use the least possible amount of tokens to reach the correct solution, thereby minimizing inference costs, ideally without compromising on accuracy. We achieve this goal by means of a modi-

¹Anonymous Institution, Anonymous City, Anonymous Region, Anonymous Country. Correspondence to: Anonymous Author <anon.email@domain.com>.

Preliminary work. Under review by the International Conference on Machine Learning (ICML). Do not distribute.

fied reinforcement learning formulation which encourages the model to produce correct answers with short chain-of-057 thoughts. To the best of our knowledge, we are among the 058 first to consider this problem, and we discuss concurrent 059 literature in Section 2. As a result, the model learns when 060 to stop thinking-rather than solving an easy math problem, 061 such as simple addition, through multiple approaches, it rec-062 ognizes when it has found the correct answer and concludes 063 its reasoning efficiently while maintaining accuracy.

064 In order to achieve this goal, we provide a reinforcement 065 learning implementation of the above principle which in-066 volves only a couple of line changes in a standard rein-067 forcement learning pipeline; this allows to directly leverage 068 existing RL codebases. Our method allows the user to con-069 trol the reduction in inference-time compute by adjusting a 070 scalar coefficient in an intuitive way. In other words, starting from a base reasoning model, our procedure allows to derive 072 a family of reasoning models, each with increased genera-073 tion efficiency (i.e., shorter chain-of-thoughts) compared to 074 the original reasoning model. 075

076 We conduct numerical experiments on two recently re-077 leased open-weight large reasoning models, DeepSeek-078 R1-Distill-Qwen-1.5B and DeepSeek-R1-Distill-Qwen-7B 079 (Guo et al., 2025) and derive models with a substantial 080 problem-dependent reduction in reasoning cost while ap-081 proximately maintaining accuracy. For the 7B model, our 082 method produces a model with a reduction of 16% tokens 083 on the competition-level benchmark American Invitational Mathematics Examination 2024 while slightly increasing 085 accuracy, and a reduction of 30% with a small reduction 086 in accuracy of 1% on the MATH dataset (Hendrycks et al., 087 2021), and a reduction of approximately 50% tokens on 088 GSM8K (Cobbe et al., 2021b) with similar accuracy, thereby 089 showing the ability of the model to dynamically reduce its 090 test-time compute budget with minimal loss in accuracy.

091 Beyond its simplicity, an attractive property of our approach 092 is its computational efficiency: although training reasoning 093 models with large scale reinforcement learning may have a 094 prohibitive cost (Guo et al., 2025), our procedure shows that 095 training them to reason efficiently is highly viable even with 096 modest academic resources: our models are obtained with 097 only 100 reinforcement learning steps (approximately 200 098 gradient updates). The fact that we achieve a performance 099 comparable to that of the original reasoning model with a 100 short training is surprising, because in few RL steps the model needs to optimize for reasoning patterns that are shorter and more efficient than the original model.

2. Related Work

104

105

106

108

109

Improving model capabilities with test-time compute Several techniques have been developed to enhance LLM reasoning through more test-time compute. Chain of thoughts (Wei et al., 2022) can be seen as one such fundamental method. Prompt engineering is a broadly applicable technique (White et al., 2023) which can be used to elicit specific abilities that are thought to be useful to reach a solution, such as thinking step by step, exploring multiple solution paths, and double-checking the answer. However, it does not scale because it does not train the model to use these strategies effectively. Self consistency (Wang et al., 2022) on the other hand is one of the most effective ways to enhance test-time performance when test-time verifiers are not available. The method generates multiple final answers and then returns the mode of their empirical distribution. As the mode of the empirical distribution converges to the mode of the population level distribution of the model answers, the method does not scale well with the number of samples, and moreover, it is only effective when the answers can be clustered together, such as in math problems. This limitation can be bypassed by Best-of-N, a simple, general purpose and effective search technique. It relies on sampling multiple responses from the model and then selecting the best at test time according to the scoring function; however, it critically relies on the availability of an accurate test-time scoring function (Gao et al., 2023). A more sophisticated search technique is Monte Carlo Tree Search, because it directs the compute budget to the most promising directions of the search space. It was a critical component to the development of AlphaGo (Silver et al., 2017). However the algorithm is not directly applicable outside of structured search frameworks. Tree-of-thoughts (Yao et al., 2024) and its extension (Gandhi et al., 2024; Besta et al., 2024) can be seen as implementing search in natural language but they are limited by their bespoke nature. Process reward models (Lightman et al., 2024) provide step by step numerical guidance on the progress of the chain-of-thought, but they have not been as effective to build large scale reasoning systems. Finally, self-correction (Kumar et al., 2024) trains the LLMs to fact-check itself; however, it implements a specific technique within a scripted framework rather than being a general purpose technique to enhance the reasoning capabilities with more test-time compute.

While the above mentioned techniques can be highly effective in specialized scenarios, modern large scale reasoning models, which we discuss next, are trained with reinforcement learning and rely on 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 and general purpose reasoning. They are trained with large scale reinforcement learning (Guo et al., 2025), which leads them to develop emerging abilities, such as branching, verification and backtracking. Our approach aims at making these models more efficient to deploy.

111

112 Efficient serving While we focus on developing reason-113 ing models that can be served efficiently, our approach is 114 orthogonal to existing methods from the literature of ef-115 ficient LLMs; see Zhou et al. (2024) for a recent survey. 116 For example, system-level techniques build a system to 117 accelerate inference. Some examples include speculative 118 decoding (Leviathan et al., 2023) and batch engines like 119 vLLM (Kwon et al., 2023a); both can be directly com-120 bined with our method. Model-based techniques, on the 121 other hand, act directly on the model to accelerate infer-122 ence. Some examples include weight pruning (Liu et al., 123 2018) and quantization (Lin et al., 2024), which can also 124 be combined with our methodology. In contrast, our ap-125 proach leverages reinforcement learning to train the model 126 for computational efficiency, making it applicable whenever 127 the chain of thought is not required in the final answer. 128

129 Concurrent works To our knowledge, the first open-130 weight LLM that can be classified as a 'reasoning' model-131 producing long monolithic chain of thoughts-is the 32 132 billion parameter model QwQ-preview, which was released 133 on November 28 on the Hugging Face. As these models are 134 very recent, we are not aware of prior studies on efficiently 135 training these models to reason efficiently except for some 136 concurrent work, which we review below. 137

Chen et al. (2024) investigate the overthinking phenomena 138 and propose methods to mitigate it by using heuristics such 139 as First-Correct Solutions (FCS) and Greedy Diverse So-140 lutions (GDS) to generate preference data which is then 141 used for offline policy optimization. However, this method 142 doesn't allow easily tuning the model to the user's compute 143 budget. The concurrent technical report of Kimi k1.5 (Team 144 et al., 2025) also reports a method to shorten the chain-of-145 thought using a length penalty in the reward function while doing online RL, a procedure similar in principle but not 147 identical to ours. We note that their procedure does not 148 appear to have a tunable parameter which allows to obtain a 149 family of models-each with varying trade-offs-as we do. 150

152 **3. Setup**

153 154 Let p be a language model. When provided with a prompt x, 155 the language model produces a response $y = (y^1, y^2, ..., y^t)$, 156 where y^i represents the i-th token in the response and t is 157 the total number of tokens in the response sequence. More 158 precisely, the generation is *auto-regressive*, meaning that 159 given the prompt x and the tokens $y^{\leq k} = (y^1, y^2, ..., y^k)$ 160 generated so far, the next token y^{k+1} is generated from the 161 conditional model

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

151

164

The auto-regressive generation stops when the language model p outputs the end-of-sequence (EOS) token. Therefore, if $y = (y^1, y^2, ..., y^t)$ 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 $y = (y^1, y^2, ..., y^t)$ 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 thoughts" 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 z denote the full output of the model $z \sim p(x)$.

Objective function and reinforcement learning We consider problems where the responses generated from an LLM can be evaluated by a scoring function $f(x, z) \mapsto \mathbb{R}$, 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)

$$f(x,z) = 1\{z = z^{\star}(x)\}$$
(2)

where $z^{\star}(x)$ is the correct answer to the math problem x. Since z is the full output of the model, including the chain of thought, the relation $z = z^{\star}(x)$ 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

$$\operatorname{ACCURACY}(p) = \mathbb{E}_{x \sim \rho} \mathbb{E}_{z \sim p(x)} \big[1\{z = z^{\star}\} \big]. \quad (3)$$

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 the prompt distribution.

4. Method

We aim to design a method that trains models to use the minimum 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

(1)



Figure 1. 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.

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.

191 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

$$\mathbb{E}\Big[1\{z=z^{\star}(x)\}(1-\alpha f(\operatorname{LEN}(z))\Big]$$
(4)

199 using a monotonic function f of the input and a tunable 200 parameter $\alpha \in [0, 1)$. The choice $\alpha = 0$ yields the rein-201 forcement learning objective (3); increasing α increases the 202 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

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

where

183

186

196

197

203

204

205

206

208 209 210

211

212

213 214

215

216

217

218

219

$$\begin{split} \mathrm{MEAN}(x) &= \mathop{\mathbb{E}}_{\substack{y \sim p(x), \\ \mathrm{s.t.} \ 1\{z=z^{\star}\}=1}} [\mathrm{LEN}(z)] \\ \mathrm{STD}(x) &= \sqrt{\mathop{\mathrm{Var}}_{\substack{z \sim p(x), \\ \mathrm{s.t.} \ 1\{z=z^{\star}\}=1}} [\mathrm{LEN}(z)] \end{split}$$

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 generations, and that correct responses, even if long, are always preferred to incorrect ones.

4.1. Optimizing the objective with Reinforcement Learning

Since optimizing Equation (4) involves sampling from the model auto-regressively, the objective function is nondifferentiable; 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

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

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

269

270

271

272

273

274

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* generation for prompt *x*. We then simply use the sequence level advantage as the token level advantage, namely $\mathcal{A}(y^{\leq t}, x) = \mathcal{A}(y, x)$. In essence, we use PPO with the RLOO advantage estimator.

5. Experiments

We seek to evaluate our method through numerical experiments. In particular, we aim to answer the following questions:

- What are the trade-offs between accuracy and inference cost?
- What are simple baselines that are relevant in this setting?

We first discuss the setup, then introduce some baselines, present the empirical results and associated trade-offs, and finally discuss some ablations.

5.1. Setup

Initial unsuccessful experiments In our initial experiments, we performed distillation from QwQ-32B-Preview to Qwen2.5-3B-Instruct and Qwen2.5-1.5B-Instruct so as to elicit strong reasoning skills in these two models; these distilled models would have served as a starting point for our method.

However, to our surprise, the distilled models showed a regression in performance on common benchmarks such as MATH and AIME 2024 compared to the instruct model, despite using much longer chain-of-thoughts with qualitatively more advanced reasoning patterns. Although our method is still effective in reducing the length of the chain-of-thought of the distilled model, these experiments do not accurately reflect the trade-off between inference-cost and accuracy when the instruct model is also taken into account.

Models and Datasets We revisited our method following the release of the reasoning models DeepSeek-R1-Distill-Qwen-1.5B and DeepSeek-R1-Distill-Qwen-7B (Guo et al., 2025). These models were distilled from the more powerful DeepSeek-R1 using industry-grade techniques. 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. Notably, they demonstrate impressive performance on challenging benchmarks such as AIME 2024.

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. We use the same dataset across all baselines to ensure consistency.

Evaluation We evaluate the models on three datasets, ordered by increasing difficulty:

- GSM8K (Cobbe et al., 2021a), which contains gradeschool-level math problems,
- MATH (Hendrycks et al., 2021) which is a standard benchmark containing harder problems than GSM8K,
- the The American Invitational Mathematics Examination (AIME) 2024, a competition-level set of challenging mathematical problems.

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 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

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

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



Figure 2. This figures describes the results of our training on the MATH500 test set. The green triangle in the left corner represents the desirable trend of models where higher accuracy is achieved with a lower number of tokens.

every iteration, 32 prompts are selected from the dataset 303 and the global batch size is set to 128 which leads to 2 304 gradient steps per RL iteration. For the 1.5B, the learning 305 rate is set to $5 \cdot 10^{-6}$ and for the 7B, it is set to $2 \cdot 10^{-6}$. 306 For all experiments, Adam (Kingma & Ba, 2017) is used as 307 the standard optimizer. We experiment with 4 values of α : 308 0.05, 0.1, 0.2 and 0.4. For all RL experiments, the value of 309 the KL coefficient is set to 10^{-3} . We use the same prompt 310 template for all models which can be found in Appendix C. 311

5.2. Baselines

298 299

300

301 302

312

313

314

315

316

317

318

319

327

328

329

To the best of our knowledge, apart from the concurrent work discussed in Section 2, no prior studies have explored this setting. Alongside our method, we introduce and implement simple baseline approaches that help balance inference cost and accuracy.

- 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, 20,000, 24,000, and 32,000.
 - 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. 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, whereas the other rates do not achieve any reduction.

5.3. Numerical 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, and report the results on MATH500 in Figures 2 and 3. The re-



Figure 3. Results on AIME2024 dataset

sults for GSM8K can be found in Appendix B due to space reason.

As shown in Figure 2, our method enables smooth tradeoffs 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 364 on the MATH dataset decreases by 30% (from 4000 to 2800 365 tokens) while the accuracy loss is only 1%. Similarly, in 366 the AIME dataset (Figure 3), setting $\alpha = 0.2$ reduces token 367 usage by 30% (from 14,000 to 9,000) while incurring only a 2% accuracy drop compared to the DeepSeek-R1-Distill-369 Qwen-7B. 370

We offer several remarks: 371

353 354

355

357

358 359

360

361

362

363

- We prompt the original DeepSeek-R1-Distill-Qwen-7B and one of our models about a simple question "How 374 much is 1+1?". While DeepSeek-R1-Distill-Qwen-7B 375 reasoning model expends several tokens (more than a 376 page in the appendix) to arrive at the correct solution, 377 the model trained with our method quickly reaches the 378 same conclusion within few tokens. 379
- 380 • The models trained with our procedures adapt the 381 length of the chain of thought to the difficulty of the 382 problem. For example, $\alpha = 0.2$ brings a token saving 383 of 22% on AIME24 and of 77% on GSM8K compared 384

to doing RL at $\alpha = 0$.

- Even without any length penalty (i.e., $\alpha = 0$), we observe a reduction in response length on both the MATH and AIME datasets. We hypothesize that this occurs because these models have not been previously trained with reinforcement learning (RL) and have only undergone a single round of distillation from R1.
- The SFT and DPO baselines appear to perform worse than early-stopping the vLLM generation.
- The model with the highest $\alpha = 0.4$ experiences a larger performance drop compared to the others. In Figure 4, we visualize the training dynamics by plotting accuracy every 10 RL iterations. The figure illustrates how $\alpha = 0.4$ induces a rapid reduction in response length, likely preventing the model from adapting effectively, ultimately leading to lower performance.

5.4. Ablations

We perform an ablation to study a highly 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 standard practice (e.g., (Shao et al., 2024)) to



Figure 4. Evolution of the model during training with varying values of α . Checkpoints are created after 10 iterations.

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 417 correct. In that case, all rewards will be distributed within 418 $[1 - \alpha, 1]$. Assume that the reward distribution is uniformly 419 distributed in $[1 - \alpha, 1]$. In that case, the mean reward 420 is $1 - \frac{\alpha}{2}$ and the standard deviation is $\frac{\alpha}{\sqrt{12}}$. The normal-421 422 ized advantage value for a correct response with maximum 423 value r = 1 (i.e., the shortest correct response) becomes $\frac{1-(1-\alpha/2)}{\alpha} = \sqrt{3}$ which is independent of α ! In other words, 424 425 the advantage normalization, under certain conditions, can 426 bring a length decrease independent of α . The resulting 427 length decrease is generally too substantial for the model 428 to absorb, and this leads to a sharp drop in accuracy during 429 training, as can be seen in Figure 5. . 430

6. Limitations

405

406

407

408

409

410

411

412

413

414

415

416

431

432

433

434

435

436

437

438

439

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



Figure 5. Advantage normalization rapidly decreases the response length alongside accuracy.

applications. We leave such exact controllability as future work.

7. Conclusion

In this work, we introduced a novel 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.

A key advantage of our framework is its ability to adapt computational resources based on problem difficulty. This suggests that rather than training separate models targeting various inference-time compute trade-offs, a single model can adjust its inference budget dynamically. This property holds promise for applications requiring scalable, cost-effective AI solutions that are highly efficient without compromising on accuracy.

8. Impact Statement

This paper presents work whose goal is to advance the field of Machine Learning. There are many potential societal consequences of our work, none which we feel must be specifically highlighted here.

References

440 441

442

443

444

445

446

447

448

449

450

451

452

453

466

- 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.
- 461 Chen, X., Xu, J., Liang, T., He, Z., Pang, J., Yu, D., Song,
 462 L., Liu, Q., Zhou, M., Zhang, Z., Wang, R., Tu, Z., Mi,
 463 H., and Yu, D. Do not think that much for 2+3=? on
 464 the overthinking of o1-like llms, 2024. URL https:
 465 //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.
- 482 Gao, L., Schulman, J., and Hilton, J. Scaling laws for reward
 483 model overoptimization. In *International Conference on*484 *Machine Learning*, pp. 10835–10866. PMLR, 2023.
- Guo, D., Yang, D., Zhang, H., Song, J., Zhang, R., Xu, R.,
 Zhu, Q., Ma, S., Wang, P., Bi, X., et al. Deepseek-r1: Incentivizing reasoning capability in llms via reinforcement learning. *arXiv preprint arXiv:2501.12948*, 2025.
- Hendrycks, D., Burns, C., Kadavath, S., Arora, A., Basart,
 S., Tang, E., Song, D., and Steinhardt, J. Measuring mathematical problem solving with the math dataset. *arXiv preprint arXiv:2103.03874*, 2021.

- Hu, J., Wu, X., Zhu, Z., 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.
- 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=rllgTGL5DE.
- 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, Z., Sun, M., Zhou, T., Huang, G., and Darrell, T. Rethinking the value of network pruning. *arXiv preprint arXiv:1810.05270*, 2018.
- 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 *Thirtyseventh Conference on Neural Information Processing Systems*, 2023. URL https://openreview.net/ forum?id=HPuSIXJaa9.
- 511 Rajbhandari, S., Rasley, J., Ruwase, O., and He, Y. Zero:
 512 Memory optimizations toward training trillion parame513 ter models, 2020. URL https://arxiv.org/abs/
 514 1910.02054.
- 515
 516
 517
 518
 518
 519
 519
 519
 510
 511
 512
 513
 514
 515
 514
 515
 516
 517
 517
 518
 519
 519
 518
 519
 519
 519
 510
 510
 510
 510
 510
 510
 511
 512
 514
 514
 515
 516
 517
 516
 517
 518
 519
 518
 519
 519
 519
 519
 510
 510
 510
 510
 510
 510
 510
 510
 510
 510
 510
 510
 510
 510
 510
 510
 510
 510
 510
 510
 510
 510
 510
 510
 510
 510
 510
 510
 510
 510
 510
 510
 510
 510
 510
 510
 510
 510
 510
 510
 510
 510
 510
 510
 510
 510
 510
 510
 510
 510
 510
 510
 510
 510
 510
 510
 510
 510
 510
 510
 510
 510
 510
 510
 510
 510
 510
 510
 510
 510
 510
 510
 510
 510
 510
 510
 510
 510
 510
 510
 510
 510
 510
 510
 510
 510
 510
 510
 510
 510
 510
 510
 510
 510
 510
 510
 510
 510
 510
 510
 510
 510
 510
 510
 510
 510
 510
 510
 510
- 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.
- Sutton, R. S. and Barto, A. G. *Reinforcement learning: An introduction*. MIT press, 2018.

- 533 Team, K., Du, A., Gao, B., Xing, B., Jiang, C., Chen, C., 534 Li, C., Xiao, C., Du, C., Liao, C., Tang, C., Wang, C., 535 Zhang, D., Yuan, E., Lu, E., Tang, F., Sung, F., Wei, G., 536 Lai, G., Guo, H., Zhu, H., Ding, H., Hu, H., Yang, H., 537 Zhang, H., Yao, H., Zhao, H., Lu, H., Li, H., Yu, H., 538 Gao, H., Zheng, H., Yuan, H., Chen, J., Guo, J., Su, J., 539 Wang, J., Zhao, J., Zhang, J., Liu, J., Yan, J., Wu, J., 540 Shi, L., Ye, L., Yu, L., Dong, M., Zhang, N., Ma, N., 541 Pan, Q., Gong, Q., Liu, S., Ma, S., Wei, S., Cao, S., 542 Huang, S., Jiang, T., Gao, W., Xiong, W., He, W., Huang, 543 W., Wu, W., He, W., Wei, X., Jia, X., Wu, X., Xu, X., 544 Zu, X., Zhou, X., Pan, X., Charles, Y., Li, Y., Hu, Y., 545 Liu, Y., Chen, Y., Wang, Y., Liu, Y., Qin, Y., Liu, Y., 546 Yang, Y., Bao, Y., Du, Y., Wu, Y., Wang, Y., Zhou, Z., 547 Wang, Z., Li, Z., Zhu, Z., Zhang, Z., Wang, Z., Yang, Z., 548
- 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.
- 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.
- White, J., Fu, Q., Hays, S., Sandborn, M., Olea, C., Gilbert, H., Elnashar, A., Spencer-Smith, J., and Schmidt, D. C. A prompt pattern catalog to enhance prompt engineering with chatgpt. arXiv preprint arXiv:2302.11382, 2023.
- 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. Qualitative examples of the model before 550 and after training 552

We prompt the Distilled 7B model and our model trained 553 554 with $\alpha = 0.4$ on a simple prompt.

555 Prompt: Please reason step by step, and put your final an-556 swer within boxed. Question: What is 1+1?

Model (before training): 558

551

557

559 <think> Okay, so I need to figure out 560 what 1 plus 1 is. Hmm, let me think 561 about this. I remember in school, they 562 taught us that adding numbers is like 563 combining groups of things. So, if 564 I have one apple and someone gives 565 me another apple, how many apples do 566 I have? I guess that's two apples. 567 Wait, let me try to visualize it. 568 Imagine I have one block here and 569 another block there. If I put them 570 together, how many blocks do I have 571 in total? That's right, two blocks. 572 So, 1 plus 1 equals 2. But hold on, 573 is this always the case? Like, in 574 different number systems or something? 575 I mean, I know in binary, 1 plus 1 576 is 10, but I think that's a different 577 topic. The question is about regular 578 addition, right? So, in standard 579 arithmetic, 1 plus 1 is definitely 580 2. I can also think about the number 581 line. If I start at 1 and move one 582 step forward, where do I land? That's 583 2. So, that makes sense too. Maybe I 584 can use my fingers to check. If I hold 585 up one finger and then hold up another 586 finger, how many fingers am I showing? 587 Two fingers. Yep, that confirms it. 588 Is there a more mathematical way to 589 look at this? I remember something 590 about addition being the simplest 591 arithmetic operation. It's the 592 basis for more complex calculations. 593 So, 1 plus 1 is the first step in 594 understanding addition. I also recall 595 that in set theory, the sum of two 596 sets with one element each results in 597 a set with two elements. So, if set 598 A has one element and set B has one 599 element, then the union of A and B has 600 two elements. That's another way to 601 see that 1 plus 1 equals 2. Wait, but 602 in some programming languages, like 603 Python, what happens if I write 1 + 604

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 Like in different mathematical 2? 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 plus 1 equals 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 plus 1 equals 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

605 about the commutative property of 606 addition, which states that the order 607 of numbers being added doesn't change 608 the result. So, 1 + 1 is the same 609 as 1 + 1, which is still 2. And the 610 associative property, which allows us 611 to add more than two numbers together 612 without changing the result. Although, 613 in this case, we're only adding two 614 numbers, so that might not apply 615 directly, but it's good to know for 616 more complex problems. Additionally, 617 thinking about money could help. If I 618 have one dollar and I receive another 619 dollar, how much money do I have? Two 620 dollars. That's another real-world 621 example supporting the idea that 1 plus 622 1 equals 2. I suppose another way to 623 view it is through tally marks. If 624 I make one tally mark and then make 625 another tally mark, how many marks do 626 I have? Two tally marks. So, that's 627 another visual method to confirm it. Ι 628 also recall that in some early number 629 systems, like Roman numerals, 1 is 630 represented by 'I', and adding another 'I' gives you 'II', which stands for 631 two. So, even in different numeral 632 633 systems, the concept holds. Thinking 634 about the history of mathematics, 635 the concept of 1 plus 1 equals 2 was 636 established long ago. Euclid defined 637 numbers and their operations in his 638 "Elements," laying the foundation 639 for arithmetic as we know it today. 640 So, this basic principle has been 641 consistent throughout history. In 642 conclusion, after considering all these 643 various perspectives | counting, number 644 lines, sets, programming, Boolean algebra, axioms, real-life examples, 645 646 geometry, and even the history of 647 math | it's clear that 1 plus 1 equals 648 2 in standard arithmetic. There's 649 no ambiguity or alternative result 650 in these different contexts. So, I 651 can confidently say that the answer is 652 </think> 1 plus 1 equals 2. 2. This 653 is consistently true across various 654 mathematical contexts, from basic 655 arithmetic to more complex structures. 656 |2|657 658

659

Model (after training, $\alpha = 0.4$): <think> 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. </think> 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.

B. Results on GSM8K

The results on the GSM8K test set are in Figure 6. Its interesting that the DeepSeek model performs worse than the Instruct model on the GSM8K test set. However, with our RL training, we are able to get it to the same performance. However, it is hard to make it more compute efficient than the Instruct model in this dataset.

C. 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 $\begin{tabular}{ll} boxed{}. \\ Question: QUESTION \end{array}$



Figure 6. This figures describes the results of our training on the GSM8k test set. The green triangle in the left corner represents the desirable trend of models.