

Correct, Concise and Complete: Multi-stage Training For Adaptive Reasoning

Anonymous ACL submission

Abstract

The reasoning capabilities of large language models (LLMs) have improved substantially through increased test-time computation, typically in the form of intermediate tokens known as chain-of-thought (CoT). However, CoT often becomes unnecessarily long, increasing computation cost without actual accuracy gains or sometimes even degrading performance, a phenomenon known as “*overthinking*”. We propose a multi-stage efficient reasoning method that combines supervised fine-tuning—via rejection sampling or reasoning trace reformatting—with reinforcement learning using an adaptive length penalty. We introduce a lightweight reward function that penalizes tokens generated after the first correct answer but encouraging self-verification only when beneficial. We conduct a holistic evaluation across seven diverse reasoning tasks, analyzing the accuracy–response length trade-off. Our approach reduces response length by an average of 28% for 8B models and 40% for 32B models, while incurring only minor performance drops of 1.6 and 2.5 points, respectively. Despite its conceptual simplicity, it achieves a superior trade-off compared to more complex state-of-the-art efficient reasoning methods, scoring 76.6, in terms of the area under the Overthinking-Adjusted Accuracy curve (AUC_{OAA})—5 points above the base model and 2.5 points above the second-best approach.

1 Introduction

Large language models (LLMs) achieve stronger performance on reasoning-intensive tasks, such as math and code generation, by increasing test-time computation (Snell et al., 2025; OpenAI, 2024; DeepSeek-AI, 2025; OpenAI, 2025). Accuracy often improves as the model generates longer chains of thought (CoT). However, reasoning traces can also become unnecessarily long and often repetitive, yielding no additional gains, and in some cases even reducing accuracy, a phenomenon known as

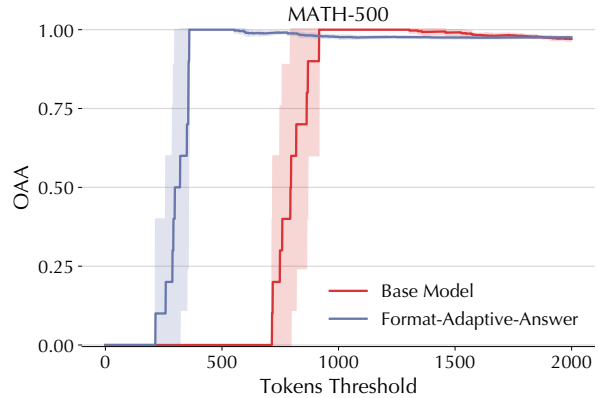


Figure 1: Overthinking-Adjusted Accuracy (OAA) (Aggarwal et al., 2025) as a function of the response length threshold on MATH-500 for Qwen3-8B. Our approach achieves similar accuracy with fewer tokens, leading to a larger area under the curve.

“*overthinking*” (Chen et al., 2025; Wu et al., 2025; Yang et al., 2025c).

To mitigate this, existing methods often impose a predefined thinking budget, truncating the reasoning trace once the budget is reached (Yang et al., 2025c) or enforcing it as a hard constraint during reinforcement learning (RL) training (Hou et al., 2025). However, such non-adaptive methods are unable to optimally balance accuracy and efficiency with respect to the response length (Snell et al., 2025; Wu et al., 2025; Yang et al., 2025c).

We introduce a multi-stage efficient reasoning framework that adaptively shortens response length while maintaining the base models’ accuracy. Our method consists of supervised fine-tuning (SFT), followed by reinforcement learning with a length penalty. We construct the training dataset for the SFT stage using two approaches: rejection sampling, selecting the shortest correct response for each problem, and reformatting reasoning traces to omit additional summaries and provide the final answer. For the RL stage, we design a reward function that penalizes tokens generated after the

067 first correct answer in the trace, encouraging con- 117
068 cise yet complete reasoning traces that lead to the 118
069 correct answer. It also incentivizes the model to 119
070 perform self-verification only when necessary, as 120
071 we show in our analysis. 121

072 We evaluate our methods on models of different 122
073 sizes from the Qwen3 and DeepSeek families, us- 123
074 ing a wide range of reasoning benchmarks, includ- 124
075 ing mathematics, science, code generation, ques- 125
076 tion answering, and long-context tasks. Our ap- 126
077 proach significantly reduces response length while 127
078 maintaining high accuracy. Furthermore, when 128
079 measured using the area under the Overthinking- 129
080 Adjusted Accuracy curve (OAA; Aggarwal et al. 130
081 (2025)), a unified metric that accounts for over- 131
082 thinking (see Figure 1), our methods consistently 132
083 improve both over the base models and state-of- 133
084 the-art efficient reasoning approaches. Our main 134
085 contributions are as follows: 135

- 086 • We propose a multi-stage efficient reasoning 136
087 framework combining SFT via rejection sam- 137
088 pling or via trace reformatting and RL with a 138
089 length-penalizing reward that penalizes tokens 139
090 generated after the first correct answer. This re- 140
091 duces response length by 28% for Qwen3-8B 141
092 with only a 1.6-point accuracy drop and 40% for 142
093 Qwen3-32B with a 2.5-point drop. 143
094 • We compare our approach with state-of-the-art 144
095 efficient reasoning methods and demonstrate con- 145
096 sistent improvements using the OAA curve, a 146
097 unified metric that accounts for overthinking. 147
098 • We analyze the trade-off between response length 148
099 and accuracy, and study how the trained models 149
100 adapt their chains of thought (CoT) for problems 150
101 of varying difficulty. 151

102 2 Methodology

103 To obtain optimal LLM reasoning traces, we pro- 152
104 pose a multi-stage training framework based on: 153
105 supervised fine-tuning followed by reinforcement 154
106 learning with an adaptive length penalty. This ap- 155
107 proach follows the paradigm originally used to train 156
108 reasoning LLMs (DeepSeek-AI, 2025).

109 **Supervised Fine-Tuning** This first stage serves 157
110 as a warm-up for RL training that also improves its 158
111 convergence. We construct our supervised training 159
112 datasets using the following approaches: 160

- 113 1. **Rejection sampling:** For each problem, we gen- 161
114 erate multiple continuations and select the short- 162
115 est correct one. While rejection sampling has 163

117 been explored in prior work as a baseline or 118
119 stand-alone method (Yang et al., 2025c), in con- 120
121 trast, we use it as the initial stage to bias the 122
123 model toward concise reasoning traces. 124

2. **Reformatting:** This approach modifies the for- 125
126 mat of model-generated reasoning traces. Rea- 127
128 soning models typically produce a structured 129
130 trace in which the intermediate reasoning (often 131
132 enclosed within `<think></think>`) is followed 133
134 by a summary and then the final answer. We 135
136 construct the dataset by removing the summary 137
138 and retaining only the final answer, encouraging 139
140 the model to generate direct solutions without 141
142 redundant reformulations. 143

131 **RL with Adaptive Length Penalty** After SFT, 132
133 we further improve efficiency through RL with an 134
135 adaptive length penalty. Specifically, we design a 136
137 verifiable reward function and use group relative 138
139 policy optimization (GRPO; (Shao et al., 2024)) 140
141 for training. In addition to the standard correctness 142
143 reward, we apply a length penalty to encourage 144
145 shorter, input-dependent reasoning traces, penaliz- 146
147 ing tokens generated after the first correct answer. 148
149 Prior methods truncate or prune traces at the token 150
151 or sentence level (Cui et al., 2025; Xia et al., 2025), 152
153 which can disrupt the reasoning flow. In contrast, 154
155 our reward function promotes responses that are 156
157 concise, complete, and correct. 158

145 The penalty is defined as the proportion of *to-* 146
147 *kens after the first correct answer* relative to the 148
149 full trace. Formally, let y denote the generated to- 149
150 ken sequence, y_{first} the subsequence up to the first 151
152 correct answer (empty if none is produced), and 153
154 L a function returning the number of tokens in a 154
155 sequence. The length penalty is: 156

$$152 R_L(y) = \begin{cases} \frac{L(y) - L(y_{\text{first}})}{L(y)} & \text{if the answer is correct,} \\ 0 & \text{otherwise,} \end{cases} \quad 152$$

153 where y_{first} denotes the prefix ending at the first 153
154 correct answer. Let $R_C(y)$ denote the correctness 154
155 and format reward. The overall reward is 155

$$156 R(y) = R_C(y) - \lambda R_L(y), \quad 156$$

157 where λ controls the trade-off between correct- 157
158 ness and reasoning efficiency; in our experiments, 158
159 we set $\lambda = 1$. 159

160 We locate the first correct answer using normal- 160
161 ized matching. If no correct answer is produced, 161
162 $y_{\text{first}} = \emptyset$, yielding zero length penalty. This dis- 162
163 courages redundant self-verification while allowing 163

self-correction: if the model initially produces an incorrect answer but later revises it correctly, no penalty is applied.

We refer to the method using rejection sampling during SFT as *Adaptive-Answer*, and the method using trace reformatting during SFT as *Format-Adaptive-Answer*.

3 Experimental Setup

Models. We use Qwen3-8B (Yang et al., 2025a) as the main model in our experiments. We further validate our method on Qwen3-1.7B, Qwen3-32B and DeepSeek-R1-Qwen-7B-distilled (DeepSeek-AI, 2025). Qwen3-32B was directly fine-tuned with reinforcement learning with verifiable reward, while the other models were trained via supervised fine-tuning on reasoning traces generated by a larger model.

Training Dataset. We train models on a sample of 13K problems from DeepScaleR (Luo et al., 2025), a collection of math datasets with problems drawn from AIME 1983-2023, AMC, Omni-Math (Gao et al., 2025), and STILL (Min et al., 2024).

Although training exclusively on math datasets, we evaluate our models across diverse domains, including science QA, commonsense reasoning, code generation, and long-context tasks. Our results show that the effects of our adaptive length penalty—reducing redundant self-verification and avoiding unnecessary continuation once correctness is reached—are domain-agnostic properties of reasoning traces. We see consistent reductions in response length with minimal accuracy loss across non-math tasks outside of our training domain.

Evaluation. We evaluate the models on a diverse set of datasets covering mathematics, coding, question answering, and long-context reasoning (details in Appendix A): AIME 24, AIME 25 (AIME, 2025), MATH-500 (Lightman et al., 2024), LiveCodeBenchv6 (Jain et al., 2024), GPQA-Diamond (Rein et al., 2024), LongBenchv2 (Bai et al., 2024), and CommonsenseQA (Talmor et al., 2019). We use the following decoding hyperparameters as recommended in Yang et al. (2025a) for the Qwen3 models: temperature = 0.7, top-p = 0.8, top-k = 20, and presence penalty = 1.5. The maximum number of output tokens is set to 32,768, except for MATH-500, AIME 24, and AIME 25, where it is set to 40,000. For each question, we sample N times and report the average accuracy as the final score, using

$N = 64$ for AIME 24 and AIME 25, and $N = 10$ for the remaining datasets.

Implementation Details. For rejection sampling, we generate 8 continuations for each problem. During the SFT stage, we train for 2 epochs with a batch size of 256 and a learning rate of 1e-5. Regarding the RL stage, we use GRPO (Shao et al., 2024) as implemented by the Verl framework (Sheng et al., 2024). We fine-tune the models with a group size of 8 and a global batch size of 256 for 50 iterations. We use the Adam optimizer with a learning rate of 1e-6, KL regularization with $\beta = 0.001$. For all experiments, including the baselines, we set the maximum number of output tokens to 16,384.

Metrics. We report both accuracy and response length (number of generated tokens) to characterize the performance-efficiency trade-off. Not all points can be directly compared using these two metrics. Therefore, we also report the area under the Overthinking-Adjusted Accuracy curve (AUC_{OAA} ; Aggarwal et al. (2025)). OAA_t measures the accuracy of the model when using fewer than t tokens:

$$OAA_t = \frac{1}{n} \sum_{i=1}^n (\text{Accuracy}_i \cdot \mathbb{I}[\#\text{Tokens}_i < t])$$

AUC_{OAA} is the area under the OAA_t curve, where the x-axis represents the token threshold t and the y-axis represents the corresponding OAA_t score, as illustrated in Figure 1.

$$AUC_{OAA} = \int_0^{t_{\max}} \frac{OAA_t}{t_{\max}} dt \approx \sum_0^{t_{\max}} \frac{OAA_t}{t_{\max}}$$

where t_{\max} is a predefined maximum number of tokens. Setting t_{\max} to a very large value is equivalent to using regular accuracy, which does not account for shorter traces. Therefore, for each dataset, we set t_{\max} to the mean number of tokens generated by the original base model.

Baselines. We compare our methods with existing state-of-the-art efficient reasoning approaches. We select a representative set of methods to cover a broad range of techniques:

- **No Thinking:** We disable thinking following the original Qwen3 paper (Yang et al., 2025a).
- **Supervised Fine-tuning (SFT):** For each problem in the training dataset, we generate 8 continuations and retain the shortest correct answer. We then fine-tune the model on the resulting dataset.

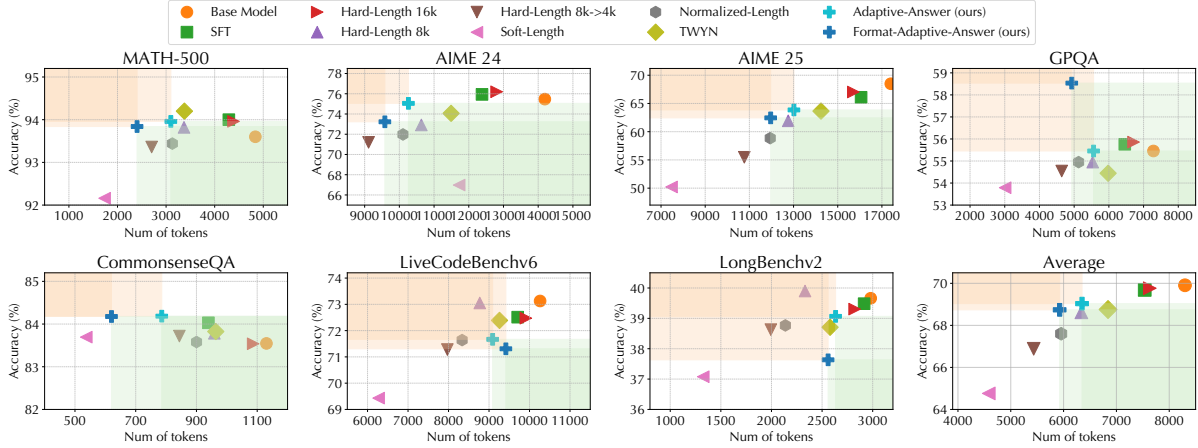


Figure 2: Average accuracy versus number of tokens for each method using Qwen3-8B. Points in the green region are dominated by *Adaptive-Answer* or *Format-Adaptive-Answer*, while points in the orange region dominate them (higher accuracy, fewer tokens).

- **RL with Hard Length Penalty (Hou et al., 2025):** Traces are truncated if they exceed a predefined maximum length. We set this threshold to 16k tokens, the maximum used in all our experiments, and 8k tokens, the average response length on the training set. We also report a curriculum variant that first trains with an 8k cutoff before lowering the threshold to 4k.
- **RL with Soft Length Penalty (Yu et al., 2025):** In addition to a hard cutoff L_{\max} , a second threshold L_{cache} introduces a gradually increasing penalty once the response length exceeds it. We set $L_{\max} = 10\text{k}$ and $L_{\text{cache}} = 8\text{k}$.
- **RL with Normalized Length Penalty (Team et al., 2025):** The length penalty is normalized using the minimum and maximum response lengths sampled within the same GRPO group.
- **RL with TWYN (Yang et al., 2025b):** Think When You Need (TWYN) is an adaptive method where rewards are based on pairwise comparisons: shorter correct responses receive higher rewards, while all incorrect responses receive equally low rewards.

4 Experimental Results

Response Length Reduction. Figure 2 shows accuracy as a function of response length across datasets when applying different efficient reasoning methods to Qwen3-8B (see Table 4 in Appendix B for absolute values). The green region indicates points dominated by *Adaptive-Answer*, while the orange region indicates points that dominate *Adaptive-Answer*.

We can see that our methods substantially re-

duce response length while maintaining accuracy on most datasets. The degree of reduction varies across tasks; however, even the less aggressive length-reduction variant, *Adaptive-Answer*, dominates other alternatives (i.e., there are almost no points in the orange area in the figures). More precisely, *Adaptive-Answer* dominates most methods on MATH-500, AIME 24, and CommonsenseQA, and is only dominated in two cases: (i) by *Hard-Length* and *Soft-Length* on LiveCodeBench, and (ii) by *Hard-Length* on LongBenchv2. *Format-Adaptive-Answer* dominates almost all other methods on the math and QA datasets, but is dominated on LiveCodeBench and LongBenchv2. We attribute the smaller reduction in response length across all efficient reasoning approaches on these two datasets to training primarily on math datasets.

In relative terms, the largest reductions—without any performance degradation—are observed on MATH-500 (36% for *Adaptive-Answer* and 50% for *Format-Adaptive-Answer*) and CommonsenseQA (30% and 45%, respectively). On GPQA Diamond, AIME 24, and AIME 25, response lengths decrease by about 25% for *Adaptive-Answer* and 32% for *Format-Adaptive-Answer*. For LiveCodeBench and LongBenchv2, both methods show only minor accuracy drops (less than two points) with smaller length reductions—11% and 12% for *Adaptive-Answer*, and 8% and 14% for *Format-Adaptive-Answer*. On average, *Adaptive-Answer* shortens responses by 28% and *Format-Adaptive-Answer* by 33%, with only a one-point decrease in accuracy. We must note, that even though training is performed only on math datasets,

	MATH-500	AIME 24	AIME 25	GPQA Diamond	Common- sense QA	LiveCode- Bench	Long- Benchv2	Avg.
Base model	81.7	68.6	75.8	<u>69.4</u>	72.1	94.3	39.4	71.6
No-Thinking	81.3	21.9	16.0	40.8	48.0	41.0	30.0	39.8
SFT	87.9	73.7	76.9	68.7	76.8	93.9	40.2	74.0
Hard-Length 16k	83.4	70.7	76.3	68.9	72.5	<u>94.1</u>	40.9	72.3
Hard-Length 8k	87.9	73.3	77.2	63.5	73.3	93.1	40.7	72.7
Hard-Length 8k \rightarrow 4k	89.1	75.1	69.8	61.0	76.5	92.7	38.9	71.8
Soft-Length	87.9	72.4	72.9	62.9	73.3	93.0	39.4	71.6
Normalized-Length	91.9	77.4	61.0	58.3	<u>78.6</u>	90.0	36.6	70.5
TWYN	89.6	74.5	<u>79.6</u>	67.1	74.8	93.9	38.8	74.0
<i>Adaptive-Answer</i> (ours)	90.3	75.8	80.0	68.8	81.4	93.1	39.4	<u>75.5</u>
<i>Format-Adaptive-Answer</i> (ours)	<u>91.2</u>	81.8	80.0	71.6	81.3	93.6	37.4	76.6

Table 1: AUC_{OAA} of all approaches applied to Qwen3-8B across datasets. On average, *Format-Adaptive-Answer* achieves the best performance, followed by *Adaptive-Answer*.

	Accuracy	#Tokens	AUC_{OAA}
Base Model	69.9	8295	71.6
SFT (Rejection Sampling)	69.7	7536	74.0
SFT (Formatting)	70.1	7475	74.7
RL (no SFT)	68.8	6303	73.2
<i>Adaptive-Answer</i>	69.0	6344	75.5
<i>Format-Adaptive-Answer</i>	68.7	5918	76.6

Table 2: Average accuracy, response length, and AUC_{OAA} of the original model, rejection sampling-based SFT, format-based SFT, RL with adaptive length penalty (without SFT), rejection sampling-based SFT followed by RL (*Adaptive-Answer*), and format-based SFT followed by RL (*Format-Adaptive-Answer*).

our methods also shorten responses on science, coding, QA, and long-context reasoning datasets.

We must highlight that not all methods are directly comparable using accuracy and response length. For example, on all datasets except AIME 24, *Soft-Length* neither dominates nor is dominated by other methods. Therefore, we also compare the AUC_{OAA} of all approaches applied to Qwen3-8B across datasets (see Table 1). This confirms the effectiveness of our methods: on average, *Format-Adaptive-Answer* outperforms all other methods, followed by *Adaptive-Answer*. *Format-Adaptive-Answer* achieves the highest score on all math and question answering datasets, except for MATH-500 where it ranks second. Interestingly, on LiveCodeBench and LongBenchv2, simple baselines such as *SFT* and *Hard-Length 16k*—equivalent to RL without a length penalty—outperform all other efficient reasoning alternatives.

Component Ablations. To evaluate the contribution of each component of our approach, we

perform ablations for each component individually. Table 2 reports the average accuracy, response length, and AUC_{OAA} for several configurations: the original model, rejection sampling-based SFT, format-based SFT, RL with adaptive length penalty (without SFT), rejection sampling-based SFT followed by RL (*Adaptive-Answer*), and format-based SFT followed by RL (*Format-Adaptive-Answer*).

Adding the rejection sampling-based SFT (*SFT (Rejection Sampling)*) stage does not yield clear improvements when examining accuracy or response length alone—*Adaptive-Answer* achieves slightly higher accuracy than *RL (no SFT)* but produces longer responses. Hence, we argue that AUC_{OAA} is a more suitable metric for ranking models when no model clearly dominates another. AUC_{OAA} clearly highlights the benefit of the SFT phase. We observe a sizable improvement of 3 AUC points over the base model.

Format-based SFT (*SFT (Formatting)*) reduces response length by 10% on average without loss in accuracy. This suggests that the summary generated before the final answer does not contribute to performance, as the model reaches the correct answer by the end of the trace. Adding the RL stage further improves AUC_{OAA} by 3 points and cuts response length 18%, with only a minor drop in accuracy. However, combining rejection sampling with formatting during SFT does not yield additional improvements over formatting alone.

Finally, SFT is a crucial stage for RL performance: although *RL (no SFT)* produces the second-shortest reasoning traces, it incurs a performance penalty and achieves a lower AUC_{OAA} compared to the full approaches (*Adaptive-Answer* and *Format-Adaptive-Answer*).

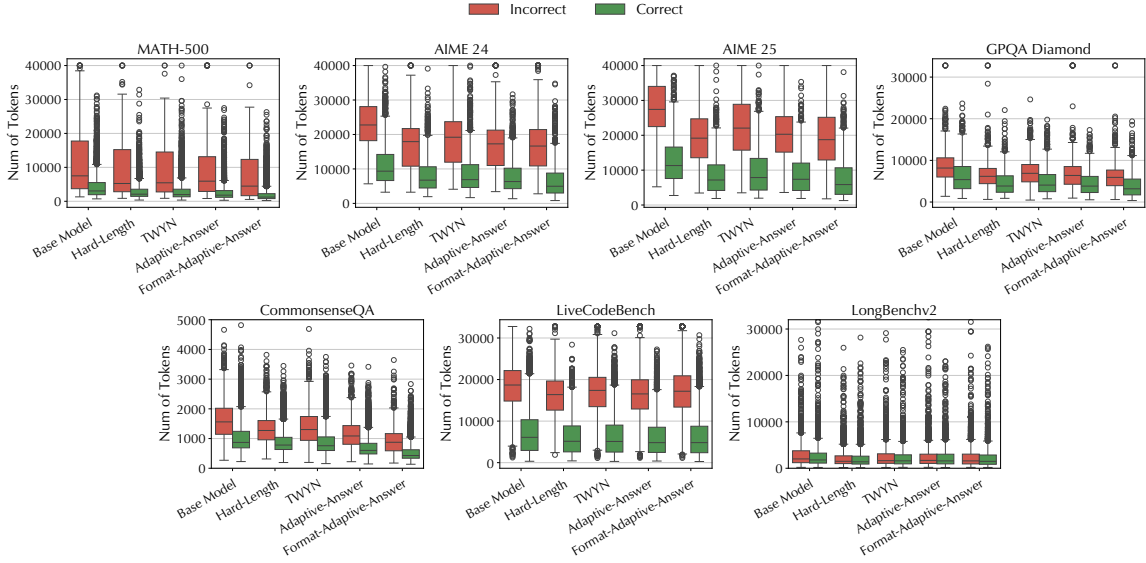


Figure 3: Response length distributions of some representative efficient reasoning methods applied to Qwen3-8B. We separate the correct and incorrect responses.

	Accuracy	#Tokens	AUC _{OAA}
Qwen3-1.7B			
Base Model	50.9	7,619	65.4
Adaptive-Answer	49.1	5,884 (-22%)	62.1
Format-Adaptive-Answer	48.3	5,918 (-22%)	62.1
Qwen3-8B			
Base Model	69.9	8,298	71.6
Adaptive-Answer	69.0	6,344 (-23%)	75.5
Format-Adaptive-Answer	68.7	5,918 (-28%)	76.6
Qwen3-32B			
Base Model	74.8	7,294	69.2
Adaptive-Answer	72.1	4,280 (-41%)	72.5
Format-Adaptive-Answer	72.2	4,372 (-40%)	72.3
DeepSeek-R1-Qwen-7B-Distill			
Base Model	50.3	6,272	62.1
Adaptive-Answer	50.4	5,133 (18%)	59.6
Format-Adaptive-Answer	50.7	4,612 (26%)	59.7

Table 3: Average accuracy, response length and AUC_{OAA} of our methods applied to Qwen3.1.7B, Qwen3-8B, Qwen-32B and DeepSeek-R1-Qwen-7B-distilled.

Efficiency and Model Size. We investigate how our approaches scale with model size and training regimes. For this set of experiments, we evaluate Qwen3- $\{1.7B, 8B, 32B\}$, and DeepSeek-R1-Qwen-7B-Distilled.

Across all models, the performance drop after applying our methods remains under 2.5 points—and is negligible for DeepSeek-R1-7B (see Table 3). Notably, the reduction in generated tokens increases with model size: *Adaptive-Answer* shortens responses by 22% on Qwen3-1.7B, 23% on Qwen3-

8B, and 40% on Qwen3-32B, demonstrating that larger models benefit more from efficient reasoning. However, AUC_{OAA} does not always align perfectly with the accuracy–response length trade-offs, particularly for smaller models, highlighting that efficiency gains can sometimes come at a subtle cost to overall reasoning effectiveness. For instance, the fine-tuned DeepSeek-R1 dominates the base model in absolute accuracy but achieves a slightly lower AUC_{OAA} score. Overall, these results indicate that our methods are model-agnostic, consistently effective across different model sizes, and operating without significant performance loss.

5 Analysis

Response Length Distribution. Figure 3 shows the response length distributions for a representative set of efficient reasoning methods applied to Qwen3-8B. Across all three datasets, incorrect answers tend to have longer traces than correct ones, highlighting a correlation between excessive reasoning and errors. Importantly, our methods effectively shift the response length distribution for both correct and incorrect answers, showing that the models adapt traces consistently, regardless of the final answer. This indicates that our approach encourages concise reasoning across all outputs, not just the correct ones.

Intermediate Answers. The RL stage encourages the model to minimize unnecessary self-verification. To analyze this, we report the num-

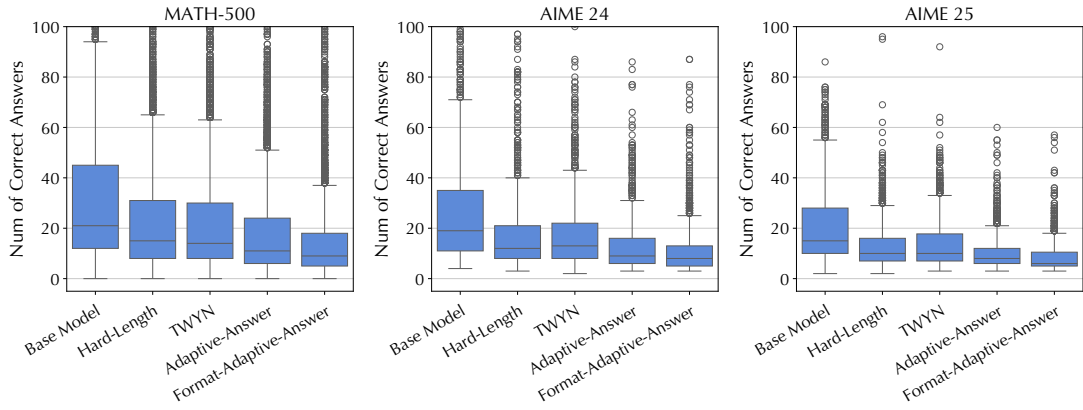


Figure 4: Distributions of the number of correct answers in the traces of some representative efficient reasoning methods applied to Qwen3-8B for MATH-500, AIME 24 and AIME 25.

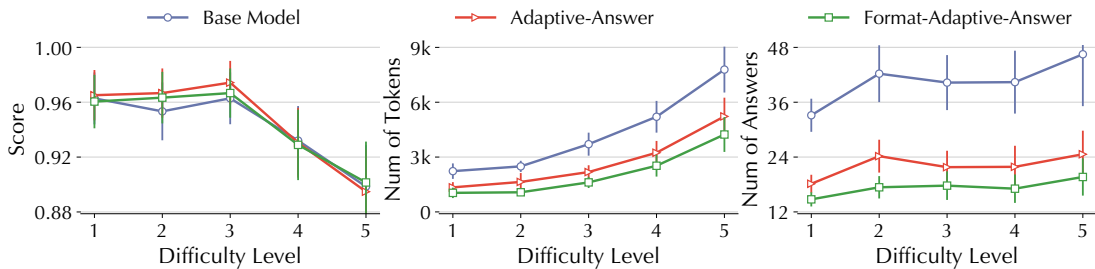


Figure 5: Accuracy, response length, and count of intermediate correct steps across difficulty levels on MATH-500.

424 ber of correct answers appearing in each reasoning
 425 trace for MATH-500, AIME 24, and AIME 25 (Fig-
 426 ure 4). While this metric is a coarse proxy—since
 427 answers may be repeated or paraphrased—it pro-
 428 vides qualitative insight into verification behav-
 429 ior. Both *Adaptive-Answer* and *Format-Adaptive-*
 430 *Answer* shift the distribution toward fewer inter-
 431 mediate correct answers, indicating reduced redun-
 432 dancy in reasoning.

433 **Difficulty Analysis.** We examine how problem
 434 difficulty affects accuracy, response length, and
 435 intermediate correct answers. Each MATH-500
 436 problem is assigned a difficulty level from 1 to 5.
 437 As shown in Figure 5, our approaches maintain the
 438 base model’s accuracy across all levels. Response
 439 length adapts to difficulty, increasing for harder
 440 problems, reflecting the need for more reasoning.
 441 Although both response length and intermediate
 442 correct answers rise with difficulty, they remain
 443 shorter than the base model, demonstrating our
 444 methods’ efficiency even on challenging problems.

445 **Qualitative Analysis.** Figure 6 compares reason-
 446 ing traces for a math problem from AIME24 pro-
 447 duced by: (i) the base model, (ii) *Adaptive-Answer*,
 448 and (iii) *Format-Adaptive-Answer* (long responses

449 are trimmed). We see that the base model per-
 450 forms seven unnecessary self-verifications after
 451 producing the first correct answer (116). In con-
 452 trast, *Adaptive-Answer* reduces this to three, while
 453 *Format-Adaptive-Answer* produces an optimal trace
 454 with no self-verifications, directly generating the
 455 final answer without summarizing the reasoning.

456 6 Related Work

457 **Test-Time Scaling.** Large Language Models per-
 458 form better on reasoning-heavy tasks such as math,
 459 problem-solving, and coding by increasing test-
 460 time computation (Wei et al., 2022; Wang et al.,
 461 2023; Snell et al., 2025). Models generate inter-
 462 mediate tokens in parallel—by sampling multiple
 463 traces (Wang et al., 2023)—or sequentially, by ver-
 464 ifying and correcting their own outputs (Madaan
 465 et al., 2023; Kumar et al., 2024). Recent works
 466 use reinforcement learning with verifiable rewards
 467 to further enhance reasoning capabilities, leading
 468 to longer CoT as well as self-verification and self-
 469 correction behaviors (OpenAI, 2024; DeepSeek-AI,
 470 2025; OpenAI, 2025; Yang et al., 2025a).

471 **Efficient Reasoning.** While reinforcement learn-
 472 ing improves reasoning ability, it often comes at the
 473 cost of efficiency. In some cases, reasoning traces

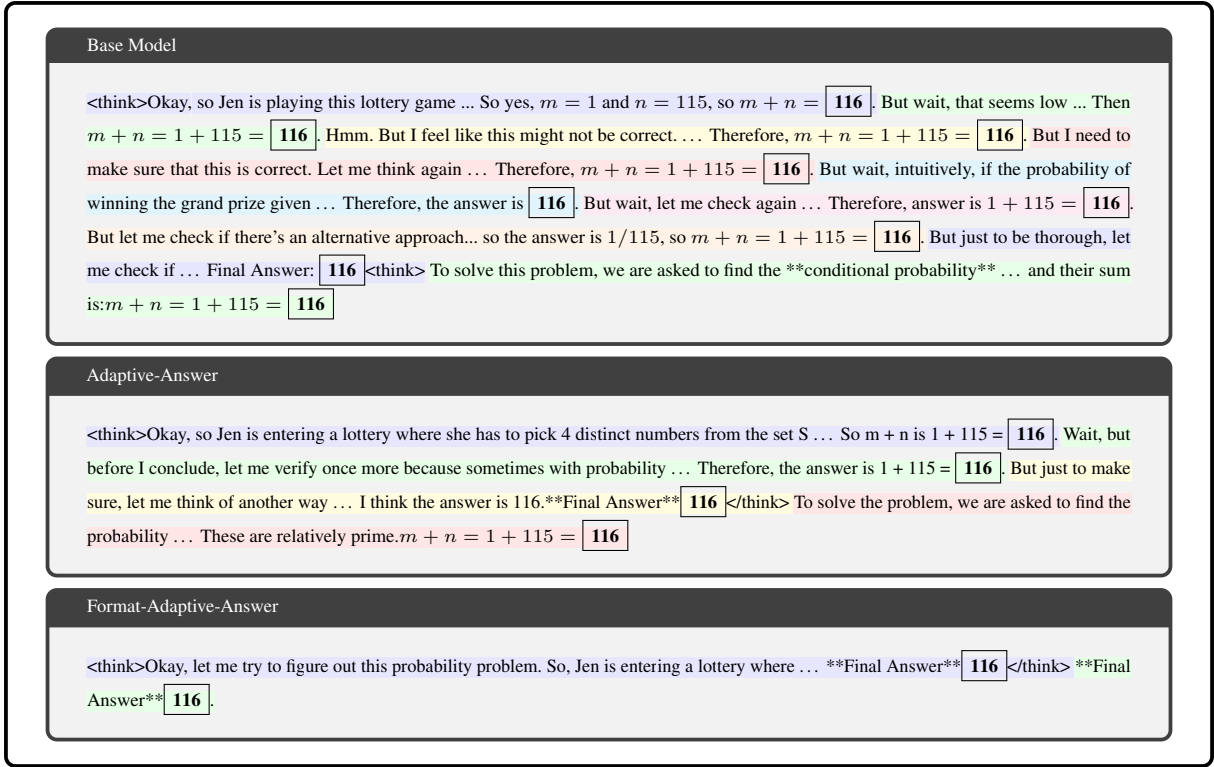


Figure 6: Reasoning traces of Qwen3-8B on AIME 24 Problem 10 before and after fine-tuning. The base model performs seven self-verifications after arriving at the correct answer, whereas *Adaptive-Answer* performs only two and *Format-Adaptive-Answer* performs none.

474 become excessively long, increasing computation
 475 without improving accuracy—and sometimes even
 476 harming it, a phenomenon known as “overthinking”
 477 (Chen et al., 2025; Yang et al., 2025c; Wu et al.,
 478 2025). Several methods have been proposed to
 479 address this issue. The most direct approach, bud-
 480 get forcing (Muennighoff et al., 2025; Yang et al.,
 481 2025a), interrupts generation once a predefined
 482 threshold is exceeded. Other methods (Yang et al.,
 483 2025c; Xia et al., 2025; Cui et al., 2025) construct
 484 synthetic datasets by shortening model-generated
 485 reasoning traces (via rejection sampling or prun-
 486 ing) and then perform supervised fine-tuning. A
 487 different line of work, to which our method be-
 488 longs, uses reinforcement learning with a length
 489 penalty in addition to the correctness reward (Lou
 490 et al., 2025; Zhang et al., 2025). The length con-
 491 straint can be either hard, applied once the CoT
 492 length exceeds a fixed threshold (Hou et al., 2025),
 493 or soft, where the penalty increases gradually as
 494 the trace length approaches the threshold (Yu et al.,
 495 2025; Aggarwal and Welleck, 2025). Instead of ap-
 496 plying penalties independently per example, some
 497 approaches (Team et al., 2025; Yang et al., 2025b)
 498 define them relative to the length and correctness

of other traces within the same GRPO group.

Unlike prior methods, which rely on a manually
 fixed “thinking budget” shared across all inputs, we
 train models to produce short yet complete reason-
 ing traces while preserving accuracy. Our approach
 incentivizes the model to adaptively infer an input-
 dependent budget.

7 Conclusion

Large Language Models (LLMs) often perform
 better on reasoning-intensive tasks by producing
 longer chains of thought. However, these chains
 are often unnecessarily long, increasing inference
 costs without improving accuracy. To address this,
 we propose a multi-stage efficient reasoning frame-
 work that consists of supervised fine-tuning—via
 rejection sampling or reformatting—followed by
 reinforcement learning with an adaptive length
 penalty. Our approach effectively shortens re-
 sponse length (28% for Qwen3-8B and 40% for
 Qwen3-32B) with only minor performance drops
 (up to 2.5 points accuracy) and outperforms exist-
 ing state-of-the-art efficient reasoning methods by
 2.5 points when evaluated with the unified metric
 AUC_{OAA} .

523 Limitations

524 Although we evaluate our methods on datasets from
525 multiple domains, our training is performed exclu-
526 sively on math datasets. Extending training to a
527 more diverse set of tasks could yield a better ac-
528 curacy–response length trade-off. Moreover, the
529 efficient reasoning methods we propose are post
530 hoc interventions; we do not explore incorporating
531 the adaptive length penalty directly during the ini-
532 tial RL training. Additionally, due to resourcing
533 constraints, we focused our experiment scope to
534 models of different sizes within the Qwen family
535 based on a dense architecture. Future explorations
536 can consider extending our approach to more model
537 families and other architectures such as Mixture-
538 of-Experts. Finally, we focus exclusively on the
539 model performance on reasoning tasks, and we do
540 not measure the change in performance on other
541 task groups that are less dependent on the CoT
542 quality.

543 Ethics Statement

544 One of our methods removes the summary that
545 the model produces at the end of its thinking con-
546 tent. Although this results in shorter responses, this
547 might also reduce the legibility of the reasoning
548 traces.

549 References

550 Pranjal Aggarwal, Seungone Kim, Jack Lanchantin,
551 Sean Welleck, Jason E Weston, Ilya Kulikov, and
552 Swarnadeep Saha. 2025. [Optimalthinkingbench:](#)
553 [Evaluating over and underthinking in LLMs.](#) In
554 *NeurIPS 2025 Workshop on Efficient Reasoning.*

555 Pranjal Aggarwal and Sean Welleck. 2025. [L1: Con-](#)
556 [trolling how long a reasoning model thinks with rein-](#)
557 [forcement learning.](#) *ArXiv preprint*, abs/2503.04697.

558 AIME. 2025. AIME problems and solutions.
559 [https://artofproblemsolving.com/wiki/](https://artofproblemsolving.com/wiki/index.php/AIME_Problems_and_Solutions)
560 [index.php/AIME_Problems_and_Solutions.](https://artofproblemsolving.com/wiki/index.php/AIME_Problems_and_Solutions)

561 Yushi Bai, Shangqing Tu, Jiajie Zhang, Hao Peng, Xi-
562 aozhi Wang, Xin Lv, Shulin Cao, Jiazheng Xu, Lei
563 Hou, Yuxiao Dong, Jie Tang, and Juanzi Li. 2024.
564 [Longbench v2: Towards deeper understanding and](#)
565 [reasoning on realistic long-context multitasks.](#) *ArXiv*
566 *preprint*, abs/2412.15204.

567 Xingyu Chen, Jiahao Xu, Tian Liang, Zhiwei He,
568 Jianhui Pang, Dian Yu, Linfeng Song, Qiuzhi
569 Liu, Mengfei Zhou, Zhuosheng Zhang, Rui Wang,
570 Zhaopeng Tu, Haitao Mi, and Dong Yu. 2025. [Do](#)
571 [NOT think that much for 2+3=? on the overthinking](#)

[of long reasoning models.](#) In *Forty-second Interna-*
tional Conference on Machine Learning. 572
573

574 Yingqian Cui, Pengfei He, Jingying Zeng, Hui Liu,
575 Xianfeng Tang, Zhenwei Dai, Yan Han, Chen Luo,
576 Jing Huang, Zhen Li, and 1 others. 2025. [Stepwise](#)
577 [perplexity-guided refinement for efficient chain-of-](#)
578 [thought reasoning in large language models.](#) *ArXiv*
579 *preprint*, abs/2502.13260.

580 DeepSeek-AI. 2025. [Deepseek-r1: Incentivizing rea-](#)
581 [soning capability in llms via reinforcement learning.](#)
582 *ArXiv preprint*, abs/2501.12948.

583 Bofei Gao, Feifan Song, Zhe Yang, Zefan Cai, Yibo
584 Miao, Qingxiu Dong, Lei Li, Chenghao Ma, Liang
585 Chen, Runxin Xu, Zhengyang Tang, Benyou Wang,
586 Daoguang Zan, Shanghaoran Quan, Ge Zhang, Lei
587 Sha, Yichang Zhang, Xuancheng Ren, Tianyu Liu,
588 and Baobao Chang. 2025. [Omni-MATH: A univer-](#)
589 [sal olympiad level mathematic benchmark for large](#)
590 [language models.](#) In *The Thirteenth International*
591 *Conference on Learning Representations.*

592 Bairu Hou, Yang Zhang, Jiabao Ji, Yujian Liu,
593 Kaizhi Qian, Jacob Andreas, and Shiyu Chang.
594 2025. [Thinkprune: Pruning long chain-of-thought](#)
595 [of llms via reinforcement learning.](#) *ArXiv preprint*,
596 abs/2504.01296.

597 Naman Jain, King Han, Alex Gu, Wen-Ding Li, Fanjia
598 Yan, Tianjun Zhang, Sida Wang, Armando Solar-
599 Lezama, Koushik Sen, and Ion Stoica. 2024. [Live-](#)
600 [codebench: Holistic and contamination free eval-](#)
601 [uation of large language models for code.](#) *ArXiv*
602 *preprint*, abs/2403.07974.

603 Aviral Kumar, Vincent Zhuang, Rishabh Agarwal,
604 Yi Su, John D Co-Reyes, Avi Singh, Kate Baumli,
605 Shariq Iqbal, Colton Bishop, Rebecca Roelofs, and
606 1 others. 2024. [Training language models to self-](#)
607 [correct via reinforcement learning.](#) *ArXiv preprint*,
608 abs/2409.12917.

609 Hunter Lightman, Vineet Kosaraju, Yuri Burda, Harri-
610 son Edwards, Bowen Baker, Teddy Lee, Jan Leike,
611 John Schulman, Ilya Sutskever, and Karl Cobbe.
612 2024. [Let’s verify step by step.](#) In *The Twelfth In-*
613 *ternational Conference on Learning Representations,*
614 *ICLR 2024, Vienna, Austria, May 7-11, 2024.* Open-
615 Review.net.

616 Chenwei Lou, Zewei Sun, Xinnian Liang, Meng Qu,
617 Wei Shen, Wenqi Wang, Yuntao Li, Qingping Yang,
618 and Shuangzhi Wu. 2025. [Adacot: Pareto-optimal](#)
619 [adaptive chain-of-thought triggering via reinfor-](#)
620 [cement learning.](#) *arXiv preprint arXiv:2505.11896.*

621 Michael Luo, Sijun Tan, Justin Wong, Xiaoxiang Shi,
622 William Y. Tang, Manan Roongta, Colin Cai, Jeffrey
623 Luo, Li Erran Li, Raluca Ada Popa, and Ion Stoica.
624 2025. [Deepscaler: Surpassing o1-preview with a](#)
625 [1.5b model by scaling rl.](#) *Notion Blog.*

626 Aman Madaan, Niket Tandon, Prakhar Gupta, Skyler
627 Hallinan, Luyu Gao, Sarah Wiegrefe, Uri Alon,

628	Nouha Dziri, Shrimai Prabhunoye, Yiming Yang, Shashank Gupta, Bodhisattwa Prasad Majumder, Katherine Hermann, Sean Welleck, Amir Yazdanbakhsh, and Peter Clark. 2023. Self-refine: Iterative refinement with self-feedback . In <i>Advances in Neural Information Processing Systems 36: Annual Conference on Neural Information Processing Systems 2023, NeurIPS 2023, New Orleans, LA, USA, December 10 - 16, 2023</i> .	4149–4158, Minneapolis, Minnesota. Association for Computational Linguistics.	685 686
637	Yingqian Min, Zhipeng Chen, Jinhao Jiang, Jie Chen, Jia Deng, Yiwen Hu, Yiru Tang, Jiapeng Wang, Xiaoxue Cheng, Huatong Song, and 1 others. 2024. Imitate, explore, and self-improve: A reproduction report on slow-thinking reasoning systems . <i>ArXiv preprint</i> , abs/2412.09413.	Kimi Team, Angang Du, Bofei Gao, Bowei Xing, Changjiu Jiang, Cheng Chen, Cheng Li, Chenjun Xiao, Chenzhuang Du, Chonghua Liao, and 1 others. 2025. Kimi k1. 5: Scaling reinforcement learning with llms . <i>ArXiv preprint</i> , abs/2501.12599.	687 688 689 690 691
643	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 <i>Proceedings of the 2025 Conference on Empirical Methods in Natural Language Processing</i> , pages 20275–20321, Suzhou, China. Association for Computational Linguistics.	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 <i>The Eleventh International Conference on Learning Representations, ICLR 2023, Kigali, Rwanda, May 1-5, 2023</i> . OpenReview.net.	692 693 694 695 696 697 698
651	OpenAI. 2024. Openai o1 system card . <i>ArXiv preprint</i> , abs/2412.16720.	Jason Wei, Xuezhi Wang, Dale Schuurmans, Maarten Bosma, Brian Ichter, Fei Xia, Ed H. Chi, Quoc V. Le, and Denny Zhou. 2022. Chain-of-thought prompting elicits reasoning in large language models . In <i>Advances in Neural Information Processing Systems 35: Annual Conference on Neural Information Processing Systems 2022, NeurIPS 2022, New Orleans, LA, USA, November 28 - December 9, 2022</i> .	699 700 701 702 703 704 705 706
653	OpenAI. 2025. Openai o3 and o4-mini system card . https://cdn.openai.com/pdf/2221c875-02dc-4789-800b-e7758f3722c1/o3-and-o4-mini-system-card.pdf .	Yuyang Wu, Yifei Wang, Ziyu Ye, Tianqi Du, Stefanie Jegelka, and Yisen Wang. 2025. When more is less: Understanding chain-of-thought length in llms . <i>ArXiv preprint</i> , abs/2502.07266.	707 708 709 710
657	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 <i>First Conference on Language Modeling</i> .	Heming Xia, Chak Tou Leong, Wenjie Wang, Yongqi Li, and Wenjie Li. 2025. Tokenskip: Controllable chain-of-thought compression in llms . <i>ArXiv preprint</i> , abs/2502.12067.	711 712 713 714
662	Zhihong Shao, Peiyi Wang, Qihao Zhu, Runxin Xu, Junxiao Song, Xiao Bi, Haowei Zhang, Mingchuan Zhang, YK Li, Yang Wu, and 1 others. 2024. Deepseekmath: Pushing the limits of mathematical reasoning in open language models . <i>ArXiv preprint</i> , abs/2402.03300.	An Yang, Anfeng Li, Baosong Yang, Beichen Zhang, Binyuan Hui, Bo Zheng, Bowen Yu, Chang Gao, Chengen Huang, Chenxu Lv, and 1 others. 2025a. Qwen3 technical report . <i>ArXiv preprint</i> , abs/2505.09388.	715 716 717 718 719
668	Guangming Sheng, Chi Zhang, Zilingfeng Ye, Xibin Wu, Wang Zhang, Ru Zhang, Yanghua Peng, Haibin Lin, and Chuan Wu. 2024. Hybridflow: A flexible and efficient rlhf framework . <i>arXiv preprint arXiv:2409.19256</i> .	Junjie Yang, Ke Lin, and Xing Yu. 2025b. Think when you need: Self-adaptive chain-of-thought learning . <i>ArXiv preprint</i> , abs/2504.03234.	720 721 722
673	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 <i>The Thirteenth International Conference on Learning Representations</i> .	Wenkai Yang, Shuming Ma, Yankai Lin, and Furu Wei. 2025c. Towards thinking-optimal scaling of test-time compute for llm reasoning . <i>ArXiv preprint</i> , abs/2502.18080.	723 724 725 726
678	Alon Talmor, Jonathan Herzig, Nicholas Lourie, and Jonathan Berant. 2019. CommonsenseQA: A question answering challenge targeting commonsense knowledge . In <i>Proceedings of the 2019 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies, Volume 1 (Long and Short Papers)</i> , pages	Qiyang Yu, Zheng Zhang, Ruofei Zhu, Yufeng Yuan, Xiaochen Zuo, Yu Yue, Weinan Dai, Tiantian Fan, Gaohong Liu, Lingjun Liu, and 1 others. 2025. Dapo: An open-source llm reinforcement learning system at scale . <i>ArXiv preprint</i> , abs/2503.14476.	727 728 729 730 731
684		Jiajie Zhang, Nianyi Lin, Lei Hou, Ling Feng, and Juanzi Li. 2025. Adaptthink: Reasoning models can learn when to think . <i>arXiv preprint arXiv:2505.13417</i> .	732 733 734 735

A Datasets

We evaluate on the following datasets that come from diverse domains including math, science, coding, question answering and long context reasoning:

MATH-500 (Lightman et al., 2024): A representative subset of 500 problems from the MATH benchmark. Each problem is assigned a difficulty level ranging from 1 to 5.

AIME 24 (AIME, 2025): 30 math problems from the 2024 edition of the American Invitational Mathematics Examination, a prestigious high school mathematics competition known for its challenging mathematical problems.

AIME 25 (AIME, 2025): 30 math problems from the 2025 edition of the American Invitational Mathematics Examination, a prestigious high school mathematics competition known for its challenging mathematical problems.

GPQA Diamond (Rein et al., 2024): a subset of 198 expert-written, graduate-level questions in biology, physics, and chemistry, designed to test the true reasoning abilities of Large Language Models (LLMs) without reliance on easily found internet answers

CommonsenseQA (Talmor et al., 2019): a dataset consisting of 1221 multiple-choice questions that require commonsense knowledge to predict the correct answers. Each question has one correct answer and four distractor answers.

LiveCodeBench (Jain et al., 2024): A holistic and contamination-free benchmark to evaluate the coding capabilities of LLMs. We use the sixth version of the dataset which contains 055 problems.

LongBenchv2 (Bai et al., 2024): a dataset of 503 challenging multiple-choice questions, with contexts ranging from 8k to 2M words. It contains the following categories: single-document QA, multi-document QA, long in-context learning, long-dialogue history understanding, code repo understanding, and long structured data understanding.

B Accuracy-Response Length Trade-off

Table 4 quantifies the accuracy-length trade-off. Tight hard-length constraints reduce average response length from 8.3k tokens (Base Model) to

5.4k, but incur a 3.0-point average accuracy drop and a severe degradation on AIME 25 (68.5 \rightarrow 55.5). Normalized-Length achieves the shortest outputs (4.6k tokens on average) but suffers the largest performance loss (69.9 \rightarrow 64.8). In contrast, adaptive methods preserve accuracy more effectively: Adaptive-Answer reduces average length by 23.5% (8.3k \rightarrow 6.3k) with only a 0.9-point accuracy decrease, while Format-Adaptive-Answer achieves a 28.6% reduction (5.9k tokens) with a 1.2-point drop. Among all efficient reasoning strategies, our proposed methods consistently occupy the Pareto-optimal region, yielding the best overall accuracy-efficiency trade-offs across benchmarks. These results indicate that instance-level length adaptation yields a substantially better efficiency-accuracy trade-off than fixed or normalized constraints.

AI use disclosure: we used AI for assistance in code writing and in manuscript typesetting.

Model	MATH-500		AIME 24		AIME 25		GPQA Diamond		Common-senseQA		LiveCode-Bench		Long-Benchv2		Average	
	Acc.↑	#Tok.↓	Acc.↑	#Tok.↓	Acc.↑	#Tok.↓	Acc.↑	#Tok.↓	Acc.↑	#Tok.↓	Acc.↑	#Tok.↓	Acc.↑	#Tok.↓	Acc.↑	#Tok.↓
Base Model	93.6	4,837	75.5	14,191	68.5	17,402	55.5	7,284	83.5	1,130	73.1	10,261	39.7	2,980	69.9	8,298
SFT	94.0	4,292	75.9	12,381	66.1	16,058	55.8	6,459	84.0	939	72.5	9,708	39.5	2,915	69.7	7,536
Hard-Length 16k	94.0	4,388	76.2	12,803	67.0	15,707	55.9	6,712	83.5	1,087	72.5	9,909	39.3	2,822	69.8	7,633
Hard-Length 8k	93.8	3,369	72.9	10,633	61.9	12,758	54.9	5,535	83.8	959	73.1	8,772	39.9	2,332	68.6	6,337
Hard-Length 8k → 4k	93.4	2,703	71.2	9,114	55.5	10,770	54.5	4,642	83.7	844	71.3	7,972	38.6	1,994	66.9	5,434
Soft-Length	93.4	3,129	72.0	10,105	58.9	11,944	54.9	5,128	83.6	900	71.6	8,335	38.8	2,137	67.6	5,954
Normalized-Length	92.2	1,734	67.0	11,723	50.2	7,475	53.8	3,011	83.7	537	69.4	6,273	37.1	1,326	64.8	4,583
TWYN	94.2	3,377	74.1	11,491	63.6	14,243	54.4	5,978	83.8	964	72.4	9,259	38.7	2,579	68.8	6,841
Adaptive-Answer	94.0	3,098	75.1	10,261	63.9	13,017	55.5	5,560	84.2	786	71.7	9,089	39.1	2,634	69.0	6,349
Format-Adaptive-Answer	93.8	2,403	73.2	9,583	62.4	11,965	58.5	4,931	84.2	620	71.3	9,416	37.6	2,559	68.7	5,925

Table 4: Accuracy (Acc.↑) and response length (#Tok.↓) of all approaches applied to Qwen3-8B. Best values per column are in bold.