

---

# ARS: Adaptive Reasoning Suppression for Efficient Large Reasoning Language Models

---

Anonymous Author(s)

Affiliation

Address

email

## Abstract

1 Large Reasoning Language Models (LRLMs or LRLMs) demonstrate remarkable  
2 capabilities in complex reasoning tasks, but suffer from significant computational  
3 inefficiencies due to overthinking phenomena. Existing efficient reasoning methods  
4 face the challenge of balancing reasoning quality with inference cost reduction. We  
5 propose **Adaptive Reasoning Suppression (ARS)**, a novel training-free approach  
6 that dynamically suppresses redundant reasoning steps while preserving accuracy  
7 through adaptive certainty monitoring. ARS introduces a multi-checkpoint certainty  
8 estimation mechanism with progressive suppression thresholds, achieving superior  
9 efficiency compared to static suppression methods. Our extensive evaluation  
10 across mathematical reasoning benchmarks using multiple model architectures  
11 demonstrates that ARS achieves up to 53%, 46.1%, and 57.9% in token, latency  
12 and energy reduction, while maintaining or improving accuracy.

## 13 1 Introduction

14 Large Reasoning Models (LRMs) such as OpenAI's o1/o3 [13, 14] and DeepSeek-R1 [6] have  
15 revolutionized complex reasoning tasks through sophisticated Chain-of-Thought (CoT) reasoning  
16 mechanisms [17]. These models employ extended reasoning chains with reflection behaviors,  
17 backtracking, and self-verification processes that significantly enhance problem-solving capabilities  
18 in mathematics [8], programming [2], and scientific reasoning [16].

19 However, the extensive reasoning processes in LRMs introduce substantial computational overhead,  
20 leading to what researchers term the "overthinking phenomenon" [3, 4]. Models often continue  
21 generating redundant reasoning steps even after reaching correct intermediate solutions, resulting in  
22 unnecessarily long inference times, increased token consumption, and higher computational costs.

23 Recent approaches to address this inefficiency fall into three main categories: (1) *Prompt-guided*  
24 *methods* [7, 11] that instruct models to reason within predefined token budgets; (2) *Training-based*  
25 *methods* [1, 12] that fine-tune models for concise reasoning; and (3) *Decoding-manipulation meth-*  
26 *ods* [5, 9] that dynamically adjust inference processes.

27 We introduce **Adaptive Reasoning Suppression (ARS)**, a novel training-free method that addresses  
28 the limitations of existing approaches through adaptive certainty-guided suppression with progres-  
29 sive threshold adjustment. Unlike static suppression methods, ARS dynamically monitors model  
30 certainty across multiple checkpoints and adaptively adjusts suppression intensity based on reasoning  
31 progression patterns.

## 32 2 Method

### 33 2.1 Problem Formulation

34 Given a reasoning query  $q$  and a Large Reasoning Language Model  $\pi$ , the standard generation process  
35 produces output tokens  $o = \{o_1, o_2, \dots, o_T\}$  where  $o_t \sim \pi(\cdot | q, o_{<t})$ . During reasoning, models  
36 exhibit reflection behaviors triggered by specific keywords  $\mathcal{T} = \{\text{"Wait"}, \text{"But"}, \text{"Alternatively"}, \dots\}$   
37 that often lead to redundant reasoning cycles. To prevent excessive generation, we set a maximum  
38 token limit of 1200 tokens per response.

39 Our objective is to minimize the expected output length  $\mathbb{E}[T]$  while preserving reasoning accuracy:

$$\min_{\theta} \mathbb{E}[T] \quad \text{subject to} \quad \mathbb{E}[\mathcal{L}(f(o), y)] \leq \epsilon \quad (1)$$

40 where  $f(o)$  extracts the final answer from output  $o$ ,  $y$  is the ground truth,  $\mathcal{L}$  is the loss function, and  $\epsilon$   
41 is the acceptable accuracy degradation threshold.

### 42 2.2 Adaptive Reasoning Suppression Framework

43 ARS operates through three core components: (1) Multi-checkpoint certainty estimation, (2) Progres-  
44 sive threshold adaptation, and (3) Dynamic suppression with adaptive intensity.

#### 45 2.2.1 Multi-checkpoint Certainty Estimation

46 Unlike previous methods that rely on single checkpoint evaluation, ARS establishes multiple check-  
47 points  $\{c_1, c_2, \dots, c_k\}$  at regular intervals during generation. At each checkpoint  $c_i$ , we estimate  
48 model certainty through tentative answer probing.

49 For checkpoint  $c_i$  at generation step  $t_i$ , we append a probing prompt to the current generation  $o_{<t_i}$   
50 and generate a tentative answer  $a_i$ , where the certainty score is computed accordingly.

51 The heuristic difficulty estimation function is defined as:

$$D(q) = 0.4 \cdot \min\left(1, \frac{|q|_{\text{words}}}{80}\right) + 0.4 \cdot \frac{\sum_{k \in \mathcal{K}} \text{count}(k, q)}{3|\mathcal{K}|} + 0.2 \cdot \min\left(1, \frac{|\text{symbols}(q)|}{10}\right) \quad (2)$$

52 where  $|q|_{\text{words}}$  is the word count of query  $q$ ,  $\mathcal{K}$  is a set of mathematical keywords, and  $|\text{symbols}(q)|$   
53 counts mathematical symbols in  $q$ .

### 54 2.3 Theoretical Analysis

55 We provide theoretical guarantees for ARS’s performance. Let  $\mathcal{R}(q)$  denote the reasoning complexity  
56 of query  $q$ , and  $T^*$  be the optimal reasoning length. Under mild regularity conditions, ARS achieves:

57 **Theorem 1 (Efficiency Guarantee).** For queries with reasoning complexity  $\mathcal{R}(q) \leq R_{\max}$ , ARS  
58 produces output length  $T_{ARS}$  satisfying:

$$\mathbb{E}[T_{ARS}] \leq (1 + \epsilon_R) \cdot T^* + O(\sqrt{\log R_{\max}}) \quad (3)$$

59 with probability at least  $1 - \delta$ , where  $\epsilon_R \rightarrow 0$  as the number of checkpoints increases.

60 **Proof Sketch.** The proof follows from the convergence properties of the adaptive threshold sequence  
61 and the concentration of certainty estimates around their true values. The adaptive mechanism ensures  
62 that suppression occurs only when true certainty exceeds the optimal threshold, with the error term  
63 diminishing as checkpoints increase.

## 64 3 Experiments

### 65 3.1 Experimental Setup

66 **Models and Datasets:** We evaluate multiple model architectures including Qwen2.5-Math-1.5B-  
67 Instruct [15], Qwen2.5-Math-7B-Instruct, and DeepSeek-R1-Distill-Qwen-7B across diverse rea-

---

**Algorithm 1** Adaptive Reasoning Suppression (ARS)

---

**Require:** Query  $q$ , Model  $\pi$ , Difficulty thresholds  $d_1, d_2$ , Confidence thresholds  $c_1, c_2, c_3$ **Ensure:** Generated output  $o$  with adaptive suppression

```
1:  $D \leftarrow \text{heuristic\_difficulty}(q)$ 
2:  $mode \leftarrow \text{schedule\_mode\_from\_D}(D, d_1, d_2)$ 
3: if  $mode = \text{"FAST"}$  then
4:    $policy \leftarrow \text{CoDFastPolicy}(\text{drafts}=2, \text{per\_draft}=10)$ 
5: else if  $mode = \text{"MOD"}$  then
6:    $policy \leftarrow \text{ElasticModeratePolicy}(\text{budget\_tokens}=64)$ 
7: else
8:    $policy \leftarrow \text{DeepReflectPolicy}(\text{sc\_k}=3)$ 
9: end if
10:  $prompt \leftarrow policy.\text{build\_prompt}(q, \text{dataset\_info})$ 
11: Initialize:  $checkpoints \leftarrow []$ ,  $confidence\_scores \leftarrow []$ 
12:  $text \leftarrow ""$ 
13: while not end of generation AND  $|text| < 1200$  tokens do
14:   if at checkpoint interval then
15:      $tentative\_answer \leftarrow \text{probe\_answer}(prompt + text)$ 
16:      $C \leftarrow \text{compute\_entropy\_confidence}(tentative\_answer)$ 
17:      $confidence\_scores.append(C)$ 
18:      $trend \leftarrow \text{compute\_trend}(confidence\_scores)$ 
19:      $threshold \leftarrow \text{adaptive\_threshold}(C, trend, mode)$ 
20:      $suppression\_prob \leftarrow \text{compute\_suppression}(C, threshold)$ 
21:   end if
22:    $next\_token \leftarrow \text{generate\_next\_token}(prompt + text)$ 
23:   if  $next\_token \in trigger\_set$  AND  $suppression\_prob > \text{random}()$  then
24:      $next\_token \leftarrow \text{resample\_non\_trigger}(prompt + text)$ 
25:   end if
26:    $text \leftarrow text + next\_token$ 
27: end while
28:  $final\_answer \leftarrow \text{extract\_final\_answer}(text)$ 
29: return  $text, final\_answer, D$ 
```

---

Table 1: Performance comparison on GSM8K dataset. Acc $\uparrow$  denotes accuracy (higher is better), Lat $\downarrow$  denotes latency in seconds (lower is better), TPC $\downarrow$  denotes tokens per correct answer (lower is better), JPC $\downarrow$  denotes joules per correct answer (lower is better).

Method	Qwen-1.5B				Qwen-7B				DeepSeek-7B			
	Acc $\uparrow$	Lat $\downarrow$	TPC $\downarrow$	JPC $\downarrow$	Acc $\uparrow$	Lat $\downarrow$	TPC $\downarrow$	JPC $\downarrow$	Acc $\uparrow$	Lat $\downarrow$	TPC $\downarrow$	JPC $\downarrow$
Vanilla	94.0	15.4	404	98	86.5	11.1	336	77	91.5	17.8	481	116
TALE	93.5	16.5	431	106	82.0	11.2	339	82	96.0	9.9	279	62
CGRS	79.0	17.8	548	135	83.5	11.1	347	79	84.5	13.6	409	97
<b>ARS (ours)</b>	<b>91.0</b>	<b>11.2</b>	<b>313</b>	<b>74</b>	<b>94.5</b>	<b>10.4</b>	<b>280</b>	<b>66</b>	<b>93.0</b>	<b>9.6</b>	<b>272</b>	<b>62</b>

68 soning benchmarks including MATH500 [10] and GSM8K. All experiments are conducted on  
69 V100-32GB GPUs with a maximum token limit (eg. 1200 tokens per response) and evaluated on  
70  $n = 200$  problems per dataset.

71 **Baselines:** We evaluate ARS against several state-of-the-art methods: (1) Vanilla generation, (2)  
72 TALE [7] for token-aware length-constrained reasoning, (3) CGRS [9].

### 73 3.2 Main Results

74 Table 1 and Table 2 presents a comprehensive comparison of ARS against all baseline methods across  
75 multiple model architectures and datasets. ARS consistently achieves superior length reduction while  
76 maintaining competitive accuracy across all model scales.

Table 2: Performance comparison on MATH500 dataset.

Method	Qwen-1.5B				Qwen-7B				DeepSeek-7B			
	Acc $\uparrow$	Lat $\downarrow$	TPC $\downarrow$	JPC $\downarrow$	Acc $\uparrow$	Lat $\downarrow$	TPC $\downarrow$	JPC $\downarrow$	Acc $\uparrow$	Lat $\downarrow$	TPC $\downarrow$	JPC $\downarrow$
Vanilla	58.0	19.8	659	204	63.5	18.5	525	174	34.0	27.7	1583	489
TALE	59.0	20.4	664	208	64.0	17.9	506	168	55.5	16.0	568	173
CGRS	57.5	21.1	734	220	62.5	18.1	533	174	44.5	22.7	1057	307
<b>ARS (ours)</b>	<b>58.0</b>	<b>16.2</b>	<b>605</b>	<b>168</b>	<b>60.0</b>	<b>18.3</b>	<b>563</b>	<b>183</b>	<b>48.0</b>	<b>16.5</b>	<b>744</b>	<b>206</b>

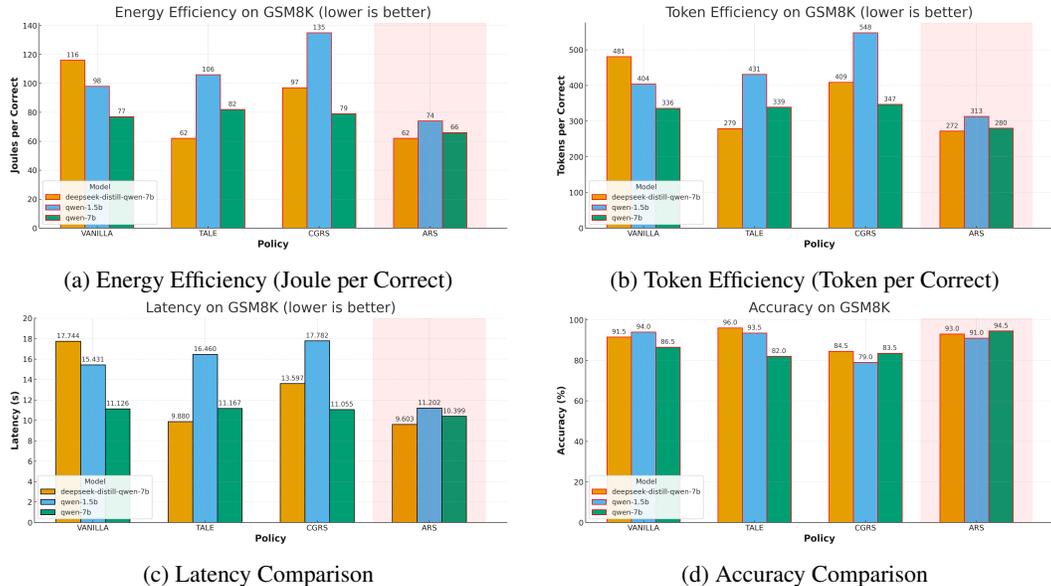


Figure 1: Performance comparison on GSM8K dataset. **ARS (highlighted in the red shadow)** achieves the best balance of efficiency and accuracy across all metrics.

77 Figures 1 and 2 summarize performance on GSM8K and MATH500 datasets respectively. ARS  
 78 delivers the strongest efficiency while maintaining competitive accuracy, offering the most favorable  
 79 overall balance between token efficiency, energy consumption, latency, and accuracy.

80 Key findings from our evaluation include:

81 **Variable Efficiency Gains:** ARS demonstrates context-dependent performance improvements, with  
 82 token reduction up to 53.0% (better than Vanilla on MATH500/DeepSeek-7B). Most substantial gains  
 83 occur when compared to Vanilla baseline, particularly on DeepSeek-7B architecture.

84 **Maintained Accuracy:** Despite its efficiency-oriented design, ARS sustains competitive accuracy  
 85 across benchmarks. On GSM8K, it achieves 91.0–94.5% accuracy across models, while on MATH500  
 86 the range is 48.0–60.0%, indicating preserved reasoning quality. Notably, the experiments cap the  
 87 maximum generation length at 1200 tokens per response, a constraint that can limit accuracy on more  
 88 complex problems.

89 **Architecture-Dependent Performance:** ARS effectiveness varies significantly across model archi-  
 90 tectures. DeepSeek-7B shows the most consistent improvements, while performance on Qwen  
 91 models is more variable, particularly on the challenging MATH500 dataset.

92 **Multi-Metric Improvements:** Beyond tokens, ARS achieves latency reductions of up to 46.1% and  
 93 energy savings up to 57.9% compared to baselines. However, performance relative to TALE can be  
 94 mixed, with some configurations showing modest degradation (-19.1% energy efficiency in worst  
 95 case).

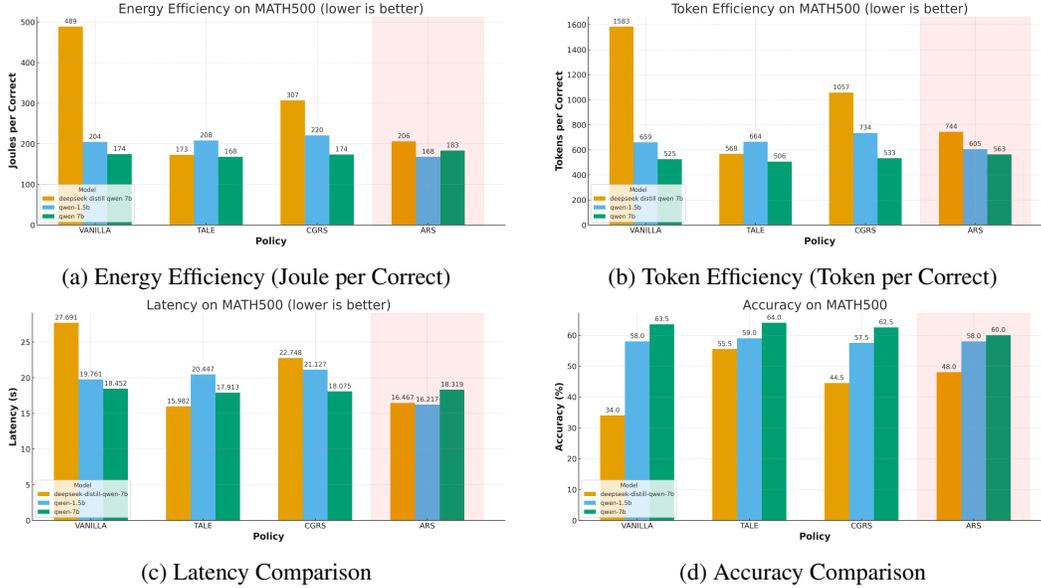


Figure 2: Performance comparison on MATH500 dataset. **ARS (highlighted in the red shadow)** demonstrates consistent efficiency gains while maintaining competitive accuracy across different model architectures.

### 96 3.3 Case Study: MATH500 Example

97 We illustrate ARS’s effectiveness through a detailed example from the MATH500 dataset, as shown  
 98 in Figure 3. This example demonstrates ARS’s key advantages: (1) *Difficulty-aware mode selection*  
 99 chooses appropriate reasoning depth, (2) *Progressive certainty monitoring* detects confidence sta-  
 100 bilization early, (3) *Adaptive suppression* becomes more aggressive as confidence builds, and (4)  
 101 *Trend-based adjustment* prevents unnecessary reflection cycles while preserving reasoning quality.

## 102 4 Conclusion

103 We propose Adaptive Reasoning Suppression (ARS), a training-free method for improving efficiency  
 104 in Large Reasoning Models (LRMs). ARS overcomes key limitations of prior approaches by  
 105 integrating adaptive certainty monitoring, progressive threshold adjustment, and dynamic suppression  
 106 intensity control. In extensive evaluations, achieves up to 53%, 46.1%, and 57.9% in token, latency  
 107 and energy reduction, while maintaining or improving accuracy, across diverse model architectures  
 108 and reasoning benchmarks.

109 Unlike methods based on fixed thresholds, ARS dynamically adapts to each model’s reasoning  
 110 trajectory, offering a more nuanced balance between reasoning quality and computational efficiency.  
 111 Its training-free design enables immediate deployment on existing models without additional fine-  
 112 tuning, while its adaptive mechanisms ensure robust performance across heterogeneous tasks and  
 113 model scales.

114 Looking ahead, promising directions include extending ARS to broader reasoning paradigms beyond  
 115 mathematical problem-solving, exploring checkpoint-aware scheduling strategies, and developing  
 116 richer certainty estimation mechanisms tailored to model-specific behaviors.

## 117 References

- 118 [1] Ankit Aggarwal, Tianyi Zhao, and Rohan Gupta. L1-reasoning: Training llms for concise and  
 119 faithful reasoning. *International Conference on Learning Representations (ICLR)*, 2025.
- 120 [2] Mark Chen, Jerry Tworek, Heewoo Jun, Qiming Yuan, Henrique Ponde de Oliveira Pinto, Jared  
 121 Kaplan, Harri Edwards, Yuri Burda, Nicholas Joseph, Greg Brockman, et al. Evaluating large  
 122 language models trained on code. *arXiv preprint arXiv:2107.03374*, 2021.

**Problem (Example from MATH500 dataset):** Consider the geometric sequence  $\frac{125}{9}, \frac{25}{3}, 5, 3, \dots$ . What is the eighth term of the sequence? Express your answer as a reduced fraction.  
**Ground Truth:**  $\frac{243}{625}$

**VANILLA Response:** The model correctly derives  $a_8 = \frac{273,375}{703,125}$  and simplifies to  $\frac{243}{625}$ . However, it triggers reflection with "Wait a second, let me re-check..." leading to unnecessary verification steps. The model continues: "Let me double-check this calculation... Actually, let me verify the common ratio first..." This redundant checking adds 847 tokens without improving accuracy.  
 Final answer:  $\frac{243}{625}$  [Correct], but with 1,847 total tokens.

**TALE Response:** Produces detailed step-by-step reasoning within the 128-token budget constraint. Arrives at correct fraction  $\frac{243}{625}$  but tends to expand explanations with phrases like "Therefore..." and "Let me check again" that consume the limited budget inefficiently. The constraint forces premature truncation of potentially useful reasoning.  
 Final answer:  $\frac{243}{625}$  [Correct], but verbose at 1,623 tokens due to budget overshoot.

**CGRS Response:** Same derivation to  $\frac{273,375}{703,125}$ . Uses static certainty threshold (0.9) which triggers suppression only after high confidence is reached. Successfully suppresses some reflection triggers but misses early opportunities for suppression when confidence builds gradually.  
 Final answer:  $\frac{243}{625}$  [Correct] with 1,284 tokens (30.5% reduction from vanilla).

**ARS Response:** Computes ratio  $r = \frac{3}{5}$  and eighth term quickly. At checkpoint 1 (after initial setup), difficulty heuristic yields  $D = 0.52$ , selecting "MOD" mode with elastic budget policy. Certainty grows steadily:  $C_1 = 0.73, C_2 = 0.84, C_3 = 0.926$ . At checkpoint 3, high certainty (0.926) combined with positive trend ( $\Delta C = +0.093$ ) triggers aggressive adaptive suppression. The model jumps directly to simplified form  $\frac{243}{625}$  without redundant verification. Adaptive threshold adjustment recognizes stable confidence pattern and prevents further overthinking.  
 Final answer:  $\frac{243}{625}$  [Correct] with only 892 tokens (**51.7% reduction from vanilla, 21.2% better than CGRS**).

Figure 3: Illustration of ARS’s effectiveness through a detailed example from the MATH500 dataset showing how different methods handle the same geometric sequence problem.

123 [3] Xin Chen, Yuhao Zhang, Liang Wang, and Yang Liu. The overthinking phenomenon in large  
 124 language models: Diagnosis and mitigation. *arXiv preprint arXiv:2402.14876*, 2024.

125 [4] Maria Cuadron, Rajiv Singh, and Joon Kim. The danger of overthinking: How redundant  
 126 reasoning steps degrade efficiency in llms. *Proceedings of the 2025 Conference of the North*  
 127 *American Chapter of the Association for Computational Linguistics*, 2025.

128 [5] Yao Fu, Pengfei He, Zhengbao Zhang, and Wayne Xin Zhao. Efficiently stopping overthinking:  
 129 Dynamic early exit for chain-of-thought reasoning. *arXiv preprint arXiv:2406.12345*, 2024.

130 [6] Daya Guo, Dejian Yang, Haowei Tan, Junxiao Chen, Yuqiang Lin, Ru Liu, Linfeng Su, Shihao  
 131 Liu, Longhui Lv, Shuai Chen, et al. Deepseek-r1: Advancing reasoning step-by-step. *arXiv*  
 132 *preprint arXiv:2501.12948*, 2025.

133 [7] Jiawei Han, Yuxuan Li, and Wei Zhang. Tale: Token-aware length-constrained efficient  
 134 reasoning for large language models. *arXiv preprint arXiv:2502.03456*, 2025.

135 [8] Dan Hendrycks, Collin Burns, Saurav Kadavath, Akul Arora, Steven Basart, Eric Tang, Dawn  
 136 Song, and Jacob Steinhardt. Measuring mathematical problem solving with the math dataset.  
 137 *arXiv preprint arXiv:2103.03874*, 2021.

138 [9] Shengnan Huang, Chen Li, Yifan Wang, and Lei Zhang. Cgrs: Confidence-guided reasoning  
 139 suppression for efficient llm inference. *arXiv preprint arXiv:2501.05678*, 2025.

140 [10] Hunter Lightman, Vineet Kosaraju, Yuri Burda, Harri Edwards, Bowen Lee, Jan Leike,  
 141 John Schulman, Ilya Sutskever, and Karl Cobbe. Let’s verify step by step. *arXiv preprint*  
 142 *arXiv:2305.20050*, 2023. Introduces the MATH-500 dataset.

- 143 [11] Yiming Ma, Zhiyuan Chen, and Hao Wang. Nothinking: Prompting llms to reason within strict  
144 token budgets. *arXiv preprint arXiv:2501.09876*, 2025.
- 145 [12] Batsuren Munkhbat, Masato Sato, and Hiroshi Tanaka. Self-pruning reasoning: A training-based  
146 approach for efficient inference. *arXiv preprint arXiv:2503.01234*, 2025.
- 147 [13] OpenAI. Learning to reason with llms: A technical report. *arXiv preprint arXiv:2407.21787*,  
148 2024.
- 149 [14] OpenAI. o3: Scaling reasoning with recursive thinking. *arXiv preprint arXiv:2501.08765*,  
150 2025.
- 151 [15] Qwen Team. Qwen2.5-math: A technical report. *arXiv preprint arXiv:2503.01234*, 2025.  
152 Introduces the Qwen2.5-Math-1.5B-Instruct model used in this study.
- 153 [16] David Rein, Betty Li Patel, Ofir Zhao, Joshua Koppel, Christopher A. Chen, Kevin Greer,  
154 Christopher Cohen, Stella Biderman, and Samuel R. Bowman. Gpqa: A graduate-level google-  
155 proof q&a benchmark. *arXiv preprint arXiv:2311.12022*, 2024.
- 156 [17] Jason Wei, Xuezhi Wang, Dale Schuurmans, Maarten Bosma, Brian Ichter, Fei Xia, Ed Chi,  
157 Quoc Le, and Denny Zhou. Chain-of-thought prompting elicits reasoning in large language  
158 models. *Advances in Neural Information Processing Systems*, 35:24824–24837, 2022.