ConciseRL: Conciseness-Guided Reinforcement Learning for Efficient Reasoning Models

Razvan-Gabriel Dumitru

University of Arizona ServiceNow AI

razvandumm@gmail.com

Vikas Yadav ServiceNow AI

Abstract

Large language models excel at complex tasks by breaking down problems into structured reasoning steps. However, reasoning traces often extend beyond reaching a correct answer, causing wasted computation, reduced readability, and hallucinations. To address this, we introduce a novel hyperparameter-free conciseness score used as a reward signal within a reinforcement learning framework to guide models toward generating correct and concise reasoning traces. This score is evaluated by a large language model acting as a judge, enabling dynamic, context-aware feedback beyond simple token length. Our method achieves state-ofthe-art efficiency-accuracy trade-offs on the MATH dataset, reducing token usage by up to 31× on simple problems while improving accuracy by 7%, and on the hardest problems, it outperforms full reasoning by +7.5% accuracy with up to $3.6 \times$ fewer tokens. On TheoremQA, our method improves accuracy by +2.2% using $12.5 \times$ fewer tokens. We also conduct ablation studies on the judge model, reward composition, and problem difficulty, showing that our method dynamically adapts reasoning length based on problem difficulty and benefits significantly from stronger judges. The code, model weights, and datasets are open-sourced at https://github.com/RazvanDu/ConciseRL.

1 Introduction

Large language models (LLMs) have recently made significant progress in solving complex multi-step tasks, such as mathematics, code generation, and symbolic reasoning. Reasoning models such as o1 (OpenAI, 2024b), DeepSeek-R1 (DeepSeek-AI et al., 2025), and S1 (Muennighoff et al., 2025) are trained to reason explicitly through intermediate steps, often using reinforcement learning (RL) to optimize for correctness and structural fidelity. While this explicit reasoning improves accuracy and interpretability, it also introduces a major challenge: reasoning traces tend to be excessively long

Darius Peteleaza

MultiversX
Lucian Blaga University of Sibiu
peteleaza.darius@gmail.com

Liangming Pan

University of Arizona

(Chen et al., 2025; Team et al., 2025). These models frequently continue reasoning well past the point where the correct answer is reached, resulting in wasted computation, degraded readability, and sometimes even contradictions or hallucinated steps. Addressing this problem is essential for reducing inference costs and improving the usability of LLM reasoning in real-world applications.

In this work, we propose a novel approach to reduce excessive reasoning in LLMs by teaching them to generate answers that are both correct and concise. Our method introduces a semantically informed **conciseness score** as a reward function for training reasoning models and an LLM as a judge to evaluate the conciseness of the reasoning trace. One key insight is that while accuracy can often be judged deterministically (e.g., matching answers or symbolic execution), conciseness is inherently subjective and better suited for LLM-based evaluation.

Previous work has explored LLM-as-a-judge setups for evaluating factual accuracy, style, and helpfulness (Li et al., 2025; Zheng et al., 2023), but existing methods for controlling reasoning length rely on token count or static heuristics, which fail to capture semantic efficiency. Our method goes beyond simple token count: concise traces are often short, but short traces are not always concise, therefore, the reward targets conciseness directly, and a shorter length appears only as a side effect. As illustrated in Figure 2, all three examples use the same number of tokens, yet the first is the most concise, receiving the highest reward. The second and third are progressively less concise, resulting in lower rewards. Unlike static length penalties, our LLM-based conciseness score is dynamic and context-aware, enabling better generalization across problem types and difficulty levels while improving explainability. Additionally, our method can be combined with correctness-based rewards, offering a flexible trade-off between conciseness and robustness.

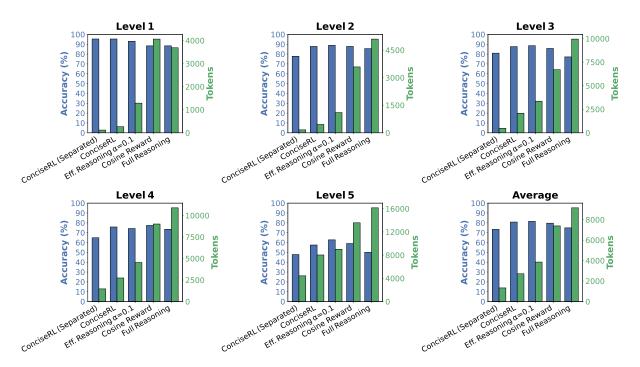


Figure 1: MATH500 histogram by difficulty level. We report both accuracy (blue, left axis) and average token length (green, right axis) for each method. All methods are based on DeepSeek-R1-Distill-Qwen-1.5B. For our method ("ConciseRL" and "ConciseRL (Separated)"), we use GPT-4.1 mini as the judge. The exact values shown in the histogram are reported in Table 4.

Our main contributions are as follows:

- We introduce a novel conciseness score and the first method that leverages an LLM to evaluate and provide a reward score based on the conciseness of reasoning traces. Unlike prior token-level heuristics that require careful hyperparameter tuning, our method captures the semantic compactness of reasoning in a context-sensitive and dynamic way, with no additional hyperparameters.
- Our method enables adaptive reasoning by adjusting response length based on problem difficulty. As shown in Figure 1 and Table 4, it allocates more tokens to harder MATH500 problems while using 31× fewer tokens on easier ones and improving accuracy by 7%. On the hardest questions, it outperforms full reasoning by 7.5% accuracy with 3.6× fewer tokens. On TheoremQA, our method achieves 2.2% higher accuracy using 12.5× fewer tokens than the full reasoning model. These results highlight a level of length-efficiency control not achieved by any static heuristic or prior method.
- We perform extensive experiments and ablations across multiple benchmarks, including

conciseness-only and accuracy-combined rewards, conditioning on problem difficulty, and analyzing how different judge models affect training dynamics, reward signals, and reasoning quality.

2 Related Work

Chain-of-Thought (CoT) (Wei et al., 2022) enhances LLM reasoning by prompting explicit intermediate steps, improving both performance and interpretability in complex tasks (Snell et al., 2025; Ouyang et al., 2022). Variants such as Tree-of-Thought (Yao et al., 2023) and Graph-of-Thought (Besta et al., 2024) extend the CoT paradigm through structured search and iterative refinement. Self-Consistency prompting (Wang et al., 2023) improves robustness by sampling and aggregating multiple reasoning paths but at the cost of increased computational overhead from longer traces.

Recent advances in reasoning models such as o1 (OpenAI, 2024b), o3 (OpenAI, 2025b), o4-mini (OpenAI, 2025b), DeepSeek-R1 (DeepSeek-AI et al., 2025), S1 (Muennighoff et al., 2025), and QwQ-32B (Team, 2025a) show strong reasoning through internal capabilities, without inference-time prompting. Some methods adopt tree-based search strategies (Yao et al., 2024), while others,

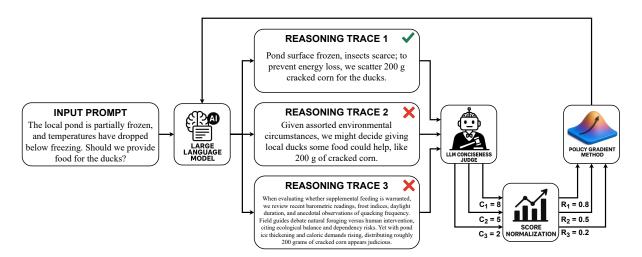


Figure 2: Given an input prompt, an LLM generates multiple reasoning traces that are evaluated by an LLM-based judge who scores each trace based on **conciseness**. Trace 1 is concise and receives the highest reward, Trace 2 has an equal length (24 tokens) but lower conciseness, while Trace 3 is the longest (71 tokens) and least concise. These rewards then guide a policy gradient update.

like DeepSeek-R1 (DeepSeek-AI et al., 2025; Shao et al., 2024), combine supervised fine-tuning with RL. Despite their performance, these methods are prone to generating verbose reasoning traces, a phenomenon known as the overthinking problem (Chen et al., 2025; Team et al., 2025).

Overthinking occurs when a model continues generating redundant reasoning steps after reaching the correct answer. This behavior wastes inference compute and risks introducing inconsistencies and logic errors (Chen et al., 2025). This issue is amplified by training objectives that reward long reasoning sequences (DeepSeek-AI et al., 2025), creating a misalignment between performance and efficiency.

To address these issues, several approaches have been proposed for CoT compression and length control. These include both reward-shaping strategies and optimization objectives designed to encourage shortness without compromising accuracy. Methods like O1-Pruner (Luo et al., 2025a) use Proximal Policy Optimization (PPO) (Schulman et al., 2017) with a length-harmonizing reward that penalizes deviations from a target length relative to a reference model. Another work (Yeo et al., 2025) applies a two-branch PPO strategy with cosine penalties based on the prediction's proximity to a maximum acceptable length. L1 (Aggarwal and Welleck, 2025), whose aim is an exact token budget, impractical when problem difficulty is unknown, employs Grouped Reinforcement Policy Optimization (GRPO) (Shao et al., 2024), explicitly conditioning the model to "think for N tokens"

and penalizing outputs exceeding learned optimal lengths. Simple Preference Optimization (SimPO)-based (Meng et al., 2024) approaches like DAST (Shen et al., 2025) use data constructed from user or synthetic feedback to learn token-length preferences, bypassing reliance on reference models. Other traditional (Han et al., 2025) and policy gradient (PG) methods (Arora and Zanette, 2025) have also been tried in this context, such as those including early fine-tuning with correctness-weighted penalties on output length.

Output length control, while previously a peripheral concern, has become central in reasoning LLMs (Sui et al., 2025; Fatemi et al., 2025; Liu et al., 2025; Zhu and Li, 2025). Earlier works used architectural tweaks or static training signals, which lacked the flexibility to handle diverse inference scenarios. However, recent RL-based models (Butcher et al., 2024; Jie et al., 2023; Hou et al., 2025) achieve better control by continuously adapting the reasoning length based on the learned reward signals. Yet many still rely on token count or heuristics as proxies for conciseness (Sui et al., 2025). Models such as C3oT (Kang et al., 2025) use LLMs as compressors to post-process reasoning traces, while continuous token methods like CCoT (Cheng and Durme, 2024) and latent-space approaches like COCONUT (Hao et al., 2024) reduce token-level redundancy by replacing explicit steps abstract representations. Although effective at shortening outputs, these approaches operate post hoc or rely on static heuristics that limit adaptability across tasks.

To address these challenges, we propose a novel approach that uses an LLM as a judge to evaluate the conciseness of reasoning traces. While LLMs have been widely used as evaluators for factuality, summarization, coherence, and style (Zheng et al., 2023; Li et al., 2025; Gu et al., 2025), leveraging them to assess reasoning conciseness remains unexplored. Our method introduces a semantically aware, dynamic reward function that goes beyond token count or static heuristics. This enables models to generate reasoning that is not only correct but also efficient and interpretable, bridging the gap between performance and inference cost.

3 Method

We build on pre-trained multi-step reasoning models, improving them to generate correct answers using the minimum semantic effort through more concise outputs. Instead of penalizing length, we reward conciseness: the ability to justify an answer with the fewest necessary and non-redundant steps. Concise traces are often short, but short traces are not always concise, therefore the reward targets conciseness directly, and a shorter length appears only as a side effect.

3.1 Conciseness and Accuracy Rewards

Given an input prompt x and a reasoning trace $y=(t_1,\ldots,t_{|y|})$ sampled from the model $p_{\theta}(\cdot\mid x)$, we compute a **conciseness score** $C(y)\in [0,1]$ by querying an external LLM judge \mathcal{J} , which receives the reasoning trace and evaluates it according to a custom system prompt (Appendix A.3). The system prompt instructs the judge to assign a discrete conciseness score s from 1 (overly verbose) to 10 (clear reasoning), disregarding correctness. The discrete score $s\in\{1,\ldots,10\}$ is then normalized as: $C(y)=\frac{s}{10}\in[0,1]$.

Accuracy Signal. Let $y^*(x)$ denote the ground-truth answer to prompt x, and let $\operatorname{Ans}(y)$ be the final answer extracted from the reasoning trace (e.g. \boxed{Ans}(y)). The accuracy is then:

$$A(y,x) = \begin{cases} 1, & \text{if } \operatorname{Ans}(y) = y^{\star}(x), \\ 0, & \text{otherwise.} \end{cases}$$
 (1)

3.2 Reward variants

We employ two reward formulations:

1. Pure conciseness

$$R_{\rm c}(y,x) = C(y). \tag{2}$$

2. Accuracy-gated conciseness (used to avoid judge calls when a trace is already wrong)

$$R_{\rm ac}(y,x) = A(y,x) \cdot C(y). \tag{3}$$

 $R_{\rm ac}$ is cheaper in terms of API calls because the judge is queried only when A(y,x)=1. On 1.5B models, $R_{\rm ac}$ costs < \$9 per training run, and there is virtually no increase in training time. The cost of prompting the model stays constant relative to the size of the training data. The judge is not used during inference, so our method has no additional cost or overhead during deployment. In the experiments, we denote the Accuracy-gated reward as "ConciseRL" and the Pure conciseness reward as "ConciseRL (Separated)". ConciseRL uses only the gated conciseness as a reward, and the Separated variant includes both conciseness and accuracy as rewards during training.

3.3 Optimization with PPO

Let ρ denote the distribution over prompts. Our objective is to maximize:

$$\mathcal{J}(\theta) = \mathbb{E}_{x \sim \rho} \, \mathbb{E}_{y \sim p_{\theta}(\cdot|x)} [R(y,x)]. \tag{4}$$

where R is a reward variant. Since expectation y is taken over sequences sampled autoregressively from p_{θ} , \mathcal{J} is not differentiable. Thus, we apply PPO (Shao et al., 2024) with a sequence-level objective.

3.4 Leave-One-Out Advantage

Based on (Arora and Zanette, 2025), for each prompt we sample n candidate traces $\{y_i\}_{i=1}^n$ and compute the trajectory return $R(y_i, x)$. The sequence-level advantage is estimated with a Reinforced Leave-One-Out baseline:

$$A(y_i, x) = R(y_i, x) - \frac{1}{n-1} \sum_{j \neq i} R(y_j, x).$$
 (5)

Using Equation (5), the clipped PPO loss is:

$$\mathcal{L}(\theta) = \mathbb{E}_{x,y} \Big[w_i \cdot \log p_{\theta}(y \mid x) \Big],$$
 (6)

$$w_i = \operatorname{clip}\left(\frac{p_{\theta}(y_i|x)}{p_{\theta,i,i}(y_i|x)}, 1-\epsilon, 1+\epsilon\right) A(y_i, x).$$
 (7)

where ϵ is the PPO clipping threshold and $p_{\theta_{\text{old}}}$ is the policy prior to the update.

4 Experiments and Results

4.1 Experimental Setup

We begin with two publicly available reasoning models as our foundation: DeepSeek-R1-Distill-Qwen-1.5B (DeepSeek-AI et al., 2025), and STILL-3-1.5B-preview (Team, 2025b; Min et al., 2024), conducting all ablations on the former with GPT-4.1 mini (OpenAI, 2025a) as the LLM judge. The judge's system prompt is shown in Figure 8. Our RL setup follows Efficient Reasoning (Arora and Zanette, 2025; LI et al., 2024), including a learning rate of 5×10^{-6} , Kullback–Leibler (KL) divergence coefficient of 10^{-3} , PPO clipping threshold of 0.2, and a maximum context window of 32K tokens, 8 rollouts per prompt, and a global batch of 128 (32 prompts). All plots are smoothed using a Gaussian kernel of width 9 for visual clarity. The Cosine Reward baseline uses hyperparameters from the Demystifying paper (Yeo et al., 2025): $L_{\text{max}}=14336$, $rc_0=2.0, rc_L=1.0, rw_0=-10.0, rw_L=0.0, and$ $r_{\text{exceed}} = -10.0$. The Efficient Reasoning baseline uses the same hyperparameters defined above, $\alpha \in 0.1, 0.2, 0.4, 0.6, 0.9$, and is implemented in vLLM (Kwon et al., 2023). We do not compare with methods that are not open-source or not reproducible. For instance, the method (Fatemi et al., 2025) is similar to ours, but their reward is based on shortening traces rather than conciseness, and their code is not publicly available. We also didn't include L1 (Aggarwal and Welleck, 2025) since it requires fixing the output length in advance, a fundamentally different setup that limits adaptability and can't learn dynamic conciseness.

Training is conducted on $4\times NVIDIA$ A100 80GB GPUs. A complete run on the 1.5B model takes ~ 20 GPU-hours. We evaluate model performance across multiple benchmarks: GSM8K (Cobbe et al., 2021), MATH500 (Hendrycks et al., 2021), TheoremQA (Chen et al., 2023), GPQAmain (Rein et al., 2024), and MMLU-Pro-1k (Wang et al., 2024).

4.2 Dynamic vs Static Rewards

We compare our method to strong baselines across five benchmarks. Baselines include Efficient Reasoning (Arora and Zanette, 2025), which introduces a length penalty modulated by a hyperparameter $\alpha \in [0,1)$. Higher values of α increase the penalty for longer generations, leading to shorter reasoning but compromising accuracy. We also compare

against Cosine Reward (Yeo et al., 2025), which assigns rewards using a cosine function that favors correct answers, modulating reward magnitude based on CoT length, and DeepScaleR (Luo et al., 2025b). The full reasoning baseline corresponds to the model's default generation without additional training to shorten the reasoning traces. The "Separated" baseline is missing from Table 2 due to substantial training costs given our budget constraints. Unlike these methods, our approach is hyperparameter-free, making it more robust and easier to deploy.

As shown in Table 1, we group methods by their average reasoning lengths and compare their accuracy to assess the efficiency-accuracy trade-offs. For our method, ConciseRL, the text in parentheses specifies the judge model used. Our strongest setting is ConciseRL (Gemini 2.5 Flash (Gemini Team, 2025)), which reaches 46.7% average accuracy at only 12.4% of the full reasoning length. On TheoremQA, it shortens traces to only 5.0% of the length and increases accuracy from 26.3% to 33.3%. Relative to Efficient Reasoning with α =0.9, ConciseRL (Gemini 2.5 Flash) achieves +23.5% in average accuracy while using fewer tokens, and it delivers large per-dataset gains such as +61.2% accuracy on GSM8K. The ConciseRL (Separated) variant remains highly compact at 16.0% length with 44.0% average accuracy, outperforming α =0.9 by +20.8% points on average at comparable length and reaching improvements of up to +56.4% on GSM8K and +35.0% on MATH500. In the mid-length comparison, ConciseRL (GPT-4.1 mini) achieves 49.3% average accuracy at 34.4% length and improves over Efficient Reasoning with α =0.6 by 5.8% points on average at a similar budget (34.6%), while also reducing tokens by 65.6% relative to Full Reasoning and increasing average accuracy by 4.0% points. With somewhat longer outputs, ConciseRL (Gemini 2.5 Flash-Lite) reaches the best average in its group at 49.7% with 44.1% length, surpassing Efficient Reasoning with α =0.2 by 6.0% points on average at similar length, and leading on GSM8K and MATH500. Among long-trace baselines, DeepScaleR achieves 52.4% average accuracy but at 117.7% length, and Cosine Reward expands traces further to 133.8% with 47.7% average accuracy.

This pattern also appears in Table 2, evaluating STILL-3-1.5B-preview. ConciseRL delivers the best overall balance, consistently maintaining strong average accuracy with low reasoning length.

	GSM8K		MATH500		MMLU-Pro-1k		GPQA-main		TheoremQA		Average	
Technique	Acc.	Len.(%)	Acc.	Len.(%)	Acc.	Len.(%)	Acc.	Len.(%)	Acc.	Len.(%)	Acc.	Len.(%)
ConciseRL (Gemini 2.5 Flash)	77.3	16.4	73.6	9.3	18.9	15.7	30.4	15.8	33.3	5.0	46.7	12.4
ConciseRL (Separated)	72.5	16.4	68.6	16.3	19.2	19.5	31.0	19.7	28.5	8.0	44.0	16.0
Eff. Reasoning α =0.9	16.1	1.2	33.6	21.0	18.2	11.6	27.9	26.2	20.4	9.7	23.2	13.9
ConciseRL (GPT-4.1 mini)	80.9	35.8	78.0	30.8	24.5	38.5	30.4	45.9	32.5	21.0	49.3	34.4
Eff. Reasoning α =0.6	55.9	23.9	64.4	26.2	27.8	44.2	35.0	55.0	34.3	23.5	43.5	34.6
ConciseRL (Gemini 2.5 Flash-Lite)	84.2	45.4	80.6	33.7	20.2	49.9	28.6	64.6	34.9	26.8	49.7	44.1
Eff. Reasoning α =0.2	79.8	43.4	79.6	32.8	22.2	38.6	29.9	62.0	34.9	41.7	49.3	43.7
Eff. Reasoning α =0.4	74.7	41.9	76.8	32.4	21.8	45.2	31.5	52.1	33.3	31.2	47.6	40.6
Eff. Reasoning α =0.1	82.5	63.7	78.4	43.9	21.6	42.6	32.1	56.9	34.5	43.0	49.8	50.0
Cosine Reward	80.4	310.8	77.0	79.9	19.8	93.8	32.1	91.3	29.4	93.2	47.7	133.8
DeepScaleR	80.7	332.7	82.6	60.0	31.5	67.7	32.6	71.3	34.8	57.0	52.4	117.7
Full Reasoning	76.3	100.0	71.4	100.0	25.8	100.0	26.6	100.0	26.3	100.0	45.3	100.0

Table 1: Comparison of accuracy (higher is better) and token length (as % of Full Reasoning; lower is better) across datasets for DeepSeek-R1-Distill-Qwen-1.5B. We group our methods with baselines that reach similar average reasoning-trace length, and **bold the best accuracy within each such group**. This highlights that our methods achieve stronger performance at similar trace length.

	GSM8K		MATH500		MMLU-Pro-1k		GPQA-main		TheoremQA		Average	
Technique	Acc.	Len.(%)	Acc.	Len.(%)	Acc.	Len.(%)	Acc.	Len.(%)	Acc.	Len.(%)	Acc.	Len.(%)
ConciseRL	78.5	11.0	77.2	33.7	19.9	32.3	28.3	34.3	29.6	20.6	46.7	26.4
Eff. Reasoning α =0.4	60.9	4.9	68.0	27.8	18.8	31.0	31.5	36.1	33.1	20.1	42.5	24.0
Eff. Reasoning α =0.9	14.6	0.4	18.8	2.0	14.6	0.1	27.7	0.2	16.5	0.1	18.4	0.6
Eff. Reasoning α =0.6	81.3	17.9	77.6	44.1	24.3	56.3	31.2	53.9	34.1	45.5	49.7	43.5
Eff. Reasoning α =0.2	81.3	17.9	77.6	44.1	24.3	56.3	31.2	53.9	34.1	45.5	49.7	43.5
Eff. Reasoning α =0.1	85.0	34.6	82.0	56.6	20.4	68.5	30.4	82.2	36.1	52.6	50.8	58.9
Cosine Reward	43.6	598.7	44.4	298.8	17.8	216.2	28.8	138.5	13.6	199.9	29.6	290.4
DeepScaleR	80.7	138.0	82.6	82.5	31.5	75.7	32.6	72.1	34.8	80.6	52.4	89.8
Full Reasoning	83.1	100.0	78.8	100.0	32.0	100.0	31.0	100.0	29.8	100.0	50.9	100.0

Table 2: Comparison of accuracy (higher is better) and token length (as % of Full Reasoning; lower is better) across datasets for STILL-3-1.5B-preview. We group our methods with baselines that reach similar average reasoning-trace length, and **bold the best accuracy within each such group**. This highlights that our methods achieve stronger performance at similar trace length.

At comparable length, our method outperforms Efficient Reasoning $\alpha{=}0.4$ by 4.2% in average accuracy. In contrast, Efficient Reasoning displays unstable performance across α values: while $\alpha{=}0.1$, $\alpha{=}0.2$, and $\alpha{=}0.6$ retain high accuracy but produce moderately long traces, aggressive settings like $\alpha{=}0.9$ sharply reduce length but cause severe performance drops. Hence, the Efficient Reasoning method displays noticeable hyperparameter sensitivity and requires considerable fine-tuning to achieve favorable performance-length trade-offs. Our ConciseRL method avoids this issue by not relying on additional hyperparameters or manual tuning. For examples of reasoning traces from different methods, see Section A.4.

During training, our model follows a similar accuracy trajectory as Efficient Reasoning $\alpha=0.2$ but converges to response lengths comparable to $\alpha=0.4$ as can be seen in Figure 3, thus achieving the benefits of both ends of the trade-off curve. Addi-

tionally, Figure 7 shows that our separated variant achieves the shortest reasoning traces among all methods, with an accuracy above the $\alpha=0.6$ baseline. The conciseness signal from our semantic reward enables this fine-grained control, in contrast to static heuristics that uniformly penalize length. This adaptivity is also reflected in the performance across the other datasets: our method achieves the shortest reasoning traces in all benchmarks among methods with similar accuracy. These results indicate that conciseness-aware optimization offers a more efficient path to high-quality reasoning than static length constraints or cosine penalties.

Overall, static length penalties require careful tuning, and small changes often push methods to either sacrifice accuracy or use many more tokens during inference. In contrast, our dynamic conciseness reward has no additional hyperparameters and consistently yields stronger accuracy at substantially lower token budgets across all settings.

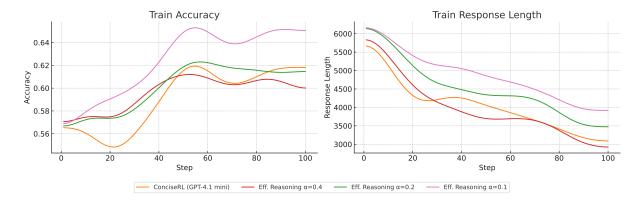


Figure 3: Training metrics across steps using DeepSeek-R1-Distill-Qwen-1.5B as the base model. The Y-axes show accuracy (higher is better) and response length in tokens (lower is better).

	GSM8K		MATH500		MMLU-Pro-1k		GPQA-main		TheoremQA		Average	
Technique	Acc.	Len.(%)	Acc.	Len.(%)	Acc.	Len.(%)	Acc.	Len.(%)	Acc.	Len.(%)	Acc.	Len.(%)
ConciseRL (Gemini 2.5 Flash-Lite)	84.2	45.4	80.6	33.7	20.2	49.9	28.6	64.6	34.9	26.8	49.7	44.1
ConciseRL (Gemini 2.5 Flash)	77.3	16.4	73.6	9.3	18.9	15.7	30.4	15.8	33.3	5.0	46.7	12.4
ConciseRL (GPT-4.1 mini)	80.9	35.8	78.0	30.8	24.5	38.5	30.4	45.9	32.5	21.0	49.3	34.4
ConciseRL (GPT-4o mini)	82.8	34.0	77.4	28.5	19.1	28.5	28.6	36.4	34.8	17.9	48.5	29.1
ConciseRL (GPT-4.1 nano)	83.6	190.0	75.6	76.6	17.5	58.7	30.6	65.8	29.3	85.8	47.3	95.4
Full Reasoning	76.3	100.0	71.4	100.0	25.8	100.0	26.6	100.0	26.3	100.0	45.3	100.0

Table 3: Comparison of accuracy (higher is better) and token length (as % of Full Reasoning; lower is better) across datasets for DeepSeek-R1-Distill-Qwen-1.5B when using different judge models.

4.3 Judge Model Comparison

The quality of the LLM judge used to assess conciseness has a significant influence on the final model's behavior. Figure 6 (appendix) shows that a more capable judge enables more effective optimization. GPT-4.1 mini (OpenAI, 2025a) and GPT-40 mini (OpenAI, 2024a) yield models that significantly reduce response length during training while maintaining or improving accuracy. In contrast, GPT-4.1 nano (OpenAI, 2025a), despite producing slightly higher training accuracy in the end, does not meaningfully shorten the reasoning traces, likely because the score has more "noise". We highlight that the accuracy curve is not necessarily indicative of final model performance: accuracy rises early as the model learns the dataset, but later decreases when the reward begins favoring conciseness more strongly and the traces become shorter.

Tables 1 and 3 show the efficiency–accuracy trade-offs introduced by each judge. ConciseRL with the Gemini 2.5 Flash-Lite judge achieves the highest average accuracy at 49.7%, utilizing 44.1% of the full reasoning tokens, and leads on GSM8K, MATH500, TheoremQA, and the overall average. ConciseRL with the Gemini 2.5 Flash judge provides the strongest compression, reach-

ing 46.7% average accuracy at only 12.4% of the length and improving TheoremQA from 26.3% to 33.3% while reducing the trace to 5.0%. GPT-40 mini reaches 48.5% accuracy at 29.1% length, substantially below the 100% length of Full Reasoning at 45.3% accuracy. GPT-4.1 mini delivers a strong efficiency—accuracy trade-off, matching the average accuracy of Efficient Reasoning at $\alpha=0.2$ (49.3%) while producing 65.6% shorter reasoning traces. On the other hand, GPT-4.1 nano offers limited trace compression (95.4% of full reasoning) despite comparable accuracy (47.3%), validating the earlier observation that noisier reward signals make optimization harder.

4.4 The Behavior of the Rewards

In Figure 4a, we compare reward values across different judges. Despite starting with the lowest average reward, the model trained with GPT-4.1 mini ends with the highest, indicating that the reward, though initially strict, was more learnable and resulted in effective policy improvement. This contrasts with GPT-4.1 nano, which yields higher early rewards but fails to produce meaningful shortening, likely due to its noisy and less discriminative signal. Notably, these values represent final reward outputs: a score of zero is assigned if the

Technique	Level 1		Level 2		Level 3		L	evel 4	Level 5		Average	
	Acc.	Len.	Acc.	Len.	Acc.	Len.	Acc.	Len.	Acc.	Len.	Acc.	Len.
ConciseRL (Separated)	95.35	118.40	77.78	159.95	80.95	476.69	64.84	1470.79	47.78	4431.25	73.34	1331.42
ConciseRL	95.35	265.58	87.77	453.94	87.62	2073.16	75.78	2738.97	57.46	8029.31	80.80	2712.19
Eff. Reasoning α =0.1	93.02	1285.63	88.88	1102.38	88.57	3355.63	74.22	4536.74	62.69	9004.73	81.48	3857.02
Full Reasoning	88.37	3686.48	85.56	5100.22	77.14	9949.32	73.44	10890.66	50.00	16156.81	74.90	9156.70

Table 4: Comparison of accuracy (higher is better) and token length (lower is better) by difficulty level on MATH500, using DeepSeek-R1-Distill-Qwen-1.5B. For ConciseRL, we use GPT-4.1 mini as the conciseness judge. **Bold** indicates the shortest reasoning traces.

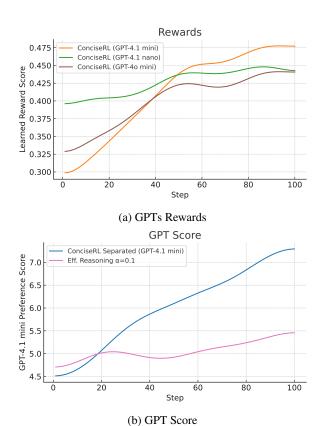


Figure 4: Training metrics across steps using different models as the judge. The Y-axes show reward values and conciseness scores assigned by the judge.

model's answer is incorrect (Equation 1), meaning the curves are also influenced by accuracy.

In Figure 4b, we isolate conciseness from correctness using our separated reward configuration. The resulting GPT score, evaluated solely on reasoning quality, improves almost linearly from roughly 4.5 to 7.5 out of 10 over the course of training. Compared to Efficient Reasoning α =0.1, we observe that although it successfully reduces trace length to nearly half (Table 2), its corresponding conciseness score only increases from 4.75 to 5.5. In contrast, our method boosts the score from 4.5 to 7.1. This shows that shorter traces are not necessarily more concise, reinforcing our argument that

brevity should be a learned outcome of semantic compactness, not a heuristic target.

4.5 Performance by Problem Difficulty

Figure 1 and Table 4 show MATH500 performance across five difficulty levels. ConciseRL (Separated) achieves up to 31× fewer tokens than the Full Reasoning baseline on Level 1 (118 vs. 3686 tokens) while improving accuracy by 7%. For the hardest problems (Level 5), our method outperforms the Full Reasoning baseline using $3.6 \times$ fewer tokens (4431 vs. 16156 tokens). Our method exhibits a strong correlation between problem difficulty and reasoning length: easier problems are solved with fewer tokens, while harder ones elicit longer reasoning. This adaptive behavior aligns with human intuition and is a desirable trait in reasoning models. In contrast, other techniques like Efficient Reasoning ($\alpha = 0.1$), Cosine Reward, or Full Reasoning either under-adapt (e.g., generating long traces even for trivial problems) or show no clear correlation.

Furthermore, as shown in Table 1, on the TheoremQA benchmark, our method achieves 2.2% higher accuracy while using 12.5× fewer tokens compared to the Full Reasoning model. These results suggest that our semantically guided reward encourages selective elaboration, producing concise reasoning where appropriate and expanding only when necessary. Also, our method produces more explainable and compressed reasoning traces, see Appendix A.4 for examples and analysis.

4.6 Kullback-Leibler Divergence

Figure 5 shows that our method (ConciseRL) induces a KL divergence profile comparable to Efficient Reasoning with $\alpha{=}0.4$, suggesting that our method modifies the policy to a similar degree while achieving significantly better efficiency–accuracy trade-offs. Also, ConciseRL results in far less KL shift than $\alpha{=}0.6$, which exhibits early instability and overshooting. Efficient Reasoning with $\alpha{=}0.9$ produces similarly short

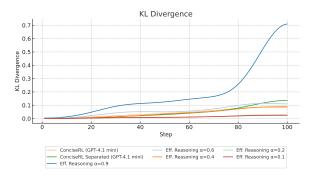


Figure 5: Kullback–Leibler divergence between the updated policy and the initial models throughout training.

reasoning traces but incurs a much higher KL divergence, around $7\times$ of our method, showing inefficient and unstable updates. In contrast, our ConciseRL (Separated) variant maintains low KL divergence throughout training, further reinforcing the benefits of our semantic reward and highlighting that trace length alone does not imply effective or stable policy learning.

5 Conclusion

In this work, we introduce a novel, hyperparameter-free conciseness score used as a reward signal within any reinforcement learning framework and evaluated by an LLM judge to encourage models to generate correct and concise reasoning traces. Our approach improves semantic density without sacrificing accuracy, yielding strong efficiency—accuracy trade-offs. At equal token length, our traces are more informative, interpretable, and explainable than those from static length-penalized baselines due to fewer filler steps and tighter logic. Overall, our results highlight the value of semantic rewards and LLM-based judges for guiding concise and correct reasoning.

In future work, we plan to evaluate our method across a wider range of model sizes to better understand how conciseness rewards scale. Additionally, while we hypothesize that using more advanced judges like GPT-40 or GPT-4.5 could further enhance performance, we could not evaluate this due to our budget limitations and API-associated costs. We are also interested in combining our conciseness optimization with orthogonal methods, such as structure-aware search or explicit length conditioning, to enhance reasoning quality and efficiency. Unlike other methods, ours is orthogonal and integrates well with approaches like L1 (Aggarwal and Welleck, 2025) (see A.2).

Limitations

Despite the strong empirical results, our study has several limitations. Due to limited hardware resources, including access to only 4×NVIDIA A100 80GB GPUs, we could not evaluate larger models, test our approach at scale, or fully compare against all relevant baselines. Training with reinforcement learning is particularly resource-intensive, which further constrained the scope of our experiments. We were also only able to partially explore combinations of our method with other recent techniques. Finally, we restricted our use of LLMs-as-judges to smaller models because of the cost of using larger models, specifically, the higher cost for long reasoning traces.

Ethical Considerations

Regarding impact, we believe that our method has positive implications for reducing the computational cost of reasoning at inference time, particularly for models deployed in environments with constrained compute budgets or latency requirements. By shifting the optimization target from token length to conciseness, we aim to improve efficiency and interpretability without sacrificing correctness.

However, we acknowledge several potential ethical risks. Using LLMs as evaluators introduces potential biases, which may reflect or amplify artifacts from their training data. If the LLM judge favors certain styles of explanation over others, this may steer the trained model toward narrow patterns of reasoning. Additionally, there is a risk that optimizing for conciseness could suppress contextually relevant reasoning steps, leading to overly truncated outputs. We mitigate these concerns by evaluating correctness alongside conciseness and analyzing behavior across problem difficulties and prompt variants. We also ensure that the LLM judge is only used during training, not at inference time, thus avoiding deployment dependencies.

Regarding safeguards, since our method builds upon publicly available pre-trained models without introducing new high-risk capabilities, we assess the risk of misuse to be low. Nevertheless, to support responsible usage, we commit to open-source the code, model weights, and datasets under a research license with clear documentation.

References

- Pranjal Aggarwal and Sean Welleck. 2025. L1: Controlling how long a reasoning model thinks with reinforcement learning. In *Second Conference on Language Modeling*.
- Daman Arora and Andrea Zanette. 2025. Training language models to reason efficiently. *Preprint*, arXiv:2502.04463.
- Maciej Besta, Nils Blach, Ales Kubicek, Robert Gerstenberger, Michal Podstawski, Lukas Gianinazzi, Joanna Gajda, Tomasz Lehmann, Hubert Niewiadomski, Piotr Nyczyk, and Torsten Hoefler. 2024. Graph of thoughts: Solving elaborate problems with large language models. Proceedings of the AAAI Conference on Artificial Intelligence, 38(16):17682–17690.
- Bradley Butcher, Michael O'Keefe, and James Titchener. 2024. Precise length control in large language models. *Preprint*, arXiv:2412.11937.
- Wenhu Chen, Ming Yin, Max Ku, Pan Lu, Yixin Wan, Xueguang Ma, Jianyu Xu, Xinyi Wang, and Tony Xia. 2023. TheoremQA: A theorem-driven question answering dataset. In *Proceedings of the 2023 Conference on Empirical Methods in Natural Language Processing*, pages 7889–7901, Singapore. Association for Computational Linguistics.
- Xingyu Chen, Jiahao Xu, Tian Liang, Zhiwei He, Jianhui Pang, Dian Yu, Linfeng Song, Qiuzhi Liu, Mengfei Zhou, Zhuosheng Zhang, Rui Wang, Zhaopeng Tu, Haitao Mi, and Dong Yu. 2025. Do NOT think that much for 2+3=? on the overthinking of long reasoning models. In *Forty-second International Conference on Machine Learning*.
- Jeffrey Cheng and Benjamin Van Durme. 2024. Compressed chain of thought: Efficient reasoning through dense representations. *Preprint*, arXiv:2412.13171.
- Karl Cobbe, Vineet Kosaraju, Mohammad Bavarian, Mark Chen, Heewoo Jun, Lukasz Kaiser, Matthias Plappert, Jerry Tworek, Jacob Hilton, Reiichiro Nakano, Christopher Hesse, and John Schulman. 2021. Training verifiers to solve math word problems. *Preprint*, arXiv:2110.14168.
- DeepSeek-AI, Daya Guo, Dejian Yang, Haowei Zhang, Junxiao Song, Ruoyu Zhang, Runxin Xu, Qihao Zhu, Shirong Ma, Peiyi Wang, Xiao Bi, Xiaokang Zhang, Xingkai Yu, Yu Wu, Z. F. Wu, Zhibin Gou, Zhihong Shao, Zhuoshu Li, Ziyi Gao, and 181 others. 2025. Deepseek-r1: Incentivizing reasoning capability in Ilms via reinforcement learning. *Preprint*, arXiv:2501.12948.
- Mehdi Fatemi, Banafsheh Rafiee, Mingjie Tang, and Kartik Talamadupula. 2025. Concise reasoning via reinforcement learning. *Preprint*, arXiv:2504.05185.
- Google Gemini Team. 2025. Gemini 2.5: Pushing the frontier with advanced reasoning, multimodality, long context, and next generation agentic capabilities. *Preprint*, arXiv:2507.06261.

- Jiawei Gu, Xuhui Jiang, Zhichao Shi, Hexiang Tan, Xuehao Zhai, Chengjin Xu, Wei Li, Yinghan Shen, Shengjie Ma, Honghao Liu, Saizhuo Wang, Kun Zhang, Yuanzhuo Wang, Wen Gao, Lionel Ni, and Jian Guo. 2025. A survey on llm-as-a-judge. *Preprint*, arXiv:2411.15594.
- Tingxu Han, Zhenting Wang, Chunrong Fang, Shiyu Zhao, Shiqing Ma, and Zhenyu Chen. 2025. Token-budget-aware LLM reasoning. In *Findings of the Association for Computational Linguistics: ACL 2025*, pages 24842–24855, Vienna, Austria. Association for Computational Linguistics.
- Shibo Hao, Sainbayar Sukhbaatar, DiJia Su, Xian Li, Zhiting Hu, Jason Weston, and Yuandong Tian. 2024. Training large language models to reason in a continuous latent space. *Preprint*, arXiv:2412.06769.
- Dan Hendrycks, Collin Burns, Saurav Kadavath, Akul Arora, Steven Basart, Eric Tang, Dawn Song, and Jacob Steinhardt. 2021. Measuring mathematical problem solving with the MATH dataset. In *Thirty-fifth Conference on Neural Information Processing Systems Datasets and Benchmarks Track (Round 2)*.
- Bairu Hou, Yang Zhang, Jiabao Ji, Yujian Liu, Kaizhi Qian, Jacob Andreas, and Shiyu Chang. 2025. Thinkprune: Pruning long chain-of-thought of llms via reinforcement learning. *Preprint*, arXiv:2504.01296.
- Renlong Jie, Xiaojun Meng, Lifeng Shang, Xin Jiang, and Qun Liu. 2023. Prompt-based length controlled generation with reinforcement learning. *Preprint*, arXiv:2308.12030.
- Yu Kang, Xianghui Sun, Liangyu Chen, and Wei Zou. 2025. C3ot: Generating shorter chain-of-thought without compromising effectiveness. *Proceedings of the AAAI Conference on Artificial Intelligence*, 39(23):24312–24320.
- Woosuk Kwon, Zhuohan Li, Siyuan Zhuang, Ying Sheng, Lianmin Zheng, Cody Hao Yu, Joseph Gonzalez, Hao Zhang, and Ion Stoica. 2023. Efficient memory management for large language model serving with pagedattention. In *Proceedings of the 29th Symposium on Operating Systems Principles*, SOSP '23, page 611–626, New York, NY, USA. Association for Computing Machinery.
- Dawei Li, Bohan Jiang, Liangjie Huang, Alimohammad Beigi, Chengshuai Zhao, Zhen Tan, Amrita Bhattacharjee, Yuxuan Jiang, Canyu Chen, Tianhao Wu, Kai Shu, Lu Cheng, and Huan Liu. 2025. From generation to judgment: Opportunities and challenges of llm-as-a-judge. *Preprint*, arXiv:2411.16594.
- Jia LI, Edward Beeching, Lewis Tunstall, Ben Lipkin, Roman Soletskyi, Shengyi Costa Huang, Kashif Rasul, Longhui Yu, Albert Jiang, Ziju Shen, Zihan Qin, Bin Dong, Li Zhou, Yann Fleureau, Guillaume Lample, and Stanislas Polu. 2024. Numinamath. [https://huggingface.co/AI-MO/NuminaMath-CoT] (https:

- //github.com/project-numina/
 aimo-progress-prize/blob/main/
 report/numina_dataset.pdf).
- Yue Liu, Jiaying Wu, Yufei He, Ruihan Gong, Jun Xia, Liang Li, Hongcheng Gao, Hongyu Chen, Baolong Bi, Jiaheng Zhang, Zhiqi Huang, Bryan Hooi, Stan Z. Li, and Keqin Li. 2025. Efficient inference for large reasoning models: A survey. *Preprint*, arXiv:2503.23077.
- Haotian Luo, Li Shen, Haiying He, Yibo Wang, Shiwei Liu, Wei Li, Naiqiang Tan, Xiaochun Cao, and Dacheng Tao. 2025a. O1-pruner: Length-harmonizing fine-tuning for o1-like reasoning pruning. In 2nd AI for Math Workshop @ ICML 2025.
- Michael Luo, Sijun Tan, Justin Wong, Xiaoxiang Shi, William Tang, Manan Roongta, Colin Cai, Jeffrey Luo, Tianjun Zhang, Erran Li, Raluca Ada Popa, and Ion Stoica. 2025b. Deepscaler: Surpassing olpreview with a 1.5b model by scaling rl. Notion Blog.
- Yu Meng, Mengzhou Xia, and Danqi Chen. 2024. Simpo: Simple preference optimization with a reference-free reward. In *Advances in Neural Information Processing Systems*, volume 37, pages 124198–124235. Curran Associates, Inc.
- Yingqian Min, Zhipeng Chen, Jinhao Jiang, Jie Chen, Jia Deng, Yiwen Hu, Yiru Tang, Jiapeng Wang, Xiaoxue Cheng, Huatong Song, Wayne Xin Zhao, Zheng Liu, Zhongyuan Wang, and Ji-Rong Wen. 2024. Imitate, explore, and self-improve: A reproduction report on slow-thinking reasoning systems. arXiv preprint arXiv:2412.09413.
- Niklas Muennighoff, Zitong Yang, Weijia Shi, Xiang Lisa Li, Li Fei-Fei, Hannaneh Hajishirzi, Luke Zettlemoyer, Percy Liang, Emmanuel Candes, and Tatsunori Hashimoto. 2025. s1: Simple test-time scaling. In Workshop on Reasoning and Planning for Large Language Models.
- OpenAI. 2024a. Openai gpt-4o system card.
- OpenAI. 2024b. Openai o1 system card.
- OpenAI. 2025a. Introducing gpt-4.1 in the api.
- OpenAI. 2025b. Openai o3 and o4-mini system card.
- Long Ouyang, Jeffrey Wu, Xu Jiang, Diogo Almeida, Carroll Wainwright, Pamela Mishkin, Chong Zhang, Sandhini Agarwal, Katarina Slama, Alex Ray, John Schulman, Jacob Hilton, Fraser Kelton, Luke Miller, Maddie Simens, Amanda Askell, Peter Welinder, Paul F Christiano, Jan Leike, and Ryan Lowe. 2022. Training language models to follow instructions with human feedback. In *Advances in Neural Information Processing Systems*, volume 35, pages 27730–27744. Curran Associates, Inc.

- David Rein, Betty Li Hou, Asa Cooper Stickland, Jackson Petty, Richard Yuanzhe Pang, Julien Dirani, Julian Michael, and Samuel R. Bowman. 2024. GPQA: A graduate-level google-proof q&a benchmark. In *First Conference on Language Modeling*.
- John Schulman, Filip Wolski, Prafulla Dhariwal, Alec Radford, and Oleg Klimov. 2017. Proximal policy optimization algorithms. *Preprint*, arXiv:1707.06347.
- Zhihong Shao, Peiyi Wang, Qihao Zhu, Runxin Xu, Junxiao Song, Xiao Bi, Haowei Zhang, Mingchuan Zhang, Y. K. Li, Y. Wu, and Daya Guo. 2024. Deepseekmath: Pushing the limits of mathematical reasoning in open language models. *Preprint*, arXiv:2402.03300.
- Yi Shen, Jian Zhang, Jieyun Huang, Shuming Shi, Wenjing Zhang, Jiangze Yan, Ning Wang, Kai Wang, and Shiguo Lian. 2025. Dast: Difficulty-adaptive slow-thinking for large reasoning models. *Preprint*, arXiv:2503.04472.
- Charlie Victor Snell, Jaehoon Lee, Kelvin Xu, and Aviral Kumar. 2025. Scaling LLM test-time compute optimally can be more effective than scaling parameters for reasoning. In *The Thirteenth International Conference on Learning Representations*.
- Yang Sui, Yu-Neng Chuang, Guanchu Wang, Jiamu Zhang, Tianyi Zhang, Jiayi Yuan, Hongyi Liu, Andrew Wen, Shaochen Zhong, Na Zou, Hanjie Chen, and Xia Hu. 2025. Stop overthinking: A survey on efficient reasoning for large language models. *Transactions on Machine Learning Research*.
- Kimi Team, Angang Du, Bofei Gao, Bowei Xing, Changjiu Jiang, Cheng Chen, Cheng Li, Chenjun Xiao, Chenzhuang Du, Chonghua Liao, Chuning Tang, Congcong Wang, Dehao Zhang, Enming Yuan, Enzhe Lu, Fengxiang Tang, Flood Sung, Guangda Wei, Guokun Lai, and 75 others. 2025. Kimi k1.5: Scaling reinforcement learning with llms. *Preprint*, arXiv:2501.12599.
- Qwen Team. 2025a. Qwq-32b: Embracing the power of reinforcement learning.
- RUCAIBox STILL Team. 2025b. Still-3-1.5b-preview: Enhancing slow thinking abilities of small models through reinforcement learning.
- Xuezhi Wang, Jason Wei, Dale Schuurmans, Quoc V Le, Ed H. Chi, Sharan Narang, Aakanksha Chowdhery, and Denny Zhou. 2023. Self-consistency improves chain of thought reasoning in language models. In *The Eleventh International Conference on Learning Representations*.
- Yubo Wang, Xueguang Ma, Ge Zhang, Yuansheng Ni, Abhranil Chandra, Shiguang Guo, Weiming Ren, Aaran Arulraj, Xuan He, Ziyan Jiang, Tianle Li, Max Ku, Kai Wang, Alex Zhuang, Rongqi Fan, Xiang Yue, and Wenhu Chen. 2024. Mmlu-pro: A more robust and challenging multi-task language understanding

benchmark. In *Advances in Neural Information Processing Systems*, volume 37, pages 95266–95290. Curran Associates, Inc.

Jason Wei, Xuezhi Wang, Dale Schuurmans, Maarten Bosma, brian ichter, Fei Xia, Ed Chi, Quoc V Le, and Denny Zhou. 2022. Chain-of-thought prompting elicits reasoning in large language models. In *Advances in Neural Information Processing Systems*, volume 35, pages 24824–24837. Curran Associates, Inc.

Huanjin Yao, Jiaxing Huang, Wenhao Wu, Jingyi Zhang, Yibo Wang, Shunyu Liu, Yingjie Wang, Yuxin Song, Haocheng Feng, Li Shen, and Dacheng Tao. 2024. Mulberry: Empowering mllm with o1-like reasoning and reflection via collective monte carlo tree search. *Preprint*, arXiv:2412.18319.

Shunyu Yao, Dian Yu, Jeffrey Zhao, Izhak Shafran, Tom Griffiths, Yuan Cao, and Karthik Narasimhan. 2023. Tree of thoughts: Deliberate problem solving with large language models. In *Advances in Neural Information Processing Systems*, volume 36, pages 11809–11822. Curran Associates, Inc.

Edward Yeo, Yuxuan Tong, Morry Niu, Graham Neubig, and Xiang Yue. 2025. Demystifying long chain-of-thought reasoning in llms. *Preprint*, arXiv:2502.03373.

Lianmin Zheng, Wei-Lin Chiang, Ying Sheng, Siyuan Zhuang, Zhanghao Wu, Yonghao Zhuang, Zi Lin, Zhuohan Li, Dacheng Li, Eric Xing, Hao Zhang, Joseph E Gonzalez, and Ion Stoica. 2023. Judging llm-as-a-judge with mt-bench and chatbot arena. In *Advances in Neural Information Processing Systems*, volume 36, pages 46595–46623. Curran Associates, Inc.

Jason Zhu and Hongyu Li. 2025. Towards concise and adaptive thinking in large reasoning models: A survey. *Preprint*, arXiv:2507.09662.

A Appendix

A.1 Extended Results

As a supplement to our main findings, we provide additional results that reinforce the effectiveness and generality of our approach. Figure 6 compares training dynamics under different LLM judges, showing that stronger judges like GPT-4.1 mini and GPT-40 mini enable more effective reasoning compression while preserving or improving accuracy. In contrast, GPT-4.1 nano yields higher variance and fails to consistently reduce reasoning length. These trends confirm that the reliability of the conciseness signal strongly depends on the evaluator's capability (OpenAI, 2025a).

Figure 7 extends out comparison with static baselines such as Efficient Reasoning (DeepSeek-AI

et al., 2025) display a trade-off between brevity and correctness, with aggressive penalties reducing token usage but at the cost of sharp accuracy drops. Our method, ConciseRL, follows a more desirable trajectory: it achieves the conciseness of the most efficient baselines while maintaining high accuracy, effectively tracing out the Pareto frontier.

Table 5 presents raw token counts on the DeepSeek-R1-Distill-Qwen-1.5B model across all benchmarks. The separated reward variant (ConciseRL Separated) yields the most compact reasoning traces, while ConciseRL strikes a strong efficiency—accuracy balance. Notably, both outperform static methods like Cosine Reward and Efficient Reasoning across all datasets.

Finally, Table 6 replicates this comparison on the STILL-3-1.5B-preview model (Team, 2025b; Min et al., 2024), further supporting the generalizability of our framework. The same trends persist: semantic conciseness rewards outperform length-based baselines, both in efficiency and accuracy. These results confirm that our approach is robust across architectures, training signals, and judge variants.

A.2 Future Work Discussion

Although our method rewards semantic conciseness, L1 uses explicit length-based RL objectives to precisely control token-level output length. This difference in objectives means that the two approaches could complement each other, enabling simultaneous control of the token budget and semantic efficiency without interference. Lengthbased scoring approaches, however, inherently conflict with L1's exact or maximum length constraints since they introduce competing reward signals focused solely on token minimization rather than strict adherence to specified token budgets. Therefore, integrating length-based scoring directly alongside L1 is impractical, highlighting the value of our conciseness approach as a complementary strategy to optimize reasoning efficiency and performance. Unfortunately, running such experiments requires significant compute, at least 8×NVIDIA A100 80GB GPUs even for a 1.5B model, which we do not have access to.

A.3 Conciseness Evaluation Prompt

To evaluate the conciseness of reasoning traces, we use an LLM-based judge that assigns a score between 1 (overly verbose) and 10 (maximally concise), independent of correctness. The system prompt provided to the judge is designed to guide it

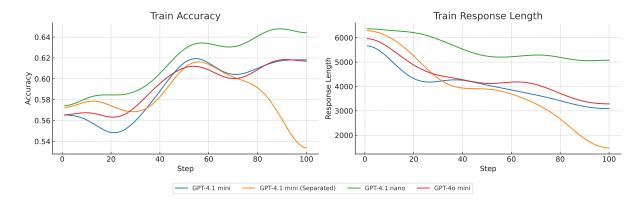


Figure 6: Training metrics across steps using DeepSeek-R1-Distill-Qwen-1.5B (DeepSeek-AI et al., 2025) as the base model and different models as the judge. The Y-axes show accuracy (higher is better; left) and response length in tokens (lower is better; right). The X-axis in both cases shows the training step.

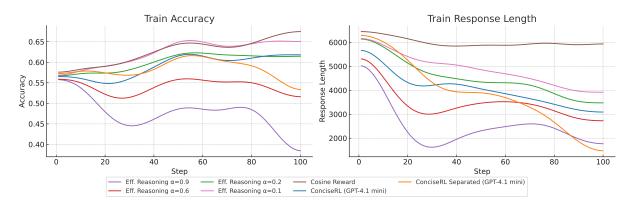


Figure 7: Training metrics across steps using DeepSeek-R1-Distill-Qwen-1.5B (DeepSeek-AI et al., 2025) as the base model and GPT-4.1 mini (OpenAI, 2025a) as the judge. The Y-axes show accuracy (higher is better; left) and response length in tokens (lower is better; right). The X-axis in both cases shows the training step.

toward assessing semantic compactness rather than surface-level brevity. The exact prompt we used can be seen in Figure 8.

A.4 Reasoning Trace Examples

To better illustrate the differences between reward strategies, we provide representative reasoning traces from ConciseRL, ConciseRL (Separated), and Efficient Reasoning with $\alpha=0.4$. We selected $\alpha=0.4$ because it achieves accuracy comparable to ConciseRL on MATH500, making it a fair baseline for direct comparison. We focus on examples where all methods correctly solve the problem to ensure the comparison of conciseness is not confounded by correctness. For more aggressive penalty settings like $\alpha=0.6$ or $\alpha=0.9$, it was difficult to find consistent examples across all levels that resulted in correct answers, so they were excluded.

The conciseness judge rewards *semantic density*: a trace must enumerate every logically necessary

step while avoiding rhetorical or computational detours. Because the reward is assigned by a large-model evaluator, it is sensitive to redundancy that naive length penalties overlook, yet agnostic to correctness, which we enforce separately. This design choice underlies all the gains discussed below.

Level-1 Example (Figure 9–Figure 10). Our CONCISERL trace solves the parity-of-students puzzle in **237** tokens, with the gated variant using only **170**. Efficient-Reasoning ($\alpha = 0.4$) needs **257**, while Cosine-Reward expands to **1278**. Most of the extra tokens in the baselines are meta-commentary or reintroductions of already stated variables; our trace states the constraints once, enumerates the only two admissible multiples of 13, and dismisses the invalid option in a single inequality check, yielding length savings without omitting any logical step.

Level-2 Example (Figure 11–Figure 12). For the "roots of unity" problem, CONCISERL finishes

	GSI	M8K	MA	MATH500		MMLU-Pro-1k		GPQA-main		TheoremQA		erage
Technique	Acc.	Len.	Acc.	Len.	Acc.	Len.	Acc.	Len.	Acc.	Len.	Acc.	Len.
ConciseRL (Separated)	72.5	248	68.6	1703	19.2	2904	31.0	4913	28.5	1534	43.96	2260
ConciseRL	80.9	543	78.0	3221	24.5	5732	30.4	11429	32.5	4034	49.26	4992
ConciseRL (GPT-4o mini)	82.8	515	77.4	2971	19.1	4240	28.6	9045	34.8	3439	48.54	4042
ConciseRL (GPT-4.1 nano)	83.6	2880	75.6	7994	17.5	8738	30.6	16363	29.3	16510	47.32	10497
Eff. Reasoning α =0.9	16.1	19	33.6	2195	18.2	1723	27.9	6509	20.4	1872	23.24	2463
Eff. Reasoning α =0.6	55.9	362	64.4	2735	27.8	6583	35.0	13675	34.3	4533	43.48	5578
Eff. Reasoning α =0.4	74.7	635	76.8	3382	21.8	6730	31.5	12973	33.3	6009	47.62	5946
Eff. Reasoning α =0.2	79.8	658	79.6	3428	22.2	5746	29.9	15415	34.9	8029	49.28	6655
Eff. Reasoning α =0.1	82.5	966	78.4	4588	21.6	6347	32.1	14157	34.5	8271	49.82	6866
Cosine Reward	80.4	4711	77.0	8348	19.8	13962	32.1	22708	29.4	17933	47.74	13532
DeepScaleR	80.7	5044	82.6	6160	31.5	10071	32.6	17742	34.8	10970	52.4	9997.4
Full Reasoning	76.3	1516	71.4	10442	25.8	14886	26.6	24880	26.3	19251	45.28	14195

Table 5: Comparison of accuracy (higher is better) and token length (lower is better) across datasets for DeepSeek-R1-Distill-Qwen-1.5B (DeepSeek-AI et al., 2025). For ConciseRL, we use GPT-4.1 mini (OpenAI, 2025a) as the conciseness judge.

	GSM8K		MATH500		MMLU-Pro-1k		GPQA-main		TheoremQA		Average	
Technique	Acc.	Len.	Acc.	Len.	Acc.	Len.	Acc.	Len.	Acc.	Len.	Acc.	Len.
ConciseRL	78.5	403	77.2	2519	19.9	4299	28.3	7670	29.6	3141	46.7	3606.4
Eff. Reasoning α =0.9	14.6	15	18.8	151	14.6	18	27.7	40	16.5	19	18.4	48.6
Eff. Reasoning α =0.6	81.3	655	77.6	676	24.3	1354	31.2	1252	34.1	214	49.7	727.8
Eff. Reasoning α =0.4	60.9	178	68.0	2080	18.8	4126	31.5	8082	33.1	3056	42.5	3504.4
Eff. Reasoning α =0.2	81.3	655	77.6	3291	24.3	7498	31.2	12078	34.1	6924	49.7	6089.2
Eff. Reasoning α =0.1	85.0	1266	82.0	4226	20.4	9117	30.4	18408	36.1	7997	50.8	8202.8
Cosine Reward	43.6	21887	44.4	22320	17.8	28788	28.8	31027	13.6	30430	29.6	26890.4
DeepScaleR	80.7	5044	82.6	6160	31.5	10071	32.6	17742	34.8	10970	52.4	9997.4
Full Reasoning	83.1	3656	78.8	7470	32.0	13310	31.0	22398	29.8	15219	50.9	12410.6

Table 6: Comparison of accuracy (higher is better) and token length (raw token count; lower is better) across datasets for STILL-3-1.5B-preview (Team, 2025b; Min et al., 2024). For ConciseRL, we use GPT-4.1 mini (OpenAI, 2025a) as the conciseness judge.

in **331** tokens (and **217** for the separated reward), whereas Efficient-Reasoning $\alpha=0.4$ requires **842**. The baselines repeatedly paraphrase the factorisation $z^6-1=(z^2-1)(z^4+z^2+1)$, reformulate roots of unity three times, and speculate about smaller n; our model invokes the factorisation exactly once, recognises that the remaining roots are primitive sixth roots, and concludes immediately. This trims the narrative overhead while keeping the algebra explicit.

Level-3 example (Figure 13–Figure 14). In the handshake-counting puzzle, CONCISERL compresses the reasoning to 115 tokens (or 69 in the separated variant) versus 535 for Efficient-Reasoning and 3035 for Cosine-Reward. The baselines re-read the prompt, rehearse elementary combinatorics, and pepper in "thinking-aloud" asides; our trace states the bipartite handshake structure in one declarative sentence and performs the 6×10 multiplication once, making the combinatorial in-

sight transparent.

Level-4 Example (Figure 15–Figure 16). When converting repeating decimals to fractions, CONCISERL requires only 130 tokens, while Efficient-Reasoning consumes 863 and Cosine-Reward 1467. All three traces apply identical fraction conversions; our model simply avoids restating intermediate results and deletes free-form reflections such as "Hmm, repeating decimals can be tricky," leading to reduction in length with no loss of clarity.

Level-5 Example (Figure 17–Figure 18). For the leaking-bucket problem, CONCISERL uses 192 tokens (or 76 separated) versus 437 for Efficient-Reasoning, 5 086 for Cosine-Reward, and 1139 for full reasoning. Baselines repeatedly explain geometric decay and rewrite the same power of $\frac{2}{3}$; our trace articulates the multiplicative factor once and streams three intermediate states inline, preserving every quantitative step while cutting the

```
System prompt:
You are an expert evaluator tasked with scoring the conciseness of a reasoning trace from an AI
Conciseness means providing clear, precise, and direct reasoning.
- High-scoring reasoning (8-10) is brief yet explicitly demonstrates the logical steps or thought
processes clearly.
 Medium-scoring reasoning (5-7) might have minor redundancy, slight verbosity, or slightly unclear
phrasing.
- Low-scoring reasoning (1-4) is either overly verbose, repetitive, vague, contains placeholders,
or is too brief.
Do NOT reward extremely short traces that only state the final answer without reasoning.
Tags like <think>, </think>, <answer>, </answer> are acceptable and should NOT affect scoring.
Evaluate ONLY conciseness. Ignore correctness or accuracy entirely.
Provide ONLY a single integer from 1 (least concise) to 10 (most concise). Do NOT include
explanations or any additional text.
Reasoning Trace:
\{response\ text\}
Conciseness Score (1-10):
```

Figure 8: System propmpt used for the LLM conciseness judge.

filler.

Emergent Pattern across Difficulty Levels.

Across Levels 1–5, CONCISERL consistently removes filler phrases, nested restatements, and speculative digressions, leaving only the algebraic core. The percentage of tokens saved grows with problem difficulty, because long chains of identical transformations are where verbosity compounds. Importantly, the resulting traces are *more* readable: each symbol is introduced once, variables are reused consistently, and the final boxed answer is reached without back-tracking. The concision gain does not come at the expense of policy stability; Appendix 4.6 shows that our KL-divergence curve mirrors that of Efficient-Reasoning $\alpha = 0.4$ while producing far shorter traces and higher average accuracy.

Why Our Model is Better. The qualitative evidence above corroborates the quantitative results from the paper: by optimising a *semantic* conciseness reward, CONCISERL (i) achieves **higher accuracy** at matched or lower token budgets, because unnecessary tokens often accompany spurious logical branches that hurt correctness; (ii) offers **superior interpretability**, as every surviving token carries deductible mathematical purpose, enabling rapid human or downstream-model verification; and (iii) incurs **lower inference cost**—up to 31x fewer generated tokens on easy problems—without

extra hyper-parameters or dataset-specific tuning. In short, our reward converts verbosity into signal, aligning the policy with human preferences for clear, economical argumentation.

MATH500 - Level 1 Question: Thirty-five students attend a math club meeting. The number of girls at the meeting is a multiple of 13, and there are more girls than boys attending the meeting. How many boys are at the meeting?

Answer: 9

Reasoning Trace - ConciseRL - 237 tokens:

```
First, let \ (b\ ) represent the number of boys and \ (g\ ) represent the number of girls. We
1. (b + g = 35)
2. \( g \) is a multiple of 13.
3. \( g > b \).
39 (but 39 exceeds 35, so only 13 and 26).
Check \( g = 26 \):
- Then \( b = 35 - 26 = 9 \).
- \( g > b \) since 26 > 9. This works.
Check \( g = 13 \):
- Then \( b = 35 - 13 = 22 \).
- Then \( b = 35 - 13 = \angle 2 \).
- \( g > b \) since 13 < 22. This does not work.
Thus, the only valid solution is (b = 9).
Reasoning Trace - ConciseRL (Separated) - 170 tokens:
Given that the total number of students is 35, and the number of girls is a multiple of 13, let's
denote the number of girls as 13k, where k is a positive integer. Then, the number of boys would
be 35 - 13k.
Since there are more girls than boys, we have:
13k > 35 - 13k
26k > 35
k > 35/26 \approx 1.346
Thus, k must be at least 2. Substituting k\!=\!2\!:
Number of girls = 26
Number of boys = 35 - 26 = 9
Therefore, there are 9 boys at the meeting.
```

Figure 9: Example reasoning traces generated by our methods ("ConciseRL" and "ConciseRL (Separated)") on a Level 1 MATH500 (Hendrycks et al., 2021) question.

MATH500 - Level 1 Question: Thirty-five students attend a math club meeting. The number of girls at the meeting is a multiple of 13, and there are more girls than boys attending the meeting. How many boys are at the meeting?

Answer: 9

```
Reasoning Trace - Efficient Reasoning \alpha = 0.4 - 257 tokens:
```

```
Let \( G \) be the number of girls and \( B \) be the number of boys. We know:
1. \( G + B = 35 \)
2. \ (G \ ) is a multiple of 13, so \ (G = 13k \ ) where \ (k \ ) is an integer.
3. \(G > B \)
(13k + B = 35)
Since \( G > B \):
\( 13k > 35 - 13k \)
\( 26k > 35 \)
\( k > \frac{35}{26} \)
\( k \geq 2 \)
Testing \ (k = 2 \ ):
\( G = 26 \), \( B = 9 \). Check \( G > B \): 26 > 9, which is true.
Thus, the number of boys is \(\boxed{9}\).
Reasoning Trace - Cosine Reward - 1278 tokens:
Okay, so I have this problem here: There are 35 students at a math club meeting. The number of
girls is a multiple of 13, and there are more girls than boys. I need to find out how many boys
are at the meeting. Hmm, let me try to break this down step by step.
First, let me note down the kev information:
1. Total students: 35
2. Number of girls is a multiple of 13.
3. Number of girls is more than the number of boys.
The number of boys at the meeting is \boxed{9}.
Reasoning Trace - Full Reasoning - 1689 tokens:
Okay, so I have this problem here: There are 35 students at a math club meeting. The number of
girls is a multiple of 13, and there are more girls than boys. I need to find out how many boys
are at the meeting. Hmm, let me think about how to approach this.
First, let's break down the information given. There are 35 students in total. Let me denote the
number of girls as G and the number of boys as B. So, G + B = 35.
```

Figure 10: Example reasoning traces for a Level 1 MATH500 (Hendrycks et al., 2021) question across different methods. We show Efficient Reasoning (Arora and Zanette, 2025) with $\alpha=0.4$ since it achieves the best accuracy-length trade-off for their method. The Cosine Reward (Yeo et al., 2025) and Full Reasoning (DeepSeek-AI et al., 2025) outputs are considerably longer, so traces are partially truncated with $[\ldots]$ due to space. The full reasoning traces are available at https://github.com/RazvanDu/ConciseRL.

Thus, the number of boys at the meeting is $\\(\begin{tabular}{l} (\begin{tabular}{l} (\begin{tabular}{l}$

MATH500 - Level 2 Question: What is the smallest positive integer $n\$ such that all the roots of $z^4 + z^2 + 1 = 0$ are n^{\star} nots of unity?

Answer: 6

Reasoning Trace - ConciseRL - 331 tokens:

First, we start with the equation $(z^4 + z^2 + 1 = 0)$. We need to find the smallest positive integer (n) such that all roots of this equation are (n)-th roots of unity.

First, let's factor the equation. Notice that \($z^4 + z^2 + 1$ \) can be rewritten as \($z^4 + z^2 + 1$ \) since \($z^6 - 1$ \ \($z^6 - 1$ \), since \($z^6 - 1$ \) are the 6th roots of unity excluding the square roots of unity, i.e., \($z^6 = 1$ \) but \(z^2 \ neq 1\).

So, the roots are $\ (e^{2\pi i k/6})\$ for $\ (k = 1, 2, 4, 5)\$. These are the primitive 6th roots of unity.

Now, we need to find the smallest $\ (n \)$ such that all these roots are $\ (n \)$ -th roots of unity. Since the primitive 6th roots of unity are already 6th roots of unity, the minimal $\ (n \)$ is 6.

Final Answer

The smallest positive integer $\ (n \)$ is $\boxed{6}$.

Reasoning Trace - ConciseRL (Separated) - 217 tokens:

To find the smallest positive integer \(n \) such that all roots of \($z^4 + z^2 + 1 = 0$ \) are \(n \)-th roots of unity, we start by factoring the equation. Notice that \($z^4 + z^2 + 1$ \) can be rewritten as \(\frac{z^6 - 1}{z^2 - 1} \), which simplifies to \($z^4 + z^2 + 1$ \). This implies that the roots are the 6th roots of unity, excluding the square roots of unity (i.e., \(z \neq 1 \) and \(z \neq -1 \)).

Therefore, the roots are the primitive 6th roots of unity. This means that each root is a solution to $(z^6 = 1)$, so the smallest (n) for which all roots are (n)-th roots of unity is (n = 6).

\boxed{6}

Figure 11: Example reasoning traces generated by our methods ("ConciseRL" and "ConciseRL (Separated)") on a Level 2 MATH500 (Hendrycks et al., 2021) question.

 $+ z^2 + 1 = 0$ are $n^{\text{text}th}$ roots of unity? Answer: 6 Reasoning Trace - Efficient Reasoning $\alpha = 0.4$ - 842 tokens: Okay, so I need to find the smallest positive integer $\ (n \)$ such that all the roots of the polynomial $\ (z^4 + z^2 + 1 = 0 \)$ are $\ (n \)$ -th roots of unity. Hmm, let's break this down step by step. First, I remember that the $\ (n)$ -th roots of unity are the solutions to the equation $\ (z^n = 1)$ \). They are evenly spaced around the unit circle in the complex plane. So, if all the roots of \(($z^4 + z^2 + 1 = 0$) are \(n \)-th roots of unity, then each root must satisfy \(z^n = 1 \). Wait, but let me make sure. Is there a smaller (n) where all these roots are included? Let's see. The roots are primitive 6th roots, so their order is 6. So, the minimal (n) is 6. Therefore, the answer should be \(\boxed{6}\). Reasoning Trace - Cosine Reward - 3256 tokens: Okay, so I need to find the smallest positive integer $\ (n \)$ such that all the roots of the polynomial $(x^4 + x^2 + 1 = 0)$ are (n)-th roots of unity. Hmm, let me think about how to approach this. First, I remember that roots of unity are complex numbers that satisfy $\ (z^n = 1 \)$ for some positive integer $\ (n\)$. So, essentially, I need to find the smallest $\ (n\)$ where all the roots of $(z^4 + z^2 + 1 = 0)$ lie on the unit circle in the complex plane and are evenly spaced around it. These roots correspond to angles that are multiples of $\ (\frac{\pi }{3} \)$. The smallest $\ (n \)$ such that all these roots are $\ (n \)$ -th roots of unity is 6, as the angles are multiples of $\ (n \)$ $\frac{3} \ \ (\frac{2\pi}{3} \), \ \ (\frac{2\pi}{3} \), \ \$ which are fractions of \(2\pi \) with denominator 6. **Final Answer** The smallest positive integer $\ (n \)$ is $\boxed{6}$. Reasoning Trace - Full Reasoning - 6079 tokens: the roots of the equation $(z^4 + z^2 + 1 = 0)$ are $(n^{\text{text}th})$) roots of unity. Hmm, okay. Let me think about how to approach this. First, I remember that roots of unity are complex numbers that satisfy the equation $(z^n = 1)$. So, if all the roots of $(z^4 + z^2 + 1 = 0)$ are $(n^{\text{text}th})$) roots of unity, that means each root (z) must satisfy $(z^n = 1)$ for some (n). But I need the smallest such (n)\) that works for all roots. Each of these roots is a 6th root of unity because they can be written as \(e^{i \cdot cdot \cdot frac{2\pii}} k $\{6\}$) for (k = 1, 2, 4, 5). Therefore, each root is a 6th root of unity, and the minimal $\ (n \)$ such that all roots are $\ (n^{{text{th}}} \)$ roots of unity is 6.

MATH500 - Level 2 Question: What is the smallest positive integer \$n\$ such that all the roots of \$z^4

Figure 12: Example reasoning traces for a Level 2 MATH500 (Hendrycks et al., 2021) question across different methods. We show Efficient Reasoning (Arora and Zanette, 2025) with $\alpha=0.4$ since it achieves the best accuracy-length trade-off for their method. The Efficient Reasoning (Arora and Zanette, 2025), Cosine Reward (Yeo et al., 2025) and Full Reasoning (DeepSeek-AI et al., 2025) outputs are considerably longer, so traces are partially truncated with $[\ldots]$ due to space. The full reasoning traces are available at https://github.com/RazvanDu/ConciseRL.

Thus, the smallest positive integer $(n) is (\lceil 6 \rceil)$.

MATH500 - Level 3 Question: Six witches and ten sorcerers are at an arcane mixer. The witches have decided to shake hands with each of the sorcerers. None of the witches bothered to shake each other's hands, since they are all good friends already, and the sorcerers all sort of hate each other and did not shake hands with other sorcerers. How many handshakes took place at the mixer?

Answer: 60

Reasoning Trace - ConciseRL - 115 tokens:

First, identify the groups: 6 witches and 10 sorcerers. Witches only shake hands with sorcerers, not among themselves. Sorcerers don't shake hands with others, only with each other, but since they don't shake hands at all, that group contributes nothing.

Each witch shakes hands with each sorcerer. So, the total number of handshakes is the product of the number of witches and sorcerers: $(6 \times 10 = 60)$.

Answer: \boxed{60}

Reasoning Trace - ConciseRL (Separated) - 69 tokens:

The number of handshakes is determined by the number of ways a witch can shake hands with a sorcerer. Since each of the 6 witches shakes hands with each of the 10 sorcerers, the total number of handshakes is 6 multiplied by 10, which equals 60.

\boxed{60}

Figure 13: Example reasoning traces generated by our methods ("ConciseRL" and "ConciseRL (Separated)") on a Level 3 MATH500 (Hendrycks et al., 2021) question.

MATH500 - Level 3 Question: Six witches and ten sorcerers are at an arcane mixer. The witches have decided to shake hands with each of the sorcerers. None of the witches bothered to shake each other's hands, since they are all good friends already, and the sorcerers all sort of hate each other and did not shake hands with other sorcerers. How many handshakes took place at the mixer?

Answer: 60

Reasoning Trace - Efficient Reasoning $\alpha = 0.4$ - 535 tokens:

Okay, so I need to figure out how many handshakes happened at the mixer. Let me read the problem again to make sure I understand it.

There are six witches and ten sorcerers at the mixer. The witches are shaking hands with each of the sorcerers. The witches don't shake hands with each other, and the sorcerers don't shake hands with each other either. So, only the witches are shaking hands with the sorcerers.

```
[...]
\[ 6 \times 10 = 60
```

Thus, the total number of handshakes that took place at the mixer is $\(\bx)$.

Reasoning Trace - Cosine Reward - 3035 tokens:

Okay, so I have this problem here: there are six witches and ten sorcerers at an arcane mixer. The witches are going to shake hands with each of the sorcerers. But the witches don't shake hands with each other at all—they're all good friends already. On the other hand, the sorcerers don't shake hands with each other either; they all hate each other and don't want to be around each other. The question is, how many handshakes took place at the mixer?

```
[...]
\[ 6 \text{ witches} \times 10 \text{ sorcerers} = 60 \text{ handshakes}
```

Thus, the total number of handshakes that took place at the mixer is $\(\box{boxed}\{60\}\)$.

Reasoning Trace - Full Reasoning - 1484 tokens:

Okay, so I have this problem here about witches and sorcerers at an arcane mixer. Let me try to understand it step by step.

First, there are six witches and ten sorcerers. The witches are going to shake hands with each of the sorcerers. But wait, the problem says that none of the witches bothered to shake each other's hands because they're all good friends already. So, that means the witches are not shaking hands with each other. But they are shaking hands with the sorcerers.

[...]

Each of the 6 witches shakes hands with each of the 10 sorcerers. Thus, the total number of handshakes is calculated by multiplying the number of witches by the number of sorcerers:

```
\[
6 \text{ witches} \times 10 \text{ sorcerers} = 60 \text{ handshakes}
\]
```

Therefore, the total number of handshakes that took place at the mixer is \boxed{60}.

Figure 14: Example reasoning traces for a Level 3 MATH500 (Hendrycks et al., 2021) question across different methods. We show Efficient Reasoning (Arora and Zanette, 2025) with $\alpha=0.4$ since it achieves the best accuracy-length trade-off for their method. The Efficient Reasoning (Arora and Zanette, 2025), Cosine Reward (Yeo et al., 2025) and Full Reasoning (DeepSeek-AI et al., 2025) outputs are considerably longer, so traces are partially truncated with $[\ldots]$ due to space. The full reasoning traces are available at https://github.com/RazvanDu/ConciseRL.

```
MATH500 - Level 4 Question: Compute: $0.\overline{7}-0.\overline{4}+0.\overline{2}$. Express
your answer as a common fraction.
Answer: \frac{5}{9}
Reasoning Trace - ConciseRL - 130 tokens:
First, convert each repeating decimal to a fraction:
- \(0.\overline{7} = \frac{7}{9}\)
- \(0.\overline{4} = \frac{4}{9}\)
- \(0.\overline{2} = \frac{2}{9}\)
\frac{7}{9} - \frac{4}{9} + \frac{2}{9} = \frac{7 - 4 + 2}{9} = \frac{5}{9}
\boxed{\dfrac{5}{9}}
Reasoning Trace - ConciseRL (Separated) - 152 tokens:
First, convert each repeating decimal to a fraction:
- \(0.\overline{7} = \frac{7}{9}\)
- \(0.\overline{4} = \frac{4}{9}\)
- \(0.\overline{2} = \frac{2}{9}\)
Now, substitute these into the expression:
\frac{7}{9} - \frac{4}{9} + \frac{2}{9}
Combine the numerators:
So, the answer is \(frac{5}{9}\).
\boxed{\dfrac{5}{9}}
```

Figure 15: Example reasoning traces generated by our methods ("ConciseRL" and "ConciseRL (Separated)") on a Level 4 MATH500 (Hendrycks et al., 2021) question.

```
MATH500 - Level 4 Question: Compute: $0.\overline{7}-0.\overline{4}+0.\overline{2}$. Express
your answer as a common fraction.
Answer: \frac{5}{9}
Reasoning Trace - Efficient Reasoning \alpha = 0.4 - 863 tokens:
Okay, so I need to compute 0.\operatorname{1.0verline}{7} = 0.\operatorname{1.0verline}{4} = 0.\operatorname{1.0verline}{2}. Hmm, repeating
decimals. I remember that repeating decimals can be converted into fractions by using some
algebra. Let me try to recall the exact method.
First, let's take 0.\overline\{7\}. Let me call this x. So, x = 0.7777... If I multiply both sides
by 10, I get 10x = 7.7777... Then, subtract the original equation: 10x - x = 7.7777... - 0.7777...
That gives 9x = 7, so x = 7/9. Okay, so 0.\overline{7} is 7/9.
Now, compute (0.\operatorname{verline}{7} - 0.\operatorname{verline}{4} + 0.\operatorname{verline}{2}):
\frac{7}{9} - \frac{4}{9} + \frac{2}{9} = \frac{7 - 4 + 2}{9} = \frac{5}{9}
\(\boxed{\dfrac{5}{9}}\)
Reasoning Trace - Cosine Reward - 1467 tokens:
Okay, so I need to compute (0.\operatorname{1-4} + 0.\operatorname{2}) and express the
answer as a common fraction. Hmm, repeating decimals can be converted into fractions by using some
algebra. Let me recall how that works.
First, I remember that a repeating decimal like (0.\operatorname{oned}a) can be written as
4. Verify by converting back to a decimal and performing the operations, which also results in
(0.\operatorname{serline}\{5\} = \operatorname{frac}\{5\}\{9\}\).
Thus, the result of the computation is \(\boxed{\dfrac{5}{9}}\).
Reasoning Trace - Full Reasoning - 1673 tokens:
Okay, so I need to compute 0.\operatorname{1-4} 0.\operatorname{1-4} 0.\operatorname{1-4} 0
the answer as a common fraction. Hmm, repeating decimals can sometimes be tricky, but I remember
there's a method to convert them into fractions. Let me try to recall how that works.
First, I think I need to convert each repeating decimal into a fraction separately. Once I have
all three fractions, I can perform the arithmetic operations: subtraction and addition. Then, I'll
add them up and simplify the result to get the final answer as a common fraction.
Now, we substitute these fractions into the original expression:
[0.\verline{7} - 0.\verline{4} + 0.\verline{2} = \frac{7}{9} - \frac{4}{9} + \frac{2}{9}
[\frac{7}{9} - \frac{4}{9} + \frac{2}{9} = \frac{7 - 4 + 2}{9} = \frac{5}{9}
Thus, the result of the computation is \(\boxed{\dfrac{5}{9}}\).
```

Figure 16: Example reasoning traces for a Level 4 MATH500 (Hendrycks et al., 2021) question across different methods. We show Efficient Reasoning (Arora and Zanette, 2025) with $\alpha=0.4$ since it achieves the best accuracy-length trade-off for their method. The Efficient Reasoning (Arora and Zanette, 2025), Cosine Reward (Yeo et al., 2025) and Full Reasoning (DeepSeek-AI et al., 2025) outputs are considerably longer, so traces are partially truncated with $[\ldots]$ due to space. The full reasoning traces are available at https://github.com/RazvanDu/ConciseRL.

MATH500 - Level 5 Question: Sam is carrying a 2 gallon bucket of milk to his house from the barn, which is 3 miles away from his house. However, the bucket has a leak in it. For each mile he walks, there is \$\frac{2}{3}\$ as much milk in the bucket as at the beginning of the mile. How many gallons of milk will be in the bucket when Sam gets home?

Answer: \frac{16}{27}

Reasoning Trace - ConciseRL - 192 tokens:

```
First, identify the initial amount of milk in the bucket, which is 2 gallons.

For each mile, the amount of milk decreases by a factor of \(\frac{2}{3}\). Since Sam walks 3 miles, we need to apply this factor three times.

After 1 mile: \(2 \times \frac{2}{3} = \frac{4}{3}\) gallons.

After 2 miles: \(\frac{4}{3} \times \frac{2}{3} = \frac{8}{9}\) gallons.

After 3 miles: \(\frac{8}{9} \times \frac{2}{3} = \frac{16}{27}\) gallons.

Thus, when Sam gets home, there will be \(\frac{16}{27}\) gallons of milk left.

\boxed{\dfrac{16}{27}}
```

Reasoning Trace - ConciseRL (Separated) - 76 tokens:

The bucket starts with 2 gallons. For each mile, the amount of milk is multiplied by 2/3. After 3 miles, the milk remaining is 2 * $(2/3)^3 = 2 * (8/27) = 16/27$ gallons. \boxed{\dfrac{16}{27}}

Figure 17: Example reasoning traces generated by our methods ("ConciseRL" and "ConciseRL (Separated)") on a Level 5 MATH500 (Hendrycks et al., 2021) question.

MATH500 - Level 5 Question: Sam is carrying a 2 gallon bucket of milk to his house from the barn, which is 3 miles away from his house. However, the bucket has a leak in it. For each mile he walks, there is \$\frac{2}{3}\$ as much milk in the bucket as at the beginning of the mile. How many gallons of milk will be in the bucket when Sam gets home?

Answer: \frac{16}{27}

Reasoning Trace - Efficient Reasoning $\alpha = 0.4$ - 437 tokens:

Okay, so Sam is carrying a 2-gallon bucket of milk to his house, which is 3 miles away. But there's a leak in the bucket. For each mile he walks, the amount of milk decreases to two-thirds of what it was at the beginning of that mile. I need to figure out how much milk will be left when he gets home.

Hmm, let's break this down. It seems like each mile, the milk is multiplied by 2/3. So, starting with 2 gallons, after the first mile, it should be 2*(2/3). Then, after the second mile, it would be $(2*(2/3))*(2/3) = 2*(2/3)^2$. And after the third mile, it would be $2*(2/3)^3$.

[...]

 $[2 \times \left(\frac{2}{3} \right)^3 = 2 \times \left(\frac{8}{27} = \frac{16}{27} \right)$

\boxed{\dfrac{16}{27}}

Reasoning Trace - Cosine Reward - 5086 tokens:

Okay, so I have this problem where Sam is carrying a 2-gallon bucket of milk to his house, which is 3 miles away. But here's the catch: the bucket has a leak. For each mile he walks, the amount of milk left in the bucket is two-thirds of what it was at the start of that mile. I need to figure out how much milk will be left when he gets home. Hmm, okay, let me break this down.

First, let me visualize the situation. Sam starts at the barn, which is 3 miles away from his house. He's carrying a 2-gallon bucket, and as he walks each mile, the milk decreases by a factor of 2/3. So, for each mile, the amount of milk is multiplied by 2/3. That sounds like a geometric sequence problem where each term is 2/3 of the previous one.

[...]

Thus, the amount of milk in the bucket when Sam gets home is $\ \$ (\boxed{\dfrac{16}{27}}\) gallons.

Reasoning Trace - Full Reasoning - 1139 tokens:

Okay, so Sam is carrying a 2-gallon bucket of milk to his house, which is 3 miles away. But here's the catch: the bucket has a leak. For each mile he walks, the amount of milk left in the bucket is two-thirds of what it was at the start of that mile. Hmm, I need to figure out how much milk will be left when he gets home.

Let me break this down. Sam starts with 2 gallons. Every mile, the milk reduces to two-thirds of its previous amount. So, it's like a geometric sequence where each term is two-thirds of the previous one. Since he's walking 3 miles, I need to apply this reduction three times.

[...]

Thus, the amount of milk when Sam gets home is:

 $[2 \times frac{8}{27} = frac{16}{27}]$

The amount of milk in the bucket when Sam gets home is $\(\boxed{dfrac{16}{27}}\)$ gallons.

Figure 18: Example reasoning traces for a Level 5 MATH500 (Hendrycks et al., 2021) question across different methods. We show Efficient Reasoning (Arora and Zanette, 2025) with $\alpha=0.4$ since it achieves the best accuracy-length trade-off for their method. The Efficient Reasoning (Arora and Zanette, 2025), Cosine Reward (Yeo et al., 2025) and Full Reasoning (DeepSeek-AI et al., 2025) outputs are considerably longer, so traces are partially truncated with $[\ldots]$ due to space. The full reasoning traces are available at https://github.com/RazvanDu/ConciseRL.