



MUR: Momentum Uncertainty guided Reasoning for Large Language Models

Hang Yan^{1,2*}, Fangzhi Xu^{1,3*}, Rongman Xu^{1,2}, Yifei Li^{1,2}, Jian Zhang^{1,3}, Haoran Luo⁴, Xiaobao Wu⁵, Anh Tuan Luu^{4,6}, Haiteng Zhao^{7†}, Qika Lin^{8†}, Jun Liu^{1,2†},

¹School of Computer Science and Technology, Xi'an Jiaotong University

²Ministry of Education Key Laboratory of Intelligent Networks and Network Security

³Shaanxi Province Key Laboratory of Big Data Knowledge Engineering

⁴Nanyang Technological University ⁵Shanghai Jiao Tong University ⁶VinUniversity

⁷Shanghai AI Laboratory ⁸National University of Singapore

hyan@stu.xjtu.edu.cn fangzhixu98@gmail.com zhaohaiteng@pku.edu.cn

qikalina@foxmail.com liukeen@xjtu.edu.cn

Abstract

Large Language Models (LLMs) have achieved impressive performance on reasoning-intensive tasks, yet optimizing their reasoning efficiency remains an open challenge. While Test-Time Scaling (TTS) improves reasoning quality, it often leads to overthinking—wasting tokens on redundant computations. This work investigates *how to efficiently and adaptively guide LLM TTS without additional training*. Inspired by the concept of momentum in physics, we propose Momentum Uncertainty guided Reasoning (MUR), which dynamically allocates thinking budgets to critical reasoning steps by tracking and aggregating step-wise uncertainty over time. To support flexible inference-time control, we introduce γ -control, a simple mechanism that tunes the reasoning budget via a single hyperparameter. We provide theoretical intuition to support the superiority of MUR as a low-pass filter. MUR is comprehensively evaluated against various TTS methods across four challenging benchmarks (MATH-500, AIME24, AIME25, and GPQA-diamond) using different sizes of recent Qwen3 models (1.7B, 4B, and 8B). Results demonstrate that MUR reduces computation by over 45% on average while improving accuracy by 0.33–3.46%.

1 Introduction

Large Language Models (LLMs) (Brown et al., 2020; Grattafiori et al., 2024) demonstrate remarkable performance in reasoning-intensive scenarios, including logic, mathematics, and game-playing tasks. A critical advancement in optimizing their

* means equal contribution.

† denotes corresponding authors.

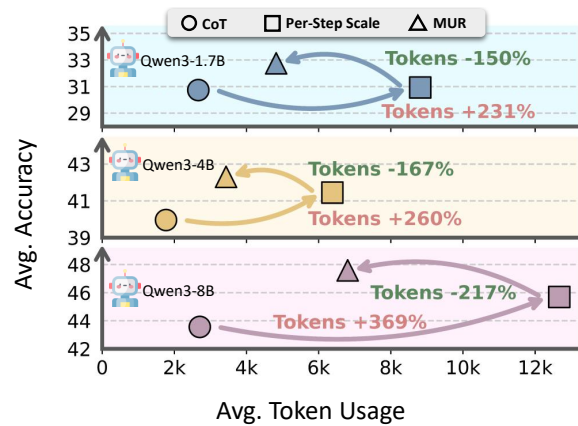


Figure 1: Comparisons of average accuracy and token usage. *Per-Step Scale* refers to TTS methods that optimize every step without compute-saving mechanisms. *MUR* is a computationally efficient approach that selectively scales only key steps. The percentage in this figure is calculated based on CoT budget without TTS.

reasoning quality is *Test-Time Scaling* (TTS). Existing methods either incentivize long thinking patterns through reinforcement learning with verifiable rewards (RLVR) (Ye et al., 2025; Jaech et al., 2024; Guo et al., 2025), or employ stepwise optimization via parallel sampling (Yao et al., 2023; Lightman et al., 2023; Wang et al., 2024b; Ma et al., 2024; Xu et al., 2025) and sequential critique (Lan et al., 2024; Li et al., 2025).

While effective, the issue of *overthinking* (Chen et al., 2024b; Sui et al., 2025) is widely observed that degrades the inference efficiency. As shown in Figure 1, the performance can even be slightly improved, despite >45% reduction in thinking tokens against *Per-Step Scale*. This demonstrates that there is significant room for improvement in making long thinking concise.

Intuitively, LLMs should spend more token bud-

gets on complex steps to deliberately enhance output quality, while generating simple steps directly to avoid overthinking. However, it still remains challenging to identify key steps and dynamically allocate computes. Recent works (Xia et al., 2025a; Jiang et al., 2025; Yang et al., 2025d; Yu et al., 2025; Yang et al., 2025c) explore training methods to adaptively allocate token usage on different steps, which introduce additional training costs and lack generalization. Off-the-shelf training-free methods (Kim et al., 2025; Xu et al., 2025; Wang et al., 2025) scale thinking tokens in a fixed manner, failing to adapt to problem complexity or on-going reasoning process.

Therefore, the pursuit of efficiently and adaptively guiding LLM test-time scaling without extra-training is both intriguing and understudied. To answer this question, we are the first to model LLM reasoning with the concept of momentum. In physics, momentum accumulates historical information over time and resists sudden changes. Based on this and the successful application of Gradient Descent with Momentum (Qian, 1999), we propose Momentum Uncertainty guided Reasoning (*MUR*), a novel approach that dynamically evaluates the overall uncertainty of a reasoning path by aggregating historical step-level uncertainties, mirroring the smooth and consistent evolution observed in physical dynamics. Without requiring any training, *MUR* selectively allocates computation only to critical steps during inference. Based on the approach, we introduce the concept of γ -control, where we can flexibly control the thinking budget and the performance, with only one hyperparameter γ . Further, this work proves that *MUR* is theoretically grounded in terms of discounted credit assignment and stability while maintaining compatibility with existing TTS methods. Extensive experiments across four challenging benchmarks and three backbone model sizes demonstrate that *MUR* reduces the thinking budget by over 45% on average while even improving accuracy by 0.33–3.46%.

The key contributions include:

- (1) Adaptive Scaling Technique.** We propose the novel concept of momentum uncertainty and offer a training-free solution *MUR* to dynamically allocate thinking budgets to key reasoning steps guided by momentum uncertainty, which is compatible with various TTS methods.
- (2) Efficiency and Performance Gains:** *MUR* reduces the thinking costs by 45% even with obvious performance gains, across a wide range of bench-

marks and model sizes. The proposed γ -control offers flexible solution to balance performance and efficiency.

- (3) Theoretical Support:** *MUR* is theoretically grounded in terms of discounted credit assignment, stability, and convergence, which support its practical superiority.

2 Related Work

2.1 Test-Time Scaling

Test-time scaling (TTS) leverages additional inference compute to enhance performance (Brown et al., 2024; Wu et al., 2024b). Existing approaches range from training-based RLVR (Ye et al., 2025; Guo et al., 2025; Wu et al., 2025, 2026b) to training-free strategies, including parallel sampling (Yao et al., 2023; Ma et al., 2024; Xu et al., 2025; Wu et al., 2026a) and sequential refinement (Wu et al., 2024a; Lan et al., 2024; Li et al., 2025). However, they often inefficiently allocate compute to simple steps. *MUR* proposes an orthogonal optimization that guides scaling specifically towards key steps, significantly reducing redundant computation.

2.2 Overthinking

Although LLMs demonstrate significant performance gains through TTS methods, they are likely to introduce computational overhead and reasoning latency (Chen et al., 2024b; Sui et al., 2025). One line of mitigating overthinking is to shorten reasoning length through post-training (Xia et al., 2025a; Jiang et al., 2025; Yang et al., 2025d; Yu et al., 2025; Yang et al., 2025c), which introduces training overhead and limits their generalization. Another line is training-free methods (Kim et al., 2025; Xu et al., 2025; Wang et al., 2025), reducing token usage in a fixed manner, which lacks adaptation to on-going reasoning process. Our work *MUR*, without training, adaptively saves unnecessary computes during the whole reasoning process.

2.3 Uncertainty Estimation

The reasoning path of LLM often contains reliability issues, like hallucinations or biased responses (Xia et al., 2025b). One line of uncertainty estimation is scaling more computes, including verbalizing methods (Tian et al., 2023; Taneru et al., 2024), consistency-based methods (Hou et al., 2024; Chen and Mueller, 2024; Gao et al., 2024), and semantic clustering methods (Kuhn et al., 2023; Farquhar et al., 2024; Nikitin et al.,

2024). Another line is utilizing the internal information during decoding (Ahdritz et al., 2024; Chen et al., 2024a; Sriramanan et al., 2024), which estimates the uncertainty of generated path through aggregating token-level probabilities, lacking the adaptation to different reasoning steps. Our method *MUR*, assigns more attention to recent steps, while reducing the impact of early steps.

3 Method

In this section, we first formulate the stepwise test-time scaling, adaptive scaling and step-level uncertainty (Sec. 3.1). Then we formally propose momentum uncertainty, followed by theoretical proof of its superiority (Sec. 3.2). Based on the momentum uncertainty, we introduce γ -control mechanism to flexibly scale inference-time scaling (Sec. 3.3). The overview of *MUR* is presented in Figure 2.

3.1 Preliminary

Stepwise test-time scaling LLM reasoning can be formulated as auto-regressively generating step a_t at each timestamp t , based on the inputs and previous steps:

$$a_t \sim p_\theta(\cdot|x, \mathbf{a}_{<t}), \quad (1)$$

where x is the concatenation of input question and instruction. $\mathbf{a}_{<t}$ represents previous steps. θ denotes the parameters of pre-trained LLM, and p_θ is the probability distribution. Notably, step division in this paper is autonomously determined by the backbone model through specific prompt structures that guide the model to generate reasoning in logical segments, which can be found in Appendix B.

To optimize the quality of the reasoning path, current methods apply test-time scaling at each step, which can be formulated as follows:

$$\hat{a}_t \sim Q(\cdot|x, \mathbf{a}_{<t}), \quad (2)$$

where \hat{a}_t is the optimized step. Q denotes the specific test-time scaling method, such as *Best-of-N* (Brown et al., 2024).

Adaptive Scaling Conventional test-time scaling methods typically apply optimization at every decoding step, leading to excessive token usage and computational overhead. However, not all steps require such enhancement, and current research on adaptive compute allocation remains limited, often overlooking this inefficiency. We therefore

pose the central question: **When should compute be scaled during inference?** To address this, we model this research question with a binary detector D that selectively activates test-time scaling based on contextual reasoning dynamics:

$$\hat{a}_t = \begin{cases} Q(\cdot|x, \mathbf{a}_{<t}) & , D(t) = \text{True} \\ a_t & , D(t) = \text{False} \end{cases} \quad (3)$$

Here, D determines whether to invoke a test-time scaling method at each step based on historical information. Our work focuses **exclusively** on designing the detector D to assess the reasoning trajectory and adaptively decide whether to allocate additional compute to the current step a_t .

Step-level Uncertainty Uncertainty estimation quantifies an LLM’s confidence in its output, where higher uncertainty implies lower confidence. For step a_t consisting of N tokens, we compute the step-level uncertainty based on token-wise probabilities. Specifically, we define the average negative log-likelihood of the tokens as:

$$m_t = \frac{1}{N} \sum_{j=1}^N -\log p_\theta(a_t^{(j)}|x, \mathbf{a}_{<t}, a_t^{(<j)}), \quad (4)$$

where m_t is the uncertainty of step t . $a_t^{(j)}$ is j -th token of step a_t . And $a_t^{(<j)}$ denotes the prefix token sequence $a_t^{(1)}, a_t^{(2)}, \dots, a_t^{(j-1)}$.

3.2 Momentum Uncertainty

LLM can maintain an uncertainty estimation M for the reasoning process, reflecting the global assessment of both input x and generated steps $\mathbf{a}_{<t}$. Ideally, this uncertainty should evolve smoothly, adapting to new steps as they are generated, while preserving a calibrated estimate of earlier steps. Inspired by the concept of momentum in physics, which retains and updates an object’s motion by accumulating past forces while resisting abrupt changes. We propose momentum uncertainty, a recursive formulation of M that dynamically tracks overall uncertainty during reasoning:

$$M_t = \alpha M_{t-1} + (1 - \alpha)m_t, \quad (5)$$

where M_t is the momentum uncertainty at timestamp t , with initial value $M_0 = 0$. And $\alpha \in (0, 1)$ is a hyper-parameter controlling the momentum changing.

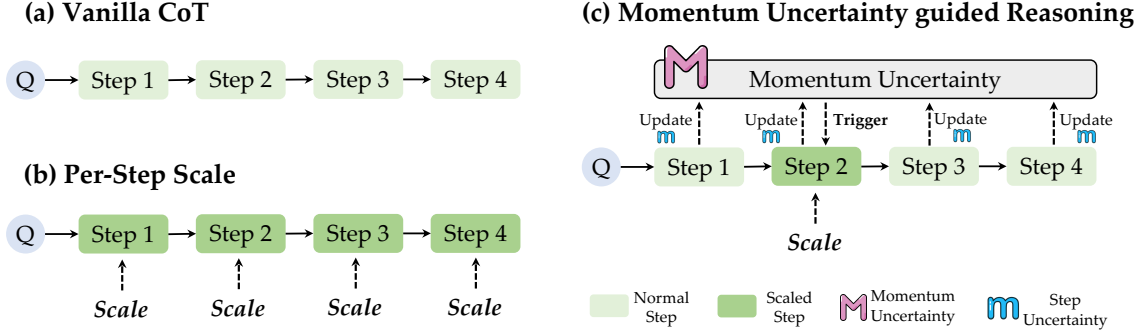


Figure 2: Comparison of reasoning methods. (a) *Vanilla CoT*: Standard stepwise reasoning without test-time scaling. (b) *Per-Step Scale*: scales computes per reasoning step. (c) *MUR*: Adaptive test-time scaling framework (ours).

With a recursive definition, momentum uncertainty aggregates all generated step-level uncertainties to represent the overall estimation of the reasoning process. Further, we introduce the excellent property of *momentum uncertainty* with theoretical and experimental analysis.

Proposition 1: *Momentum uncertainty is an exponentially weighted sum of step-level uncertainties, emphasizing recent steps and fading earlier ones.*

Proof. We provide a detailed derivation in Appendix A.1. It transforms Equation 5 into the exponential weighting of step-level uncertainties as follows:

$$M_t = (1 - \alpha) \sum_{i=1}^t \alpha^{t-i} m_i. \quad (6)$$

Through Equation 6, M_t assigns different weights α^{t-i} to historical step-level uncertainty m_i , emphasizing recent uncertainties while smoothing early fluctuations, balancing the attention among different steps. This aligns with the intuition that recent steps can better represent the reasoning uncertainty, so that momentum uncertainty can well track the evolving of uncertainty change.

Notably, We focus solely on the internal uncertainty signals of the model, disregarding the specific logical information of the output content. This is because the uncertainty signal inherently reflects the accuracy of the model’s reasoning (Xu et al., 2025; Yang et al., 2025b). \square

Proposition 2: *Acting as a low-pass filter, momentum uncertainty M_t attenuates high-frequency components while preserving low-frequency signals, leading to more stable estimates.*

Proof. LLM decoding contains unavoidable noise (Wang et al., 2024a; Zhou et al., 2024), introducing variance to uncertainty estimation. Assume each step-level uncertainty m_t contains two parts:

$$m_t = \mu_t + \epsilon_t, \quad (7)$$

where μ_t is the pure step-level uncertainty, and ϵ_t is a noise originating from training or randomly sampling, etc.

Leveraging the frequency-domain framework of (Li et al., 2024) and the convergence theory of (Liu et al., 2020), we can treat the momentum uncertainty as a low-pass filter as follows:

$$H(\omega) = \frac{1 - \alpha}{1 - \alpha e^{-j\omega}}, \quad (8)$$

where $\omega \in [0, \pi]$ denotes the normalized signal angular frequency. And the derivative of the magnitude response is as follows:

$$\frac{d|H(\omega)|}{d\omega} = -\frac{(1 - \alpha)\alpha \sin \omega}{(1 - 2\alpha \cos \omega + \alpha^2)^{3/2}} < 0, \quad (9)$$

so the magnitude response $|H(\omega)|$ decreases monotonically from 1 to $\frac{1-\alpha}{1+\alpha}$ as ω increases from 0 to π , demonstrating low-pass filter behavior, which can effectively attenuate high-frequency components ϵ_t . The high-frequency signal contains noise and sudden fluctuation of reasoning uncertainty, both of which will be filtered to smooth the estimation process of reasoning uncertainty μ_t . Detailed proof is attached in Appendix A.2. \square

While the auto-regressive nature of LLMs leads to a theoretical expectation of temporal correlation in the noise signal, our empirical findings justify the validity of independent modeling. We analyze

the autocorrelation function (ACF) of the noise signal using real data sampled from several LLMs, and results demonstrate that, we have confidence over 95% to consider real noise signal ϵ_t is not temporally correlated. Based on this critical finding, we provide further theoretical intuition and experimental analysis that momentum uncertainty is superior to naive average uncertainty method (Ren et al., 2022; Manakul et al., 2023; Dobriban et al., 2024). More details can be found in A.3. Moreover, experimental comparison is in Sec. 4.2.

3.3 Scalable Thinking with γ -control

Since momentum uncertainty captures the overall confidence in the reasoning trajectory, we propose a γ -control mechanism to identify whether the current step is incompatible with prior reasoning. This mechanism balances reasoning performance against computational cost.

Scale High-uncertainty Steps At each step, the step-level uncertainty m_t reflects the model’s confidence in the current generation a_t , while M_{t-1} aggregates uncertainty over previous steps. If $m_t > M_{t-1}$, the current step is more uncertain than the reasoning history, suggesting it may be erroneous. To address this, we introduce a checking mechanism that selectively scales uncertain steps.

To tolerate minor fluctuations while flagging significant deviations, we apply a γ -control threshold. Specifically, we define a detector D in Equation 3 as:

$$\hat{a}_t = \begin{cases} Q(\cdot|x, \mathbf{a}_{<t}) & , \exp(m_t) > \exp(M_{t-1})/\gamma \\ a_t & , \text{others} \end{cases}, \quad (10)$$

where γ is the controllable scaling rate, ranging from (0,1) in practice. The scaling factor $\frac{1}{\gamma}$ effectively raises the detection boundary, allowing slight uncertainty increases while catching large deviations. Smaller γ values result in fewer steps being scaled, enabling flexible control over the computational budget. More details can be found in Appendix.

The inequality in Equation 10 flags when a step diverges significantly from the previous reasoning, a corrective test-time scaling is triggered to improve output quality. A theoretical analysis of γ -control is provided in Appendix A.4 and empirical results of γ -control is presented in Sec. 5.1.

Orthogonal to Test-Time Scaling Methods Our momentum uncertainty-based detector D is orthog-

onal and complementary to current test-time scaling methods, such as *best-of-N* and thinking model. It identifies uncertain steps and selectively triggers compute-intensive optimization, maintaining or even improving overall performance while reducing redundancy.

4 Experiments

4.1 Experimental Setup

Benchmarks We evaluate our proposed method *MUR* on three widely adopted math reasoning benchmarks MATH-500 (Hendrycks et al., 2021), AIME24, and AIME25. In addition, we include GPQA-diamond (Rein et al., 2024) to validate the generalization to the science domain.

Metrics We adopt pass@1 rate as our **Acc.** metric. We also report the average token usage of backbone model as **#Token** for each solution, providing an aspect of efficiency evaluation. For AIME24 and AIME25, to reduce the infection of randomness, we sample 16 times for each query and report the average accuracy and token usage.

Test-Time Scaling Settings We adopt four TTS methods as the basic setting.

1) **Guided Search**. It can be viewed as step-level *Best-of-N* (Brown et al., 2024), where N candidate steps are sampled in parallel at each timestep, and the optimal one is selected. 2) **LLM As a Critic**. The LLM receives feedback after generating each step and iteratively refines its output based on the critique (Lan et al., 2024; Li et al., 2025). 3) **ϕ -Decoding** (Xu et al., 2025). It does not require external models but selects the best step from candidates using the foresight sampling strategy. 4) **Thinking Mode** (Yang et al., 2025a) Models with thinking mode generates longer reasoning path, introducing deliberate optimization to each step.

Baselines We adopt four baselines. 1) **CoT** (Wei et al., 2022). Standard stepwise reasoning without scaling. 2) **Per-Step Scale**. Test-time scaling methods that scale the computation for each step. 3) **Avg. uncertainty**. Average the uncertainty across all generated steps (Ren et al., 2022; Manakul et al., 2023; Dobriban et al., 2024) to represent the overall uncertainty of the reasoning process, then scale steps with uncertainty higher than this average. 4) **SMART**. Following the original work by Kim et al. (2025), the backbone model generates

	MATH-500		AIME24		AIME25		GPQA-diamond		Avg.			
	Acc.↑	#Tokens↓	Acc.↑	#Tokens↓	Acc.↑	#Tokens↓	Acc.↑	#Tokens↓	Acc.↑	Δ↑	#Tokens↓	Δ↓
Qwen3-1.7B												
Vanilla CoT	69.20	1,047	17.92	4,243	9.58	4,273	26.26	1,086	30.74	-	2,662	-
Guided search												
+ Per-Step Scale	70.80	3,460	17.92	17,463	10.42	16,680	27.27	6,739	31.60	-	11,086	-
+ Avg uncertainty	70.20	2,398	18.33	7,850	9.58	8,883	25.76	3,404	30.97	(-0.63)	5,634	(-49.18%)
+ SMART	70.80	3,128	17.50	8,955	8.96	10,091	24.74	3,825	30.50	(-1.10)	6,500	(-41.37%)
+ MUR (ours)	71.20	1,321	18.33	4,712	10.63	5,179	32.83	2,005	33.25	(+1.65)	3,304	(-70.19%)
LLM as a critic												
+ Per-Step Scale	70.20	1,098	16.04	3,362	10.00	3,160	28.28	892	31.13	-	2,128	-
+ Avg uncertainty	68.60	1,019	17.92	4,176	9.17	3,174	26.77	1,417	30.62	(-0.51)	2,447	(+14.97%)
+ SMART	70.40	878	18.96	3,976	8.96	3,600	28.28	1,446	31.65	(+0.52)	2,475	(+16.31%)
+ MUR (ours)	71.20	902	19.38	3,892	10.21	4,011	32.32	1,693	33.28	(+2.15)	2,625	(+26.25%)
φ-Decoding												
+ Per-Step Scale	68.00	5,501	17.50	19,612	8.96	18,550	25.76	9,261	30.06	-	13,231	-
+ Avg uncertainty	69.00	2,844	19.17	13,743	8.33	15,785	25.25	2,431	30.44	(+0.38)	8,701	(-34.24%)
+ SMART	70.20	3,848	21.04	19,437	8.13	24,113	23.23	3,338	30.65	(+0.59)	12,684	(-4.13%)
+ MUR (ours)	69.80	2,520	20.21	13,711	9.58	16,088	27.27	1,827	31.72	(+1.66)	8,537	(-35.48%)
Qwen3-4B												
Vanilla CoT	79.40	772	24.08	3,111	16.46	2,577	39.90	612	39.96	-	1,768	-
Guided search												
+ Per-Step Scale	79.80	3,048	29.38	13,761	19.17	10,663	42.42	3,517	42.69	-	7,747	-
+ Avg uncertainty	79.80	1,911	28.33	7,012	18.54	7,719	39.90	1,354	41.64	(-1.05)	4,499	(-41.93%)
+ SMART	81.60	2,476	24.58	8,515	15.42	9,375	43.43	2,116	41.26	(-1.43)	5,621	(-27.45%)
+ MUR (ours)	81.40	824	29.58	4,265	19.17	7,162	41.92	929	43.02	(+0.33)	3,295	(-57.47%)
LLM as a critic												
+ Per-Step Scale	80.80	777	25.21	3,334	17.92	3,260	40.91	737	41.21	-	2,027	-
+ Avg uncertainty	81.40	741	25.63	3,217	20.00	3,120	39.90	804	41.73	(+0.52)	1,971	(-2.79%)
+ SMART	80.60	813	26.04	3,203	17.50	3,201	43.43	724	41.89	(+0.68)	1,985	(-2.06%)
+ MUR (ours)	81.60	745	26.04	3,309	20.21	3,113	40.91	699	42.19	(+0.98)	1,967	(-2.98%)
φ-Decoding												
+ Per-Step Scale	76.80	4,690	27.08	14,394	16.46	14,109	41.41	4,263	40.44	-	9,364	-
+ Avg uncertainty	80.60	1,866	26.67	14,361	18.54	14,836	39.90	1,511	41.43	(+0.99)	8,144	(-13.03%)
+ SMART	79.40	2,776	26.25	19,327	17.71	22,807	40.40	2,195	40.94	(+0.50)	11,776	(+25.76%)
+ MUR (ours)	79.60	1,796	27.29	8,563	18.13	8,845	41.92	944	41.74	(+1.30)	5,037	(-46.21%)
Qwen3-8B												
Vanilla CoT	81.40	1,131	34.17	4,077	18.75	4,746	39.90	859	43.56	-	2,703	-
Guided search												
+ Per-Step Scale	83.20	4,069	35.83	19,805	21.67	21,586	46.46	4,252	46.79	-	12,428	-
+ Avg uncertainty	82.80	2,427	35.21	11,223	22.08	12,193	43.94	2,213	46.01	(-0.78)	7,014	(-43.56%)
+ SMART	82.60	3,502	31.04	17,055	20.00	17,705	46.97	3,797	45.15	(-1.64)	10,515	(-15.39%)
+ MUR (ours)	83.20	2,607	38.13	7,959	24.38	7,582	46.97	3,122	48.17	(+1.38)	5,318	(-57.21%)
LLM as a critic												
+ Per-Step Scale	83.40	1,022	33.13	4,846	21.04	4,818	44.44	1,172	45.50	-	2,965	-
+ Avg uncertainty	82.40	1,086	31.67	5,326	21.88	4,705	41.92	1,375	44.47	(-1.03)	3,123	(+5.35%)
+ SMART	83.20	1,167	32.92	4,737	21.46	4,780	44.95	1,069	45.63	(+0.13)	2,938	(-0.89%)
+ MUR (ours)	83.80	1,132	34.17	4,846	22.50	4,913	44.95	1,007	46.36	(+0.84)	2,975	(+0.34%)
φ-Decoding												
+ Per-Step Scale	84.20	5,841	31.88	43,212	19.58	36,669	43.43	4,726	44.77	-	22,612	-
+ Avg uncertainty	81.80	3,222	34.17	17,807	21.46	20,151	45.45	2,087	45.72	(+0.95)	10,817	(-52.16%)
+ SMART	83.20	4,782	33.13	31,942	22.08	33,123	44.44	4,167	45.71	(+0.94)	18,504	(-18.17%)
+ MUR (ours)	84.40	2,854	36.67	20,969	24.38	22,296	47.47	2,359	48.23	(+3.46)	12,120	(-46.40%)

Table 1: Main results. The best results are highlighted in bold. **Acc.** denotes pass@1 rate and **#Tokens** denotes the backbone model’s average token usage for each query, more details concerning external model token usage is in Appendix C.1. We also report the delta compared to *Per-Step Scale* baseline, including the accuracy difference and the percentage of saved tokens. **Red** indicates worse performance, while **green** indicates better performance against *Per-Step Scale*. Here, \uparrow denotes that higher values are better, whereas \downarrow means lower values are preferable.

reasoning steps autonomously. If the token-level confidence (TLC) falls below a predefined threshold, we apply TTS methods.

Implementation Details We conduct all experiments on different models from Qwen3-series (Yang et al., 2025a), including Qwen3-1.7B, Qwen3-4B, and Qwen3-8B. The hyper-parameter

α and γ are both set to 0.9 as default if no additional explanation is provided. For more implementation details, please refer to Appendix B.

4.2 Main Results

Table 1 and Table 2 report four widely adopted reasoning benchmarks across 3 sizes of models.

	MATH-500		AIME24		AIME25		GPQA-diamond		Avg.			
	Acc.↑	#Tokens↓	Acc.↑	#Tokens↓	Acc.↑	#Tokens↓	Acc.↑	#Tokens↓	Acc.↑	Δ↑	#Tokens↓	Δ↓
Qwen3-1.7B												
Vanilla CoT	69.20	1,047	17.92	4,243	9.58	4,273	26.26	1,086	30.74	-	2,662	-
Thinking Mode												
+ Per-Step Scale	87.60	5,841	41.46	16,392	29.17	17,880	38.89	6,032	49.28	-	11,536	-
+ Avg uncertainty	88.80	4,528	47.29	16,472	30.63	16,948	39.39	5,819	51.53	(+2.25)	10,942	(-5.15%)
+ SMART	89.60	5,214	47.50	17,032	29.17	17,316	38.38	7,678	51.16	(+1.88)	11,810	(+2.37%)
+ <i>MUR</i> (ours)	89.20	5,041	47.71	15,264	31.25	16,146	39.90	5,231	52.02	(+2.74)	10,421	(-9.67%)
Qwen3-4B												
Non-Thinking Mode	79.40	772	25.83	3,111	15.00	2,577	39.90	612	40.03	-	1,768	-
Thinking Mode												
+ Per-Step Scale	89.20	4,598	68.33	13,648	59.38	17,256	51.01	6,547	66.98	-	10,512	-
+ Avg uncertainty	93.80	3,846	68.75	14,832	59.79	18,131	52.53	5,561	68.72	(+1.74)	10,593	(+0.76%)
+ SMART	94.00	4,932	68.33	15,131	58.96	18,104	53.53	8,024	68.71	(+1.72)	11,548	(+9.85%)
+ <i>MUR</i> (ours)	94.00	3,607	68.13	13,009	60.21	16,156	54.04	4,801	69.10	(+2.12)	9,393	(-10.64%)
Qwen3-8B												
Non-Thinking Mode	81.40	1,131	34.17	4,077	18.75	4,746	39.90	859	43.56	-	2,212	-
Thinking Mode												
+ Per-Step Scale	94.60	5,227	72.29	13,793	61.46	17,138	56.06	6,910	71.10	-	10,767	-
+ Avg uncertainty	90.60	4,385	70.42	15,463	60.83	18,608	55.05	6,579	69.23	(-1.87)	11,259	(+4.57%)
+ SMART	93.00	5,482	68.33	16,926	55.42	20,000	54.04	8,726	67.70	(-3.40)	12,784	(+18.73%)
+ <i>MUR</i> (ours)	93.80	5,328	73.33	14,416	61.25	17,779	57.58	6,147	71.49	(+0.39)	10,918	(+1.40%)

Table 2: Results of Thinking Switch. *Vanilla CoT* represents the non-thinking mode. *Per-Step Scale* here denotes the thinking mode of Qwen3 models. Red indicates worse performance against Per-Step Scale, while green indicates better performance. Here, ↑ denotes that higher values are better, whereas ↓ means lower values are preferable.

MUR consistently outperforms strong baselines.

The main results demonstrate the superior token saving capacity of *MUR* in most scenarios, and consistently improves the accuracy against Per-Step Scale methods (from 0.33% to 3.46%). This benefits from reducing overthinking on simple steps, while keeping optimization for difficult steps.

MUR outperforms average uncertainty and SMART on both token usage and accuracy (1.66%, 1.62% for average, respectively). Although the two baselines generate fewer tokens than *MUR* in few cases, the accuracy drops even lower than Per-Step Scale. This indicates that they can’t well evaluate the reasoning process, which laterally proves the superiority of *MUR*.

External critic reduces backbone token usage.

In the *LLM as a critic* setting, *MUR* shows higher backbone token usage than baselines, such as 26.25% more than Per-Step Scale on Qwen3-1.7B. This discrepancy arises because external critics in Per-Step Scale provide hints that shorten the backbone’s output to reach the answer. However, Table 1 only reflects backbone cost. As shown in Appendix C.1, when accounting for both backbone and external model token usage, *MUR* remains the most efficient method.

MUR can generalize to LRMs. Large reasoning models (LRMs) optimize performance by generating overlong reasoning path, leading to excessive token usage. To overcome this, we directly output steps detected as needing no computes scaling by *MUR*, avoiding heavy computes introduced by thinking process. More implementation details can be found in Appendix B. Results in Table 2 demonstrate that *MUR* outperforms all three baselines, improving accuracy from 0.39% to 2.74% against Per-Step Scale baseline, which indicates that *MUR* adaptively identifies key steps during reasoning. This validates the generality of *MUR*.

5 Analysis

In this section, we firstly present scaling law of γ -control (Sec. 5.1), through which we can well control performance and budget balance. Then we analysis the number of reasoning steps and token usage (Sec. 5.2), revealing that *MUR* only scales a minor portion of steps. Finally, we randomly scale some steps (Sec. 5.3), laterally demonstrating that *MUR* can identify crucial steps. Additional analysis of the impact of hyperparameter α and case study can be found in Appendix C.

5.1 Scaling Law of γ -control

γ -control well balance performance and budget. The hyperparameter γ adjusts the detection

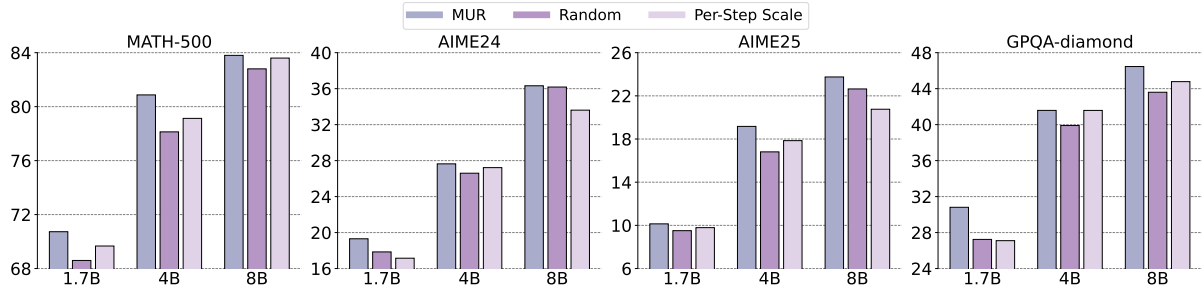


Figure 3: The Y ticks stand for accuracy. X ticks stand for different sizes of Qwen3-series models. For each dataset, we average the three test-time scaling reasoning methods (Guided search, LLM as a critic, ϕ -decoding).

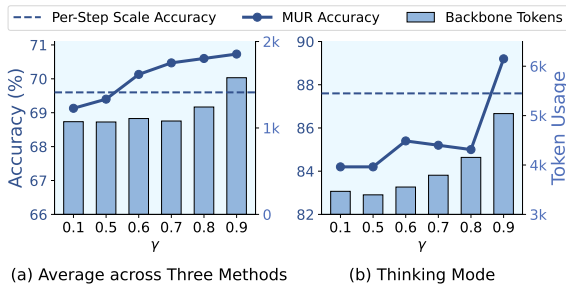


Figure 4: The scaling law of hyperparameter γ . We analyze MATH-500 based on Qwen3-1.7B. The X axis stands for different values of γ . (a) reports the average of Guided search, LLM as a critic and ϕ -Decoding. (b) reports the scaling law of thinking switch.

process in Equation 10, with a lower γ leading to stricter detection boundary condition, then we apply less scaling and less token usage. We report this in Figure 4. The accuracy improves with more token usage, indicating that we can well control the reasoning performance by only adjusting a single hyperparameter γ . It is worth noting that $\gamma = \infty$ equivalents to Per-Step Scale reasoning, whose accuracy drops lower with excessive token usage. More details can be found in Appendix C.2.

5.2 Step and Token Usage Analysis

MUR only scale a minor portion of steps. We report the number of reasoning steps and corresponding token usage under different settings in Figure 5. Under each setting, the result is the average across all the four benchmarks and the three test-time scaling reasoning methods (Guided search, LLM as a critic, ϕ -decoding). With the guidance of MUR, the backbone generates 4.38-6.49 steps for average, scaling only 0.45-0.90 steps for each query. This average steps lower than 1 indicates that for some simple questions, the backbone directly outputs the whole reasoning process, without any scaling, which is equivalent to CoT.

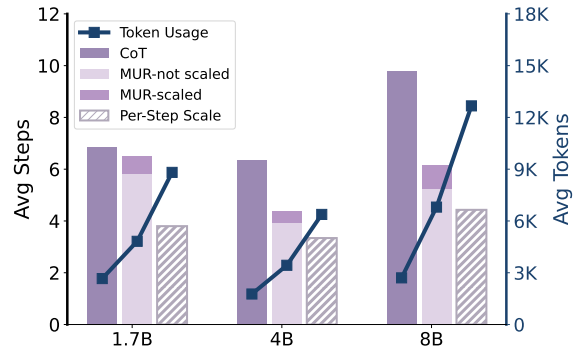


Figure 5: Average steps and token usage for each query. X ticks represent the sizes of different Qwen3-series models. For MUR, we report both scaled steps and not scaled steps.

MUR exhibits superior token efficiency. MUR significantly reduces Per-Step Scale’s token usage over 45% for average. Qwen3-4B generates the least tokens, while Qwen3-8B generates the most tokens, indicating that the former is more efficient and suitable for real-world scenarios.

Scaling reduces total number of steps. Interestingly, we observe that increased scaling leads to fewer steps. For instance, the Per-Step Scale method consumes the most tokens yet generates the least steps on average. This occurs because scaling refines intermediate reasoning, allowing the backbone to reach the solution more directly. Statistics in Appendix C.3 further show that harder benchmarks trigger a higher percentage of scaled steps, reflecting increased model uncertainty.

5.3 Random Scale Result

MUR identifies crucial steps to scale. We randomly scale several steps, keeping the same number of scaled steps as experiments of MUR in Table 1, whose details can be found in Appendix C.3. Results in Figure 3 demonstrates the average accuracy

across the three TTS settings. Random scaling performs worse than Per-Step Scale, indicating that the absence of scaling key steps leads to performance drop. However, *MUR*, which has the same number of scaled steps as random scaling, performs better than both random and Per-Step Scale (1.72% and 1.53% for average), revealing that *MUR* identifies key steps during reasoning.

6 Conclusion

In this paper, we emphasize the key insight that off-the-shelf test-time scaling methods allocate excessive token usage, leading to degradation of both effectiveness and efficiency. To address this, we propose *MUR*, a training-free reasoning framework, which can be orthogonally combined with other test-time scaling methods. We only scale key steps detected by *MUR*. Theoretical analysis and extensive experiments on both LLMs and LRMs demonstrate the superiority of *MUR*.

Limitations

- (1) A primary limitation of *MUR* is its reliance on the model’s internal calibration, which may contain bias from the pre-training phase. Addressing probability calibration is an important problem but is orthogonal to the focus of this work. Importantly, leveraging internal probability signals allows *MUR* to operate without relying on external reward models, auxiliary verifiers, or additional supervision, which significantly improves practicality, efficiency, and deployment simplicity.
- (2) Moreover, *MUR* is difficult to apply directly to unstructured reasoning tasks where the answer is continuous text rather than distinct step-by-step reasoning. However, we emphasize that structured, stepwise reasoning remains a critical frontier for advancing LLM capabilities with a wide domain coverage, particularly in domains such as mathematical problem-solving, where reasoning efficiency continues to be a major bottleneck. By focusing on this high-impact and technically challenging setting, *MUR* makes a meaningful contribution to improving test-time scaling efficiency. Moreover, we believe that the core idea of uncertainty-guided adaptive computation has the potential to be extended to broader reasoning paradigms, which we consider an important direction for future work.

Ethics Statement

We use LLM assistant to polish our writing and code. We agree to release our paper under the MIT License. All external artifacts are utilized in strict compliance with their original licenses and intended academic purposes.

Acknowledgments

This work was supported by Fundamental and Interdisciplinary Disciplines Breakthrough Plan of the Ministry of Education of China (JYB2025XDXM116), National Natural Science Foundation of China (No.62137002, 62293553, 62450005, 62477036, 62192781).

References

- Gustaf Ahdritz, Tian Qin, Nikhil Vyas, Boaz Barak, and Benjamin L Edelman. 2024. Distinguishing the knowable from the unknowable with language models. In *Proceedings of the 41st International Conference on Machine Learning*, pages 503–549.
- Bradley Brown, Jordan Juravsky, Ryan Ehrlich, Ronald Clark, Quoc V Le, Christopher Ré, and Azalia Mirhoseini. 2024. Large language monkeys: Scaling inference compute with repeated sampling. *arXiv preprint arXiv:2407.21787*.
- Tom Brown, Benjamin Mann, Nick Ryder, Melanie Subbiah, Jared D Kaplan, Prafulla Dhariwal, Arvind Neelakantan, Pranav Shyam, Girish Sastry, Amanda Askell, and 1 others. 2020. Language models are few-shot learners. *Advances in neural information processing systems*, 33:1877–1901.
- Chao Chen, Kai Liu, Ze Chen, Yi Gu, Yue Wu, Mingyuan Tao, Zhihang Fu, and Jieping Ye. 2024a. Inside: LLMs’ internal states retain the power of hallucination detection. In *ICLR*.
- Jiuhai Chen and Jonas Mueller. 2024. Quantifying uncertainty in answers from any language model and enhancing their trustworthiness. In *Proceedings of the 62nd Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*, pages 5186–5200.
- Xingyu Chen, Jiahao Xu, Tian Liang, Zhiwei He, Jianhui Pang, Dian Yu, Linfeng Song, Qiuzhi Liu, Mengfei Zhou, Zhuosheng Zhang, and 1 others. 2024b. Do not think that much for $2+3=?$ on the overthinking of o1-like llms. *arXiv preprint arXiv:2412.21187*.
- Edgar Dobriban, Hamed Hassani, Osbert Bastani, Mengxin Yu, Insup Lee, Shuo Li, Xinneng Huang, and Matteo Sesia. 2024. *Uncertainty in language models: Assessment through rank-calibration. Preprint*, arXiv:2404.03163.

- Sebastian Farquhar, Jannik Kossen, Lorenz Kuhn, and Yarin Gal. 2024. Detecting hallucinations in large language models using semantic entropy. *Nature*, 630(8017):625–630.
- Yumeng Fu, Jiayin Zhu, Lingling Zhang, Bo Zhao, Shaoxuan Ma, Yushun Zhang, Yanrui Wu, and Wenjun Wu. 2025. Geolaux: A benchmark for evaluating mllms’ geometry performance on long-step problems requiring auxiliary lines. *arXiv preprint arXiv:2508.06226*.
- Xiang Gao, Jiaxin Zhang, Lalla Mouatadid, and Kamalika Das. 2024. Spuq: Perturbation-based uncertainty quantification for large language models. In *Proceedings of the 18th Conference of the European Chapter of the Association for Computational Linguistics (Volume 1: Long Papers)*, pages 2336–2346.
- Aaron Grattafiori, Abhimanyu Dubey, Abhinav Jauhri, Abhinav Pandey, Abhishek Kadian, Ahmad Al-Dahle, Aiesha Letman, Akhil Mathur, Alan Schelten, Alex Vaughan, and 1 others. 2024. The llama 3 herd of models. *arXiv preprint arXiv:2407.21783*.
- Daya Guo, Dejian Yang, Haowei Zhang, Junxiao Song, Ruoyu Zhang, Runxin Xu, Qihao Zhu, Shitong Ma, Peiyi Wang, Xiao Bi, and 1 others. 2025. Deepseek-r1: Incentivizing reasoning capability in llms via reinforcement learning. *arXiv preprint arXiv:2501.12948*.
- 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. *arXiv preprint arXiv:2103.03874*.
- Bairu Hou, Yujian Liu, Kaizhi Qian, Jacob Andreas, Shiyu Chang, and Yang Zhang. 2024. Decomposing uncertainty for large language models through input clarification ensembling. In *Proceedings of the 41st International Conference on Machine Learning*, pages 19023–19042.
- Aaron Jaech, Adam Kalai, Adam Lerer, Adam Richardson, Ahmed El-Kishky, Aiden Low, Alec Helyar, Aleksander Madry, Alex Beutel, Alex Carney, and 1 others. 2024. Openai o1 system card. *arXiv preprint arXiv:2412.16720*.
- Lingjie Jiang, Xun Wu, Shaohan Huang, Qingxiu Dong, Zewen Chi, Li Dong, Xingxing Zhang, Tengchao Lv, Lei Cui, and Furu Wei. 2025. Think only when you need with large hybrid-reasoning models. *arXiv preprint arXiv:2505.14631*.
- Yujin Kim, Euiin Yi, Minu Kim, Se-Young Yun, and Taehyeon Kim. 2025. Guiding reasoning in small language models with llm assistance. *arXiv preprint arXiv:2504.09923*.
- Diederik P Kingma. 2014. Adam: A method for stochastic optimization. *arXiv preprint arXiv:1412.6980*.
- Lorenz Kuhn, Yarin Gal, and Sebastian Farquhar. 2023. Semantic uncertainty: Linguistic invariances for uncertainty estimation in natural language generation. *arXiv preprint arXiv:2302.09664*.
- Woosuk Kwon, Zhuohan Li, Siyuan Zhuang, Ying Sheng, Lianmin Zheng, Cody Hao Yu, Joseph E. Gonzalez, Hao Zhang, and Ion Stoica. 2023. Efficient memory management for large language model serving with pagedattention. In *Proceedings of the ACM SIGOPS 29th Symposium on Operating Systems Principles*.
- Tian Lan, Wenwei Zhang, Chen Xu, Heyan Huang, Dahua Lin, Kai Chen, and Xian-Ling Mao. 2024. Criticeval: Evaluating large-scale language model as critic. *Advances in Neural Information Processing Systems*, 37:66907–66960.
- Xianliang Li, Jun Luo, Zhiwei Zheng, Hanxiao Wang, Li Luo, Lingkun Wen, Linlong Wu, and Sheng Xu. 2024. On the performance analysis of momentum method: A frequency domain perspective. *arXiv preprint arXiv:2411.19671*.
- Yansi Li, Jiahao Xu, Tian Liang, Xingyu Chen, Zhiwei He, Qiuzhi Liu, Rui Wang, Zhuosheng Zhang, Zhaopeng Tu, Haitao Mi, and 1 others. 2025. Dancing with critiques: Enhancing llm reasoning with stepwise natural language self-critique. *CoRR*.
- Hunter Lightman, Vineet Kosaraju, Yuri Burda, Harrison Edwards, Bowen Baker, Teddy Lee, Jan Leike, John Schulman, Ilya Sutskever, and Karl Cobbe. 2023. Let’s verify step by step. In *The Twelfth International Conference on Learning Representations*.
- Xiaoyu Liu, Paiheng Xu, Junda Wu, Jiaxin Yuan, Yifan Yang, Yuhang Zhou, Fuxiao Liu, Tianrui Guan, Hao-liang Wang, Tong Yu, and 1 others. 2025. Large language models and causal inference in collaboration: A comprehensive survey. *Findings of the Association for Computational Linguistics: NAACL 2025*, pages 7668–7684.
- Yanli Liu, Yuan Gao, and Wotao Yin. 2020. An improved analysis of stochastic gradient descent with momentum. *Advances in Neural Information Processing Systems*, 33:18261–18271.
- Chang Ma, Haiteng Zhao, Junlei Zhang, Junxian He, and Lingpeng Kong. 2024. Non-myopic generation of language models for reasoning and planning. *arXiv preprint arXiv:2410.17195*.
- Xueguang Ma, Qian Liu, Dongfu Jiang, Ge Zhang, Zejun Ma, and Wenhui Chen. 2025. **General-reasoner: Advancing llm reasoning across all domains**. *arXiv:2505.14652*.
- Potsawee Manakul, Adian Liusie, and Mark Gales. 2023. Selfcheckgpt: Zero-resource black-box hallucination detection for generative large language models. In *Proceedings of the 2023 Conference on Empirical Methods in Natural Language Processing*, pages 9004–9017.

- Alexander Nikitin, Jannik Kossen, Yarin Gal, and Pekka Marttinen. 2024. Kernel language entropy: Fine-grained uncertainty quantification for llms from semantic similarities. *Advances in Neural Information Processing Systems*, 37:8901–8929.
- Ning Qian. 1999. On the momentum term in gradient descent learning algorithms. *Neural networks*, 12(1):145–151.
- 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*.
- Jie Ren, Jiaming Luo, Yao Zhao, Kundan Krishna, Mohammad Saleh, Balaji Lakshminarayanan, and Peter J Liu. 2022. Out-of-distribution detection and selective generation for conditional language models. *arXiv preprint arXiv:2209.15558*.
- Gaurang Sriramanan, Siddhant Bharti, Vinu Sankar Sadasivan, Shoumik Saha, Priyatham Kattakinda, and Soheil Feizi. 2024. Llm-check: Investigating detection of hallucinations in large language models. *Advances in Neural Information Processing Systems*, 37:34188–34216.
- Yang Sui, Yu-Neng Chuang, Guanchu Wang, Jiamu Zhang, Tianyi Zhang, Jiayi Yuan, Hongyi Liu, Andrew Wen, Shaochen Zhong, Hanjie Chen, and 1 others. 2025. Stop overthinking: A survey on efficient reasoning for large language models. *arXiv preprint arXiv:2503.16419*.
- Sree Harsha Tanneru, Chirag Agarwal, and Himabindu Lakkaraju. 2024. Quantifying uncertainty in natural language explanations of large language models. In *International Conference on Artificial Intelligence and Statistics*, pages 1072–1080. PMLR.
- Katherine Tian, Eric Mitchell, Allan Zhou, Archit Sharma, Rafael Rafailov, Huaxiu Yao, Chelsea Finn, and Christopher D Manning. 2023. Just ask for calibration: Strategies for eliciting calibrated confidence scores from language models fine-tuned with human feedback. In *Proceedings of the 2023 Conference on Empirical Methods in Natural Language Processing*, pages 5433–5442.
- Bin Wang, Chengwei Wei, Zhengyuan Liu, Geyu Lin, and Nancy Chen. 2024a. Resilience of large language models for noisy instructions. In *Findings of the Association for Computational Linguistics: EMNLP 2024*, pages 11939–11950.
- Peiyi Wang, Lei Li, Zhihong Shao, Runxin Xu, Damai Dai, Yifei Li, Deli Chen, Yu Wu, and Zhifang Sui. 2024b. Math-shepherd: Verify and reinforce llms step-by-step without human annotations. In *Proceedings of the 62nd Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*, pages 9426–9439.
- Yiming Wang, Pei Zhang, Siyuan Huang, Baosong Yang, Zhuosheng Zhang, Fei Huang, and Rui Wang. 2025. Sampling-efficient test-time scaling: Self-estimating the best-of-n sampling in early decoding. *arXiv preprint arXiv:2503.01422*.
- Jason Wei, Xuezhi Wang, Dale Schuurmans, Maarten Bosma, Fei Xia, Ed Chi, Quoc V Le, Denny Zhou, and 1 others. 2022. Chain-of-thought prompting elicits reasoning in large language models. *Advances in neural information processing systems*, 35:24824–24837.
- Jinyang Wu, Mingkuan Feng, Shuai Zhang, Feihu Che, Zengqi Wen, Chonghua Liao, and Jianhua Tao. 2024a. Beyond examples: High-level automated reasoning paradigm in in-context learning via mcts. *arXiv preprint arXiv:2411.18478*.
- Jinyang Wu, Chonghua Liao, Mingkuan Feng, Shuai Zhang, Zhengqi Wen, Haoran Luo, Ling Yang, Huazhe Xu, and Jianhua Tao. 2025. Templaterl: Structured template-guided reinforcement learning for llm reasoning. *arXiv preprint arXiv:2505.15692*.
- Jinyang Wu, Shuo Yang, Changpeng Yang, Yuhao Shen, Shuai Zhang, Zhengqi Wen, and Jianhua Tao. 2026a. Spark: Strategic policy-aware exploration via dynamic branching for long-horizon agentic learning. *arXiv preprint arXiv:2601.20209*.
- Jinyang Wu, Guocheng Zhai, Ruihan Jin, Jiahao Yuan, Yuhao Shen, Shuai Zhang, Zhengqi Wen, and Jianhua Tao. 2026b. Atlas: Orchestrating heterogeneous models and tools for multi-domain complex reasoning. *arXiv preprint arXiv:2601.03872*.
- Yangzhen Wu, Zhiqing Sun, Shanda Li, Sean Welleck, and Yiming Yang. 2024b. Scaling inference computation: Compute-optimal inference for problem-solving with language models. In *The 4th Workshop on Mathematical Reasoning and AI at NeurIPS*, volume 24.
- Heming Xia, Yongqi Li, Chak Tou Leong, Wenjie Wang, and Wenjie Li. 2025a. Tokenskip: Controllable chain-of-thought compression in llms. *arXiv preprint arXiv:2502.12067*.
- Zhiqiu Xia, Jinxuan Xu, Yuqian Zhang, and Hang Liu. 2025b. A survey of uncertainty estimation methods on large language models. *arXiv preprint arXiv:2503.00172*.
- Fangzhi Xu, Hang Yan, Chang Ma, Haiteng Zhao, Jun Liu, Qika Lin, and Zhiyong Wu. 2025. ϕ -decoding: Adaptive foresight sampling for balanced inference-time exploration and exploitation. *arXiv preprint arXiv:2503.13288*.
- Fangzhi Xu, Hang Yan, Qiushi Sun, Jinyang Wu, Zixian Huang, Muye Huang, Jingyang Gong, Zichen Ding, Kanzhi Cheng, Yian Wang, and 1 others. 2026. Odysseyarena: Benchmarking large language models for long-horizon, active and inductive interactions. *arXiv preprint arXiv:2602.05843*.

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. *arXiv preprint arXiv:2505.09388*.

Chenxu Yang, Qingyi Si, Yongjie Duan, Zheliang Zhu, Chenyu Zhu, Qiaowei Li, Minghui Chen, Zheng Lin, and Weiping Wang. 2025b. Dynamic early exit in reasoning models. *arXiv preprint arXiv:2504.15895*.

Junjie Yang, Ke Lin, and Xing Yu. 2025c. Think when you need: Self-adaptive chain-of-thought learning. *arXiv preprint arXiv:2504.03234*.

Wenkai Yang, Shuming Ma, Yankai Lin, and Furu Wei. 2025d. Towards thinking-optimal scaling of test-time compute for llm reasoning. *CoRR*.

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. *Advances in neural information processing systems*, 36:11809–11822.

Yixin Ye, Zhen Huang, Yang Xiao, Ethan Chern, Shijie Xia, and Pengfei Liu. 2025. Limo: Less is more for reasoning. *arXiv preprint arXiv:2502.03387*.

Zhaojian Yu, Yinghao Wu, Yilun Zhao, Arman Cohan, and Xiao-Ping Zhang. 2025. Z1: Efficient test-time scaling with code. *CoRR*.

Jian Zhang, Zhangqi Wang, Haiping Zhu, Kangda Cheng, Kai He, Bo Li, Qika Lin, Jun Liu, and Erik Cambria. 2026a. [Mars: Multi-agent adaptive reasoning with socratic guidance for automated prompt optimization](#). *Proceedings of the AAAI Conference on Artificial Intelligence*, 40(19):16307–16315.

Jian Zhang, Zhiyuan Wang, Zhangqi Wang, Fangzhi Xu, Qika Lin, Lingling Zhang, Rui Mao, Erik Cambria, and Jun Liu. 2026b. [Maps: Multi-agent personality shaping for collaborative reasoning](#). *Proceedings of the AAAI Conference on Artificial Intelligence*, 40(19):16316–16324.

Jian Zhao, Runze Liu, Kaiyan Zhang, Zhimu Zhou, Junqi Gao, Dong Li, Jiafei Lyu, Zhouyi Qian, Biqing Qi, Li Xiu, and 1 others. 2025. Genprm: Scaling test-time compute of process reward models via generative reasoning. *CoRR*.

Zhanke Zhou, Rong Tao, Jianing Zhu, Yiwen Luo, Zeng-mao Wang, and Bo Han. 2024. Can language models perform robust reasoning in chain-of-thought prompting with noisy rationales? *Advances in Neural Information Processing Systems*, 37:123846–123910.

A More Analysis

A.1 The Formulation of Momentum Uncertainty

Proposition 1: *Momentum uncertainty is an exponentially weighted sum of step-level uncertain-*

ties, emphasizing recent steps and fading earlier ones.

Proof. Recursive expansion of M_t :

$$\begin{aligned} M_t &= \alpha M_{t-1} + (1 - \alpha)m_t \\ &= \alpha (\alpha M_{t-2} + (1 - \alpha)m_{t-1}) + (1 - \alpha)m_t \\ &= \alpha^2 M_{t-2} + \alpha(1 - \alpha)m_{t-1} + (1 - \alpha)m_t \\ &\vdots \\ &= \alpha^t M_0 + (1 - \alpha) \sum_{i=1}^t \alpha^{t-i} m_i. \end{aligned} \quad (11)$$

Substituting $M_0 = 0$, we obtain:

$$M_t = (1 - \alpha) \sum_{i=1}^t \alpha^{t-i} m_i. \quad (12)$$

This shows M_t assigns weights α^{t-i} to historical m_i , emphasizing recent uncertainties while smoothing early fluctuations.

Let the average probability of the model’s output at step t , m_t follow $m_t = m_{t-1} - \eta g_t$, where g_t denotes the custom update term at step t . The momentum mechanism implicitly applies decayed weights $1 - \alpha^{t-i}$ to historical updates.

Define cumulative updates $m_t = m_1 - \sum_{i=1}^{t-1} g_i$. Substituting into Equation 5 and Equation 11:

$$\begin{aligned} M_t &= \alpha M_{t-1} + (1 - \alpha)m_t \\ &= \alpha^t m_1 + (1 - \alpha) \sum_{i=1}^t \alpha^{t-i} m_i \end{aligned} \quad (13)$$

$$= m_1 - \sum_{i=1}^{t-1} (1 - \alpha^{t-i}) g_i. \quad (14)$$

Compared to the baseline update $m_t = m_1 - \sum_{i=1}^{t-1} g_i$, the momentum term introduces weights $1 - \alpha^{t-i}$ that decay exponentially with step distance $t - i$. \square

From the above proof, we can easily derive the following two properties:

Property 1: *Momentum Uncertainty is the Exponential Weighting of Historical Uncertainties.*

Property 2: *Momentum Uncertainty has Gradient Descent Equivalence with Decaying Weights.*

A.2 Theoretic Intuition of Stable Estimation

Proposition 2: *Acting as a low-pass filter, the momentum uncertainty M_t attenuates high-frequency components while preserving low-frequency signals, resulting in more stable estimates.*

Proof. The momentum uncertainty M_t is defined by Equation 5 as:

$$M_t = \alpha M_{t-1} + (1 - \alpha)m_t, \quad \alpha \in (0, 1).$$

Leveraging the frequency-domain framework of (Li et al., 2024) and the convergence theory of (Liu et al., 2020), we proceed to analyze the low-pass filtering characteristics of momentum.

Applying the Z-transform to Equation 5 yields:

$$M(z) = \alpha z^{-1}M(z) + (1 - \alpha)m(z), \quad (15)$$

where $M(z)$ and $m(z)$ are Z-transforms of M_t and m_t respectively, and z^{-1} denotes the unit delay operator. Rearranging terms gives the transfer function:

$$H(z) = \frac{M(z)}{m(z)} = \frac{1 - \alpha}{1 - \alpha z^{-1}}. \quad (16)$$

The spectral characteristics are examined through evaluation of the transfer function on the unit circle via the mapping $z = e^{j\omega}$, where $\omega \in [0, \pi]$ denotes normalized angular frequency. This procedure yields the following frequency response.

$$H(\omega) = \frac{1 - \alpha}{1 - \alpha e^{-j\omega}}. \quad (17)$$

It quantifies the system's amplitude and phase variation with frequency.

The magnitude response $|H(\omega)|$ characterizes gain versus frequency:

$$\begin{aligned} |H(\omega)| &= \left| \frac{1 - \alpha}{1 - \alpha e^{-j\omega}} \right| \\ &= \frac{1 - \alpha}{\sqrt{(1 - \alpha \cos \omega)^2 + (\alpha \sin \omega)^2}} \\ &= \frac{1 - \alpha}{\sqrt{1 - 2\alpha \cos \omega + \alpha^2}}. \end{aligned} \quad (18)$$

For $\omega \in (0, \pi)$, the derivative of $|H(\omega)|$ with respect to ω is negative, confirming monotonic decrease:

$$\frac{d|H(\omega)|}{d\omega} = -\frac{(1 - \alpha)\alpha \sin \omega}{(1 - 2\alpha \cos \omega + \alpha^2)^{3/2}} < 0, \quad (19)$$

where $\omega \in (0, \pi)$. Thus, $|H(\omega)|$ decrease monotonically from 1 to $\frac{1 - \alpha}{1 + \alpha}$ as ω increases from 0 to π ,

demonstrating low-pass filter behavior according to (Li et al., 2024), which can effectively attenuate high-frequency components ϵ_t .

For example, at the low frequency where $\omega = 0$:

$$|H(0)| = \frac{1 - \alpha}{\sqrt{1 - 2\alpha + \alpha^2}} = \frac{1 - \alpha}{|1 - \alpha|} = 1.$$

At the high frequency where $\omega = \pi$:

$$|H(\pi)| = \frac{1 - \alpha}{\sqrt{1 + 2\alpha + \alpha^2}} = \frac{1 - \alpha}{1 + \alpha} < 1.$$

As $\alpha \rightarrow 1$, $|H(\pi)| \rightarrow 0$, indicating complete attenuation of high-frequency components when the smoothing factor approaches 1.

In our work, momentum uncertainty tracks the change of μ_t smoothly; we prefer low-frequency signal to a high-frequency signal, which often contains noise and sudden fluctuation of μ_t (Kingma, 2014; Li et al., 2024). Notably, when the sudden fluctuation of μ_t occurs, our scaling boundary condition will be triggered to optimize the current step, which results in a more confident step and higher accuracy (Xu et al., 2025). This process indicates that our momentum uncertainty only needs to maintain the low-frequency part of μ_t , filtering both the noise and the sudden fluctuation. \square

A.3 Momentum performs better than naive average uncertainty

Proposition 3: *Momentum uncertainty can suppress the reasoning noise and well track the evolving of μ_t , resulting in better reasoning performance than average uncertainty.*

Proof. To establish this proposition, it is necessary to impose a theoretical assumption on the distribution of the noise. Owing to the autoregressive nature of large language models, an intuitive expectation is that the noise across different reasoning steps exhibits temporal dependence, which substantially complicates the theoretical analysis. Consequently, we adopt an approximate assumption that the noise terms are independent and identically distributed. In the following, we demonstrate the plausibility of this assumption through empirical analysis.

Proposition 3.1: *Noise terms from different steps are weak correlated.*

We conduct experiments on real data collected from Qwen3 series. For each reasoning trajectory,

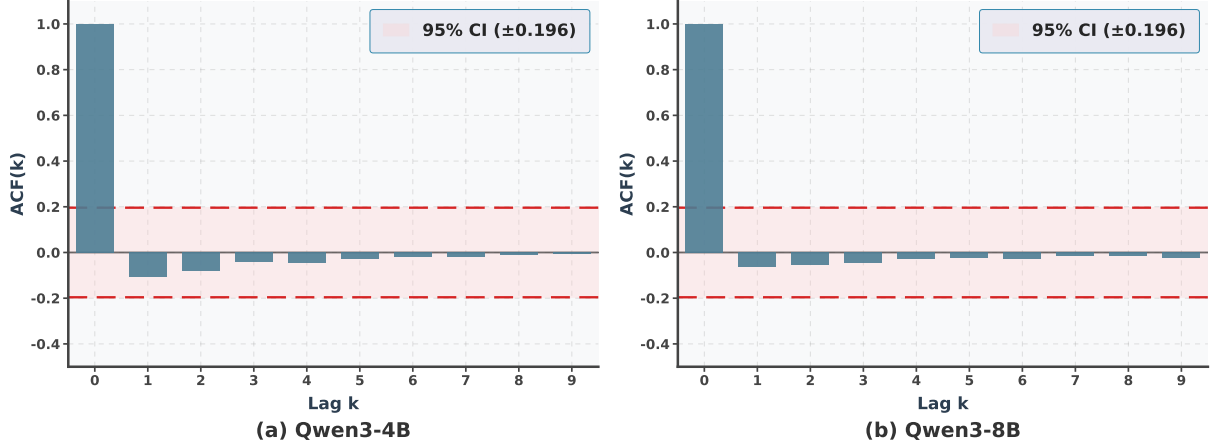


Figure 6: Autocorrelation function (ACF) of signal $\hat{\epsilon}_t$. We conduct this on both Qwen3-4B and Qwen3-8B. The max step length is set to 100 to better represent the temporal process of LLM’s reasoning.

we can get uncertainty m_t from each timestamp t . As in Equation 7, the step uncertainty contains pure uncertainty μ_t and noise ϵ_t :

$$m_t = \mu_t + \epsilon_t.$$

We employ an exponential moving average (EMA) to estimate the step-wise pure signal μ_t , using a deliberately large smoothing coefficient $\alpha_{smooth} = 0.999$. It is worth noting that this constitutes a separate EMA procedure from our momentum uncertainty update used for M_t , and in particular relies on a different choice of α_{smooth} . The estimation process is as follows:

$$\hat{\mu}_t = \alpha_{smooth}\hat{\mu}_{t-1} + (1 - \alpha_{smooth})m_t. \quad (20)$$

Notably, $\hat{\mu}_t$ is only the estimation of μ_t . Due to the low-pass capacity of EMA described in A.2, an extremely large smoothing coefficient α' only maintains the extremely low frequency signal from m_t . The filtered signal contains $\hat{\epsilon}_t$ two parts: 1) high frequency noise ϵ_t 2) part of μ_t that is not confined to the ultra-low-frequency regime.

$$\hat{\epsilon}_t = \epsilon_t + (\mu_t - \hat{\mu}_t). \quad (21)$$

Here, $(\mu_t - \hat{\mu}_t)$ is part of the signal μ_t , which exhibits pronounced autocorrelation due to the auto-regressive nature of LLMs. In other words, $(\mu_t - \hat{\mu}_t)$ is highly temporally correlated with $(\mu_{t-1} - \hat{\mu}_{t-1})$.

This yields a stronger form of hypothesis testing: if $\hat{\epsilon}_t$, a sequence contaminated by the highly correlated signal $(\mu_t - \hat{\mu}_t)$, still exhibits weak autocorrelation in a statistical sense, then it necessarily

implies that the original signal ϵ_t possesses a weak autocorrelation.

In Figure 6, we analyze the autocorrelation function (ACF) of signal $\hat{\epsilon}_t$, showing that for all lags $k \geq 1$, the values of $ACF(k)$ immediately and persistently fall within the 95% confidence interval (CI). The rapid decay and statistical insignificance jointly provide strong evidence that the sequence $\hat{\epsilon}_t$ lacks any substantial long-term or persistent serial correlation. As shown in Equation 21, we can assert with 95% confidence that $\hat{\epsilon}_t$, the sum of ϵ_t and $(\mu_t - \hat{\mu}_t)$, is not temporally correlated. Besides, due to the temporally correlated nature of $(\mu_t - \hat{\mu}_t)$, real noise signal ϵ_t is not temporally correlated with confidence over 95%. This finding is also aligned with recent research (Liu et al., 2025).

Building on **Proposition 3.1**, we posit an **approximate assumption** that each noise signal is white-noise. Notably, this analysis only provides theoretical intuition on why momentum uncertainty is better, rather than rigorous derivation. Moreover, we will provide experimental results to support our proposition.

Theoretical Intuition on the Superior of Momentum Uncertainty than Average Uncertainty.

The momentum uncertainty M_t is defined by Equation 12 as:

$$M_t = (1 - \alpha) \sum_{i=1}^t \alpha^{t-i} m_i, \quad \alpha \in (0, 1).$$

As our approximate assumption, historical uncertainties m_t contain independent noise:

$$m_t = \mu_t + \epsilon_t, \text{Var}(\epsilon_t) = \sigma^2, \quad (22)$$

where σ_t^2 is a bounded constant and μ is the ideal value without variance and bias that can represent the current reasoning and overall reasoning path status. However, it is impractical to get μ , and we can only get step-level uncertainty m which contains noise. Therefore, in our method, we aggregate each step-level uncertainty m as momentum uncertainty M to represent the overall reasoning process.

$$\begin{aligned}\text{Var}(M_t) &= (1 - \alpha)^2 \sum_{i=1}^t \alpha^{2(t-i)} \sigma^2 \\ &= (1 - \alpha)^2 \sigma^2 \sum_{i=1}^t \alpha^{2(t-i)}.\end{aligned}\quad (23)$$

Let $j = t - i$. The summation becomes a finite geometric series:

$$\begin{aligned}\sum_{i=1}^t \alpha^{2(t-i)} &= \sum_{j=0}^{t-1} \alpha^{2j} \\ &= \frac{1 - \alpha^{2t}}{1 - \alpha^2}.\end{aligned}\quad (24)$$

Substituting Equation 24 into Equation 23:

$$\text{Var}(M_t) = (1 - \alpha)^2 \frac{1 - \alpha^{2t}}{1 - \alpha^2} \sigma^2.\quad (25)$$

The vast majority of inference steps are less than twenty (as illustrated in Table 4), so t is set to $t \leq 20$. For $t \leq 20$ and $\alpha \in (0, 1)$, $\alpha^{2t} \approx 0$. Thus:

$$\text{Var}(M_t) \approx \sigma^2 \frac{(1 - \alpha)^2}{1 - \alpha^2} = \sigma^2 \frac{1 - \alpha}{1 + \alpha}.\quad (26)$$

From Equation 26, we can observe that that variance of M_t is lower than the variance of step uncertainty, which is caused by noise ϵ . We establish M_t 's superiority through the following analysis.

Let the simple average be:

$$\tilde{M}_t = \frac{1}{t} \sum_{i=1}^t m_i.\quad (27)$$

For \tilde{M}_t :

$$\text{Var}(\tilde{M}_t) = \frac{1}{t^2} \sum_{i=1}^t \sigma^2 = \frac{\sigma^2}{t}.\quad (28)$$

When $\alpha \rightarrow 1$:

$$\frac{1 - \alpha}{1 + \alpha} < \frac{1}{t} \quad \text{for } t \leq 20,\quad (29)$$

which implies $\text{Var}(M_t) < \text{Var}(\tilde{M}_t)$. Momentum achieves superior noise suppression through exponentially decaying weights. Besides, our main results in Table 1 and Table 2 laterally proves the better detection performance of momentum uncertainty.

Empirical Analysis on the Superior of Momentum Uncertainty than Average Uncertainty.

As described in A.1, momentum uncertainty M_t implements an exponentially decaying weighting scheme that assigns larger weights to recent steps and progressively smaller weights to earlier ones, thereby enabling adaptive tracking of temporal variations in the latent signal μ_t . In contrast, simple averaging assigns equal weights to all steps, which induces substantial tracking lag when the underlying signal μ_t changes, failing to adequately reflect the model's current state.

This contrast yields a stronger form of empirical validation: if momentum uncertainty can more accurately track a slowly evolving signal than average uncertainty, it is expected to exhibit even greater relative advantages in regimes where μ_t displays a mixture of slowly and rapidly varying temporal dynamics.

We provide an experimental analysis from real data to compare between momentum uncertainty and average uncertainty. Our core objective is to demonstrate that M_t provides a more stable and accurate estimation of μ_t when it evolves slowly.

We use extremely large α_{smooth} and slow-evolving estimation $\hat{\mu}$ defined in Equation 20. Besides, we define variance reduction rate as follows:

$$\Delta V = \frac{\text{Var}(\tilde{M}_t) - \text{Var}(M_t)}{\text{Var}(\tilde{M}_t)} \times 100\%,\quad (30)$$

where ΔV stands for variance reduction rate. A higher ΔV means that momentum uncertainty is better than average uncertainty.

We perform this analysis on Qwen3-series. As shown in Figure 7, in both settings, most points are above the red diagonal line, which indicates that momentum uncertainty performs much better than average uncertainty in tracking μ_t .

Under the approximate white-noise assumption, we conduct theoretical analysis on the superiority of momentum uncertainty over average uncertainty. In addition, our experiments serve as empirical evidence that supports this proposition. \square

A.4 Proof of Dynamic Compute Scaling

Proposition 4: *Optimization should be triggered with high confidence when the step-level uncertainty exhibits a significant deviation from the momentum-based uncertainty.*

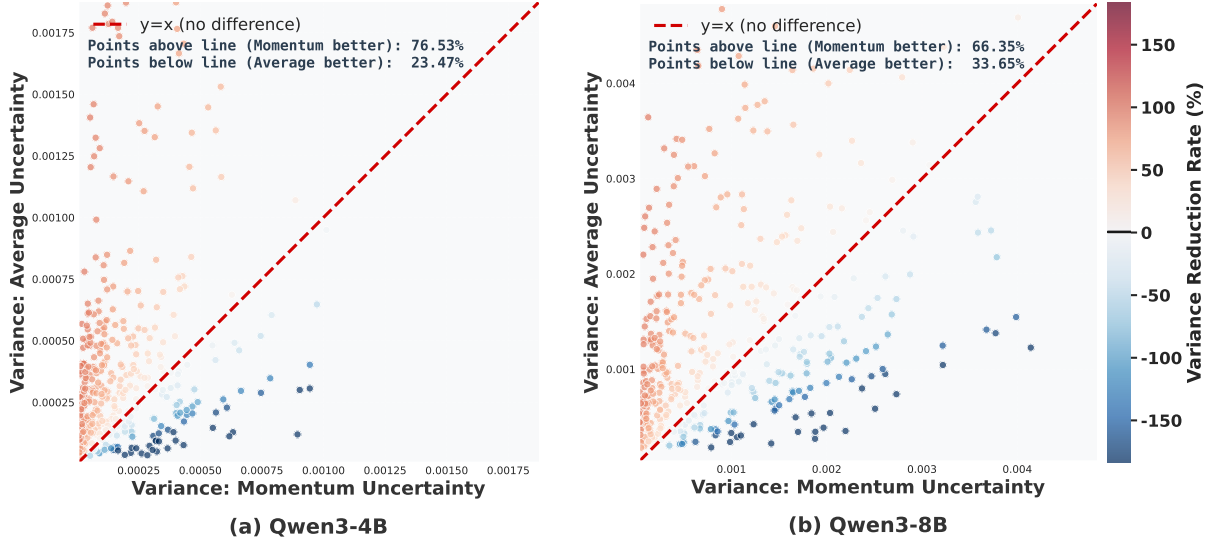


Figure 7: Variance comparison between momentum uncertainty and average. We conduct this on both Qwen3-4B and Qwen3-8B. Each point is one real reasoning path generated from LLM. Points above the red diagonal line represent that momentum uncertainty is better than average uncertainty.

Problem Formulation and Notation Let m_t denote the uncertainty of the model’s output at step $t + 1$, and M_{t-1} represent the momentum uncertainty defined as an exponentially weighted sum, and $\alpha \in (0, 1)$ be the momentum rate. The decision rule for computes scaling is formulated as:

$$\exp(m_t) > \exp(M_{t-1})/\gamma.$$

A boundary violation is flagged when this inequality holds, triggering corrective test-time scaling. We formalize the robustness guarantee below.

Based on the following two lemmas, we establish that the misjudgment probability of historical momentum uncertainty M_{t-1} exceeding the threshold $\tau_t = m_t + \ln \gamma$ approaches zero, demonstrating: When the scaling condition $\exp(m_t) > \exp(M_{t-1})/\gamma$ holds, the model identifies abnormal elevation in current uncertainty m_t with near-certain confidence, thereby efficiently triggering resource scaling.

We now provide a theoretical bound on the probability that a stable reasoning step is mistakenly flagged as uncertain.

Lemma 1: *Chernoff Bound for Single Random Variable.* By using the distribution of random variables, a more precise boundary is provided for the large deviation probability of random variables.

Let X be a real-valued random variable with moment generating function $\phi(s) = \mathbb{E}[e^{sX}]$. For any threshold $\tau \in \mathbb{R}$, the upper tail probability

satisfies:

$$\mathbb{P}(X \geq \tau) \leq \inf_{s>0} e^{-s\tau} \phi(s).$$

X has variance parameter $\hat{\sigma}_t$, $\phi(s) \leq e^{s\nu + \frac{s^2 \hat{\sigma}_t^2}{2}}$, then:

$$\mathbb{P}(X \geq \tau) \leq \exp\left(-\frac{(\tau - \nu)^2}{2\hat{\sigma}^2}\right),$$

where $\nu = \mathbb{E}[X]$.

$$\tau_t = m_t + \ln(\gamma), \quad \gamma \in (0, 1).$$

Lemma 2: *Hoeffding’s inequality.* Hoeffding’s inequality provides the upper limit of the probability that the sum of a random variable deviates from its expected value.

Assume that for each i , $X_i \in [a_i, b_i]$. Consider the sum of these random variables:

$$S_n = \sum_{i=1}^n X_i = X_1 + X_2 + X_3 + \dots + X_{n-1} + X_n.$$

Then Hoeffding’s inequality states that for all $t > 0$:

$$\mathbb{P}(S_n - \mathbb{E}[S_n] \geq t) \leq \exp\left(-\frac{2t^2}{\sum_{i=1}^n (b_i - a_i)^2}\right).$$

$$\mathbb{P}(|S_n - \mathbb{E}[S_n]| \geq t) \leq 2 \exp\left(-\frac{2t^2}{\sum_{i=1}^n (b_i - a_i)^2}\right).$$

Here $\mathbb{E}[S_n]$ denotes the expectation of S_n .

Let the momentum uncertainty sequence M_{t-1} be an exponentially weighted sum of historical step-level uncertainties $\{m_i\}_{i=1}^{t-1}$:

$$M_{t-1} = \sum_{i=1}^{t-1} \omega_i m_i, \omega_i = \alpha^{t-1-i}(1-\alpha), \sum_{i=1}^{t-1} \omega_i = 1,$$

where $m_i \in [0, 1]$ are bounded random variables. The threshold has been defined above, which is:

$$\tau_t = m_t + \ln \gamma.$$

When scaling condition $\exp(m_t) > \exp(M_{t-1})/\gamma$ holds, applying **Lemma 1**, we have:

$$\mathbb{P}(M_{t-1} \geq \tau_t) \leq \exp\left(-\frac{(\tau_t - \hat{\nu}_{t-1})^2}{2\hat{\sigma}_{t-1}^2}\right),$$

where $\hat{\nu}_{t-1} = \mathbb{E}[M_{t-1}]$, and the decay rate is controlled by α .

Proof. By the exponential smoothing definition:

$$M_{t-1} = \sum_{i=1}^{t-1} \omega_i m_i, \quad \omega_i = (1-\alpha)\alpha^{t-1-i},$$

where $m_i \in [0, 1]$ are independent or weakly dependent random variables. Define $X_i = \omega_i m_i$, which satisfies:

- $X_i \in [0, \omega_i]$.
- $b_i - a_i = \omega_i - 0 = \omega_i$.

Applying **Lemma 2**:

$$\begin{aligned} & \mathbb{P}[\exp(M_{t-1} - \mathbb{E}[M_{t-1}]) \geq \zeta] \\ & \leq \exp\left(-\frac{2\zeta^2}{\sum_{i=1}^{t-1} (b_i - a_i)^2}\right) \\ & = \exp\left(-\frac{2\zeta^2}{\sum_{i=1}^{t-1} \omega_i^2}\right). \end{aligned}$$

M_{t-1} is sub-Gaussian with parameter: $\hat{\sigma}_{t-1}^2 = \frac{1}{4} \sum_{i=1}^{t-1} \omega_i^2$. Thus:

$$\mathbb{P}(M_{t-1} - \hat{\nu}_{t-1} \geq \zeta) \leq \exp\left(-\frac{\zeta^2}{2\hat{\sigma}_{t-1}^2}\right). \quad (31)$$

Substitute $\zeta = \tau_t - \hat{\nu}_{t-1}$:

$$\begin{aligned} & \mathbb{P}(M_{t-1} \geq \tau_t) \\ & \leq \exp\left(-\frac{(\tau_t - \hat{\nu}_{t-1})^2}{2 \cdot \frac{1}{4} \sum_{i=1}^{t-1} \omega_i^2}\right) \\ & = \exp\left(-\frac{(\tau_t - \hat{\nu}_{t-1})^2}{2 \cdot \frac{1}{4} ((1-\alpha)^2 \sum_{j=0}^{t-2} (\alpha^2)^j)}\right) \\ & = \exp\left(-\frac{(\tau_t - \hat{\nu}_{t-1})^2}{2 \cdot \frac{1}{4} ((1-\alpha)^2 \cdot \frac{1-\alpha^{2(t-1)}}{1-\alpha^2})}\right) \\ & = \exp\left(-\frac{2(\tau_t - \hat{\nu}_{t-1})^2(1+\alpha)}{(1-\alpha)(1-\alpha^{2(t-1)})}\right) \\ & = \exp\left(-\frac{2(m_t + \ln \gamma - \hat{\nu}_{t-1})^2(1+\alpha)}{(1-\alpha)(1-\alpha^{2(t-1)})}\right). \end{aligned} \quad (32)$$

Since $1 - \alpha^2 = (1 - \alpha)(1 + \alpha)$, $\alpha \in (0, 1)$:

$$\sum_{i=1}^{t-1} \omega_i^2 = (1-\alpha) \cdot \frac{1-\alpha^{2(t-1)}}{1+\alpha} \leq \frac{1-\alpha}{1+\alpha}.$$

Substituting the weight sum upper bound:

$$\begin{aligned} & \mathbb{P}(M_{t-1} \geq \tau_t) \\ & \leq \exp\left(-\frac{2(m_t + \ln \gamma - \hat{\nu}_{t-1})^2(1+\alpha)}{1-\alpha}\right). \end{aligned} \quad (34)$$

(35)

□

As those in practice, we set $\alpha = 0.9$ in the probability bound here:

$$\begin{aligned} & \mathbb{P}(M_{t-1} \geq \tau_t) \\ & \leq \exp\left(-\frac{2(m_t + \ln \gamma - \hat{\nu}_{t-1})^2(1+\alpha)}{1-\alpha}\right) \\ & = \exp(-38(m_t + \ln \gamma - \hat{\nu}_{t-1})^2) \rightarrow 0. \end{aligned}$$

Define the confidence parameter ε as:

$$\varepsilon = \exp\left(-\frac{2(\ln \gamma + m_t - \hat{\nu}_{t-1})^2(1+\alpha)}{1-\alpha}\right).$$

This exponential decay ensures that deviations above $\tau_t = \ln \gamma + m_t$ are asymptotically improbable. With $\alpha = 0.9$, the bound becomes: $\varepsilon = \exp(-38(\ln \gamma + m_t - \hat{\nu}_{t-1})^2) \rightarrow 0$,

$$\mathbb{P}(M_{t-1} \geq \tau_t) = \varepsilon \rightarrow 0,$$

$$\begin{aligned} \mathbb{P}(M_{t-1} < \tau_t) &= \mathbb{P}\left(\exp(m_t) > \frac{\exp(M_{t-1})}{\gamma}\right) \\ &= 1 - \varepsilon. \end{aligned}$$

This validates the scaling decision: **The scaling condition** $\exp(m_t) > \exp(M_{t-1})/\gamma$ **holds with confidence** $1 - \varepsilon$. This result establishes generalization error control for exponential smoothing: The weighted average M_{t-1} converges to the expected uncertainty level, while the scaling condition controls abrupt deviations via tail probability analysis.

B Implementation Details

Implementation of Main Experiments Hyperparameter α and γ are set to 0.9 as default without specific claim. The temperature is set to 0.6 for all experiments. We set top-p to 0.8, top-k to 20. We set presence penalty to 1.5 and max output length to 16,384 tokens. Experiments are conducted on Nvidia A100 GPUs.

For Guided Search setting, we generate four candidates and only one verification path for each candidate. Notably, each verification contains a evaluation path and a final answer token *Yes* or *No*, indicating whether the current step is correct or not. If there is no *Yes* token in all verifications, we select the candidate with lowest probability of *No* token. Otherwise, we select the candidate with the highest probability of *Yes* token.

For LLM As a Critic setting, we prompt the critic to output whether current step is correct and the exact reason. For incorrect steps, we feed the reason path to the backbone model for better output. Specifically, we first prompt the external LLM to generate a reasoning path to judge the correctness of the generated step from the backbone model and then output token *Yes* or *No*. If the judgment token is *Yes*, we do nothing, or we will put the evaluation reasoning path to the backbone model, followed by generating an optimized reasoning step.

For ϕ -Decoding setting, we use TF-IDF metric to cluster, and we do not add the advantage term because we will not scale every step in *MUR*, which leads to the infeasibility of calculating advantage between adjacent steps. We follow the idea of foresight sampling proposed in ϕ -Decoding to use the foresight texts. In the original, the calculation of advantage is implemented by (foresight score of $step_t$ minus foresight score of $step_{t-1}$). However, as explained in *MUR*, we do not need foresight at each step. This foresight score is not available at each step in *MUR*, thus we do not include it. Notably, the remained part is also effective (Xu et al., 2025).

In practice, we do not scale the first step. Be-

cause there is no valid momentum uncertainty when identifying the first step. To achieve smoother estimation in early steps, we introduce a bias correction term following Adam (Kingma, 2014). We set the max step to 20 as default, which is well aligned with the proof in Appendix A.3.

We use General Reasoner (Ma et al., 2025) for math problem evaluation, including MATH, AIME24, AIME25. For GPQA-diamond evaluation process, we provide a python code to parse the final answer and compare it to ground truth. We adopt GenPRM (Zhao et al., 2025) as the external model for candidate selection and critic generation. We conduct all of our experiments based on vLLM (Kwon et al., 2023) reasoning tool.

Implementation of Generating One Step For generating one step, we prompt the backbone LLM to automatically define one step. Specifically, we add *Always end your solution with the phrase "the answer is" followed by your final answer. Start your solution with "Step {stepidx}:"* to the end of each input query. For the update of momentum uncertainty, we use the step-level uncertainty of optimized step. The max of each step’s length is set to 2,048 tokens.

Implementation of Thinking Switch Based on the switch interface between non-thinking mode and thinking mode provided by Qwen3-series, we propose to reduce token usage for large reasoning models with *MUR*. Specifically, we use non-thinking mode as default reasoning method, and switch to thinking mode when current step is detected as needing scaling by *MUR*. We set γ to 0.9, 0.8, 0.7 for MATH, AIME, GPQA-diamond, respectively. To avoid overthinking in each step, we limit the max thinking length to 2,048 tokens and extract all the completed sentences. Additionally, we add “Okay, so I need to” to the beginning of each prompt to correctly elicit thinking in thinking mode.

Prompt used in our experiments 1) User prompt for all settings. 2) System prompt for different datasets. We use empty system prompt for MATH-500 dataset. 3) External model prompt, in which *step_output* represents each step’s answer from the backbone model. 4) Evaluation prompt for MATH-500, AIME24 and AIME25 datasets.

Implementation of Detector The detector plays a vital role in identifying which step to scale, we implement this by maintaining and updating two

User Prompt for All Settings

INPUT QUESTION + "Always end your solution with the phrase 'the answer is' followed by your final answer. Start your solution with 'Stepstep_idx:'"

System Prompt for AIME24 and AIME25 Datasets

You are a helpful math assistant.

System Prompt for GPQA-diamond Dataset

You are a helpful assistant. Please answer "A", "B", "C", or "D".

External Model Prompt for Guided Search and LLM As a Critic

You are a teacher. Your task is to review and critique the paragraphs in solution directly. Output your judgment in the format of "`\boxed{Yes}`" if the paragraph is correct, or "`\boxed{No}`" if the paragraph is incorrect.

[Math Problem]
{problem}

[Solution]
{solution}

<paragraph_i>
{step_output}
</paragraph_i>

Evaluation Prompt for MATH-500, AIME24 and AIME25 Datasets

Question: {question}

Ground Truth Answer: {ground_truth}

Student Answer: {student_answer}

For the above question, please verify if the student's answer is equivalent to the ground truth answer. Do not solve the question by yourself; just check if the student's answer is equivalent to the ground truth answer. If the student's answer is correct, output Final Decision: Yes. If the student's answer is incorrect, output Final Decision: No.

python variables: 1) Step uncertainty, which is generated along with the reasoning text. 2) Momentum uncertainty, which is updated using step uncertainty based on Equation 5. After generating each step, we will check these two variables satisfy boundary condition in Equation 10, and trigger scaling if current step’s uncertainty is relatively higher than momentum uncertainty.

C More Experiment Results

C.1 Token Usage

We report the token usage of both the backbone and the external model in Table 3. There is no external model under ϕ -Decoding setting, so we only report the token usage under Guided Search and LLM As a Critic settings. In Table 1, *MUR* generates more tokens in some cases. This is because we only record the backbone token usage in Table 1. However, in Table 3, by adding up both backbone token usage and external model token usage, we can observe in the last column that *MUR* consistently generates fewer tokens than Per-Step Scale method, validating the token saving capacity of *MUR*. Furthermore, the trend of token usage of the Guided Search setting in Table 3 is compatible with those in Table 1.

C.2 Flexible Control with Hyperparameter γ

To further demonstrate the flexible control using hyperparameter γ , we report the detailed information concerning three model sizes and four test-time scaling methods (Guided Search, LLM As a Critic, ϕ -Decoding, thinking switch) on MATH-500 in Figure 8. It can be observed that by increasing γ , the reasoning accuracy would improve along with the token usage.

It is worth noting that in some scenarios, we observe performance degradation when we set γ to 0.9. This is consistent with our main findings: the reasoning performance drops with excessive reasoning token usage. In other words, we scale abundant steps in these scenarios. And the accuracy of Per-Step Scale method drops even lower with more token usage. Additionally, we observe that *MUR* outperforms Per-Step Scale in most scenarios. In practice, we set γ to 0.9 as the default.

C.3 Number of Steps

We report the number of steps generated by the backbone model and the number of scaled steps with *MUR* in Table 4. Additionally, we calcu-

late the percentage of scaled steps on each benchmark. For MATH-500, AIME24, AIME25, GPQA-diamond, the percentage is 8.38%, 9.34%, 12.54%, 13.75%, respectively. We can infer that among the same domain, more difficult benchmark leads to higher percentage of scaled steps. For example, AIME25 has higher scale percentage than AIME24 and MATH-500. Additionally, recent works (Fu et al., 2025; Xu et al., 2026; Wu et al., 2026b) reveals that LLM exhibits degraded performance on problems with more steps.

C.4 Impact of α

The hyperparameter α controls the update of momentum uncertainty, with a lower α leading to more intense updates. We report the impact of changing α in Figure 9. We can observe that *MUR* outperforms vanilla in most cases, which demonstrates the insensitivity and effectiveness of *MUR*. For $\alpha = 0.1$ setting, the momentum uncertainty changes too fast to well represent the overall estimation of query and generated steps, so the accuracy is relatively lower than other settings. In practice, we set $\alpha = 0.9$ as default.

C.5 Generalization to Larger Models and Different Model Family

In order to validate the generalization ability of *MUR*, we demonstrates the results on larger model (Qwen3-30B-A3B) and other model family (GLM-4-9B). From Table 5 we can observe that, *MUR* still performs better than all baselines.

C.6 Inference Time Analysis

The external detector’s implementation is based on several lines of code maintaining two float python variables, only requiring CPU computation, so its processing time is negligible. We include a comparison of the actual inference time across all baselines in Table 7, from which we can observe that *MUR* reduce time consumption in most scenarios. Besides, we analysis the p95/99 of long-tail situations in Table 6, indicating that our switching method performs comparatively on long-tail cases. Notably, our study primarily focuses on the academic exploration of achieving high reasoning performance with minimal token consumption. Many SOTA TTS methods, such as ToT (Yao et al., 2023), similarly encounter latency challenges, and real-world serving latency can be significantly mitigated through various engineering optimizations, such as

	MATH-500			AIME24			AIME25			GPQA-diamond			Avg.			$\Delta \downarrow$
	Bac \downarrow	Ext \downarrow	Sum \downarrow	Bac \downarrow	Ext \downarrow	Sum \downarrow	Bac \downarrow	Ext \downarrow	Sum \downarrow	Bac \downarrow	Ext \downarrow	Sum \downarrow	Bac \downarrow	Ext \downarrow	Sum \downarrow	
Qwen3-1.7B																
CoT	1,047	-	1,047	4,243	-	4,243	4,273	-	4,273	1,086	-	1,086	2,662	-	2,662	-
Guided search																
+ Per-Step Scale	3,460	3,186	6,646	17,463	21,607	39,070	16,680	18,212	34,892	6,739	9,258	15,997	11,086	13,066	24,151	-
+ Avg uncertainty	2,398	1,565	3,963	7,850	3,262	11,112	8,883	3,320	12,203	3,404	3,512	6,916	5,634	2,915	8,549	(-64.60%)
+ SMART	3,128	2,049	5,177	8,955	15,606	24,561	10,091	20,398	30,489	3,825	5,753	9,578	6,500	10,952	17,451	(-27.74%)
+ MUR (ours)	1,321	320	1,641	4,712	1,513	6,225	5,179	2,074	7,253	2,005	1,502	3,507	3,304	1,352	4,657	(-80.72%)
LLM as a critic																
+ Per-Step Scale	1,098	1,271	2,369	3,362	1,914	5,276	3,160	1,931	5,091	892	2,249	3,141	2,128	1,841	3,969	-
+ Avg uncertainty	1,019	1,075	2,094	4,176	542	4,718	3,174	769	3,943	1,417	2,001	3,418	2,447	1,097	3,543	(-10.73%)
+ SMART	878	670	1,548	3,976	1,241	5,217	3,600	1,486	5,086	1,446	763	2,209	2,475	1,040	3,515	(-11.44%)
+ MUR (ours)	902	337	1,239	3,892	853	4,745	4,011	828	4,839	1,693	1,282	2,975	2,625	825	3,450	(-13.09%)
Qwen3-4B																
CoT	772	-	772	3,111	-	3,111	2,577	-	2,577	612	-	612	1,768	-	1,768	-
Guided search																
+ Per-Step Scale	3,048	3,346	6,394	13,761	18,422	32,183	10,663	24,678	35,341	3,517	6,437	9,954	7,747	13,221	20,968	-
+ Avg uncertainty	1,911	1,845	3,756	7,012	4,422	11,434	7,719	4,076	11,795	1,354	2,483	3,837	4,499	3,207	7,706	(-63.25%)
+ SMART	2,476	2,212	4,688	8,515	15,623	24,138	9,375	14,199	23,574	2,116	3,409	5,525	5,621	8,861	14,481	(-30.94%)
+ MUR (ours)	824	265	1,089	4,265	2,042	6,307	7,162	13,985	21,147	929	641	1,570	3,295	4,233	7,528	(-64.10%)
LLM as a critic																
+ Per-Step Scale	777	1,373	2,150	3,334	2,040	5,374	3,260	1,885	5,145	737	2,462	3,199	2,027	1,940	3,967	-
+ Avg uncertainty	741	957	1,698	3,217	1,052	4,269	3,120	1,002	4,122	804	1,795	2,599	1,971	1,202	3,172	(-20.04%)
+ SMART	813	855	1,668	3,203	1,315	4,518	3,201	1,485	4,686	724	320	1,044	1,985	994	2,979	(-24.91%)
+ MUR (ours)	745	443	1,188	3,309	895	4,204	3,113	980	4,093	699	266	965	1,967	646	2,613	(-34.14%)
Qwen3-8B																
CoT	1,131	-	1,131	4,077	-	4,077	4,746	-	4,746	859	-	859	2,703	-	2,703	-
Guided search																
+ Per-Step Scale	4,069	3,688	7,757	19,805	23,308	43,113	21,586	23,227	44,813	4,252	7,468	11,720	12,428	14,423	26,851	-
+ Avg uncertainty	2,427	2,037	4,464	11,223	5,358	16,581	12,193	6,449	18,642	2,213	3,382	5,595	7,014	4,307	11,321	(-57.84%)
+ SMART	3,502	3,287	6,789	17,055	24,194	41,249	17,705	24,403	42,108	3,797	6,135	9,932	10,515	14,505	25,020	(-6.82%)
+ MUR (ours)	2,607	1,986	4,593	7,959	4,196	12,155	7,582	4,603	12,185	3,122	4,524	7,646	5,318	3,827	9,145	(-65.94%)
LLM as a critic																
+ Per-Step Scale	1,022	2,025	3,047	4,846	2,258	7,104	4,818	2,381	7,199	1,172	3,102	4,274	2,965	2,442	5,406	-
+ Avg uncertainty	1,086	842	1,928	5,326	1,105	6,431	4,705	1,205	5,910	1,375	1,588	2,963	3,123	1,185	4,308	(-20.31%)
+ SMART	1,167	1,160	2,327	4,737	1,547	6,284	4,780	1,945	6,725	1,069	2,366	3,435	2,938	1,755	4,693	(-13.19%)
+ MUR (ours)	1,132	783	1,915	4,846	1,014	5,860	4,913	1,237	6,150	1,007	2,211	3,218	2,975	1,311	4,286	(-20.72%)

Table 3: Token usage of both backbone and external model. **Bac** stands for backbone model, **Ext** stands for external model, and the sum of them is denoted as **Sum**. \downarrow means better for lower values.

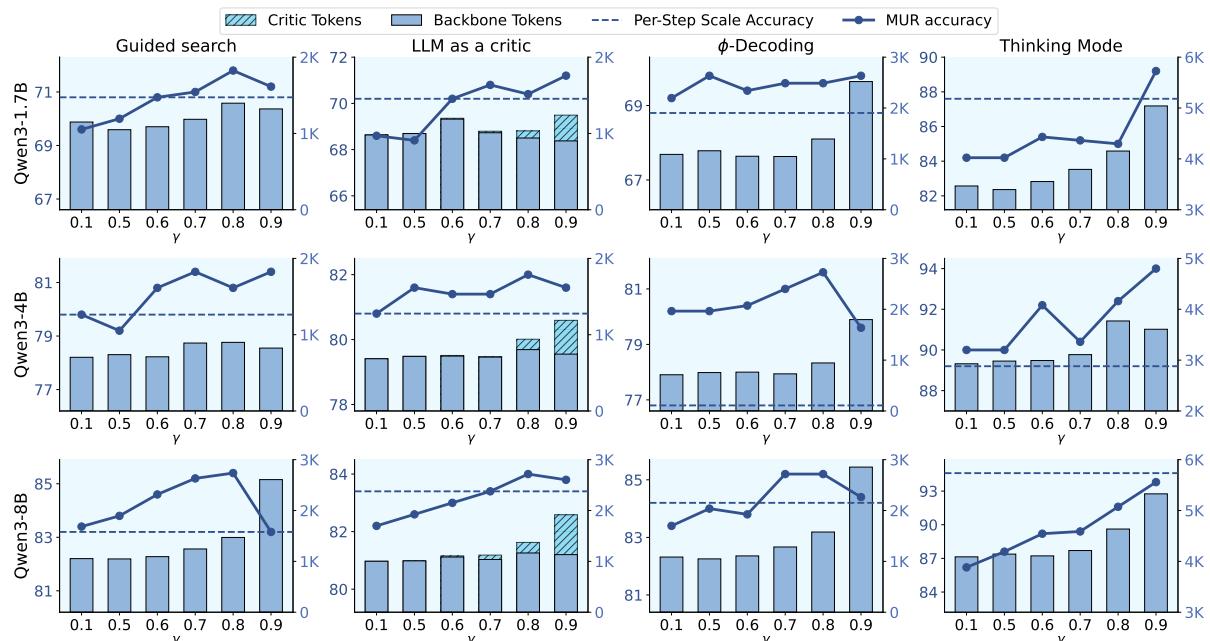


Figure 8: Detail scaling law of γ . The X axis stands for different values of γ . The Y axis stands for accuracy. Due to the reason described in Appendix C.1, we additionally report the external model token usage (denoted as Critic Tokens) under LLM as a critic setting to comprehensively reflect the overall computes.

KV-cache management and chunking, which fall beyond the scope of our paper.



Figure 9: Impact of changing α . The X axis stands for different values of α . The Y axis stands for accuracy.

C.7 Case Study

In Figure 10, we conduct a case study based on the thinking mode of Qwen3-1.7B. We analyze AIME24 and show the comparison between *MUR* and Per-Step Scale. We can observe that when *MUR* faces high uncertainty step, it triggers the thinking process, allocating more computes to optimize current step's quality. For simple steps showing low uncertainty, *MUR* directly output it without thinking. On the contrary, Per-Step Scale thinks for every step, regardless of whether the backbone is confident with the current step or not.

Datasets	MATH-500		AIME24		AIME25		GPQA-diamond		Avg	
	Total	Scaled	Total	Scaled	Total	Scaled	Total	Scaled	Total	Scaled
Qwen3-1.7B										
CoT	5.33	-	8.93	-	8.40	-	6.41	-	7.27	-
Guided search										
+ Per-Step Scale	2.89	2.89	4.20	4.20	3.37	3.37	4.55	4.55	3.75	3.75
+ MUR (ours)	5.31	0.35	6.93	0.70	7.13	0.80	7.02	0.82	6.60	0.67
LLM as a critic										
+ Per-Step Scale	4.35	4.35	3.63	3.63	3.60	3.60	3.93	3.93	3.88	3.88
+ MUR (ours)	5.86	0.40	7.57	0.60	6.87	0.83	5.77	0.81	6.52	0.66
ϕ -Decoding										
+ Per-Step Scale	2.97	2.97	5.33	5.33	2.90	2.90	3.91	3.91	3.78	3.78
+ MUR (ours)	5.80	0.39	7.47	0.93	6.57	0.60	5.59	0.76	6.35	0.67
Qwen3-4B										
CoT	5.84	-	5.70	-	5.00	-	5.57	-	5.53	-
Guided search										
+ Per-Step Scale	2.73	2.73	3.17	3.17	3.07	3.07	2.71	2.71	2.92	2.92
+ MUR (ours)	4.31	0.17	5.97	0.63	4.43	0.53	3.59	0.30	4.57	0.41
LLM as a critic										
+ Per-Step Scale	3.83	3.83	4.23	4.23	3.77	3.77	2.47	2.47	3.58	3.58
+ MUR (ours)	4.36	0.18	5.03	0.47	3.63	0.63	3.31	0.35	4.08	0.41
ϕ -Decoding										
+ Per-Step Scale	2.77	2.77	3.97	3.97	4.23	4.23	3.10	3.10	3.52	3.52
+ MUR (ours)	4.38	0.19	5.40	0.47	4.30	0.43	3.89	1.02	4.49	0.53
Qwen3-8B										
CoT	7.45	-	10.00	-	12.33	-	6.90	-	9.17	-
Guided search										
+ Per-Step Scale	3.27	3.27	4.80	4.80	4.73	4.73	3.83	3.83	4.16	4.16
+ MUR (ours)	5.32	0.40	7.80	0.80	6.13	0.97	5.20	0.69	6.11	0.71
LLM as a critic										
+ Per-Step Scale	5.01	5.01	5.13	5.13	6.10	6.10	3.92	3.92	5.04	5.04
+ MUR (ours)	5.93	0.55	7.30	0.87	7.13	0.80	5.17	0.67	6.38	0.72
ϕ -Decoding										
+ Per-Step Scale	3.20	3.20	4.33	4.33	5.33	5.33	3.45	3.45	4.08	4.08
+ MUR (ours)	4.45	1.20	8.33	0.77	6.60	1.03	4.32	2.11	5.93	1.28

Table 4: Detailed step data for scaled steps and total steps.

	MATH-500		AIME24		AIME25		GPQA-diamond		Avg.			
	Acc.↑	#Tokens↓	Acc.↑	#Tokens↓	Acc.↑	#Tokens↓	Acc.↑	#Tokens↓	Acc.↑	Δ↑	#Tokens↓	Δ↓
Qwen3-30B-A3B												
Vanilla CoT	83.40	606	26.67	3,104	20.00	3,182	43.94	579	43.50	-	1,868	-
Guided search												
+ Per-Step Scale	85.20	7,345	36.67	59,935	23.33	53,974	46.47	2,350	47.92	-	30,901	-
+ Avg uncertainty	84.20	3,779	46.67	28,139	33.33	28,933	47.47	1,044	52.92	(+5.00)	15,474	(-49.92%)
+ SMART	85.40	5,601	40.00	39,829	30.00	39,232	45.96	1,879	50.34	(+2.42)	21,635	(-29.99%)
+ MUR (ours)	84.60	1,564	53.33	27,788	30.00	27,622	49.49	855	54.36	(+6.44)	14,457	(-53.21%)
LLM as a critic												
+ Per-Step Scale	84.80	2,985	43.33	13,630	23.33	16,335	45.96	731	49.36	-	8,420	-
+ Avg uncertainty	84.40	2,502	40.00	14,553	26.67	15,143	43.94	676	48.75	(-0.61)	8,219	(-2.39%)
+ SMART	83.40	1,057	43.33	16,118	30.00	14,661	47.98	626	51.18	(+1.82)	8,116	(-3.61%)
+ MUR (ours)	86.60	1,594	43.33	14,622	30.00	13,394	45.45	706	51.35	(+1.99)	7,579	(-9.99%)
ϕ -Decoding												
+ Per-Step Scale	82.60	3,609	43.33	129,993	33.33	113,040	42.93	3,289	50.55	-	62,483	-
+ Avg uncertainty	81.20	1,990	43.33	48,869	30.00	48,246	47.98	1,162	50.63	(+0.08)	25,067	(-59.88%)
+ SMART	82.80	3,052	40.00	91,896	33.33	82,686	41.41	2,163	49.39	(-1.16)	44,949	(-28.06%)
+ MUR (ours)	86.00	1,700	43.33	51,730	30.00	54,023	48.48	1,578	51.95	(+1.40)	27,258	(-56.37%)
GLM-4-9B												
Vanilla CoT	45.20	473	6.67	714	3.33	820	28.79	608	21.00	-	654	-
Guided search												
+ Per-Step Scale	56.60	3,494	10.00	17,566	6.67	14,876	30.30	3,569	25.89	-	9,876	-
+ Avg uncertainty	48.60	1,196	13.33	9,300	3.33	8,709	27.27	1,785	23.13	(-2.76)	5,248	(-46.86%)
+ SMART	53.40	2,304	20.00	14,395	3.33	15,706	31.82	3,029	27.14	(+1.25)	8,859	(-10.30%)
+ MUR (ours)	49.60	901	20.00	11,107	10.00	10,137	32.83	1,574	28.11	(+2.22)	5,930	(-39.96%)
LLM as a critic												
+ Per-Step Scale	53.20	1,037	13.33	12,038	13.33	11,561	25.25	1,080	26.28	-	6,429	-
+ Avg uncertainty	56.00	1,008	20.00	6,082	6.67	7,694	31.31	855	28.50	(+2.22)	3,910	(-39.18%)
+ SMART	48.60	579	20.00	8,993	10.00	10,998	25.75	711	26.09	(-0.19)	5,320	(-17.25%)
+ MUR (ours)	53.20	746	23.33	8,786	13.33	10,373	29.80	890	29.92	(+3.64)	5,199	(-19.13%)
ϕ -Decoding												
+ Per-Step Scale	44.60	3,817	16.67	26,268	13.33	56,546	25.25	3,495	24.96	-	22,532	-
+ Avg uncertainty	47.80	1,316	23.33	21,796	6.67	19,333	29.80	1,466	26.90	(+1.94)	10,978	(-51.28%)
+ SMART	47.00	2,386	20.00	41,494	16.67	47,241	28.79	3,791	28.12	(+3.16)	23,728	(+5.31%)
+ MUR (ours)	48.00	1,531	23.33	26,966	10.00	24,520	31.31	2,466	28.16	(+3.20)	13,871	(-38.44%)

Table 5: Generalization results. The best results are highlighted in bold. **Acc.** denotes pass@1 rate and **#Tokens** denotes the **backbone model's** average token usage for each query. We also report the delta compared to *Per-Step Scale* baseline. **Red** indicates worse performance, while **green** indicates better performance against *Per-Step Scale*. Here, \uparrow denotes that higher values are better, whereas \downarrow means lower values are preferable.

Datasets	MATH-500		AIME24		AIME25		GPQA-diamond		Avg	
	p95	p99	p95	p99	p95	p99	p95	p99	p95	p99
Qwen3-1.7B										
CoT	9.18	18.49	62.26	96.76	30.20	34.30	16.26	96.43	29.48	80.31
Thinking Mode	84.90	110.07	207.80	211.93	198.57	208.84	61.81	79.13	138.27	152.49
Average	83.19	170.68	209.72	226.03	187.02	213.56	80.1	83.73	140.11	173.50
SMART	102.62	158.70	178.81	206.74	189.19	195.11	75.20	116.76	136.46	169.33
MUR (ours)	108.49	147.09	194.67	214.09	178.07	195.95	54.14	71.86	133.84	157.25
Qwen3-4B										
CoT	11.79	34.31	38.12	122.84	37.15	121.67	10.03	26.50	24.27	76.33
Thinking Mode	179.90	242.26	286.85	303.02	367.77	367.90	129.48	165.55	241.01	269.68
Average	153.23	181.07	280.18	311.06	292.32	327.45	178.04	219.77	225.94	259.84
SMART	138.16	185.20	365.05	373.67	326.77	332.56	139.20	214.63	242.30	276.52
MUR (ours)	105.55	147.53	264.49	324.71	308.94	343.42	82.42	128.03	190.35	235.93

Table 6: Tail latency metrics, formatted as p95 and p99.

Method	Variant	MATH	AIME24	AIME25	GPQA-diamond	Avg.	Δ
Qwen3-1.7B							
CoT		3.94	15.58	15.69	4.43	9.91	-
Guided Search	Per-Step Scale	12.80	92.47	73.34	31.82	52.61	-
	Avg uncertainty	8.68	31.67	31.67	17.62	22.41	(-30.20)
	SMART	10.57	62.77	55.03	21.11	37.37	(-15.24)
	MUR	4.96	27.25	25.11	9.89	16.80	(-35.81)
LLM As a Critic	Per-Step Scale	8.35	23.24	19.62	10.63	15.46	-
	Avg uncertainty	7.32	27.19	21.13	12.11	16.94	(+1.48)
	SMART	5.38	48.20	44.93	8.14	26.66	(+11.20)
	MUR	4.49	30.16	25.40	10.96	17.75	(+2.29)
ϕ -Decoding	Per-Step Scale	15.52	65.28	62.56	30.14	43.38	-
	Avg uncertainty	9.24	36.33	38.61	10.03	23.55	(-19.83)
	SMART	11.26	62.73	50.21	13.61	34.45	(-8.93)
	MUR	8.12	35.21	36.60	7.29	21.81	(-21.57)
Qwen3-4B							
CoT		9.46	22.60	21.90	8.57	15.63	-
Guided Search	Per-Step Scale	31.12	134.40	125.80	86.06	94.35	-
	Avg uncertainty	21.83	114.83	76.80	37.66	62.78	(-31.57)
	SMART	19.34	134.07	181.07	51.05	96.38	(+2.03)
	MUR	10.69	53.93	46.50	25.31	34.11	(-60.24)
LLM As a Critic	Per-Step Scale	18.58	104.63	70.90	42.78	59.22	-
	Avg uncertainty	15.65	64.27	71.17	40.34	47.86	(-11.36)
	SMART	15.94	49.47	61.63	21.98	37.26	(-21.96)
	MUR	12.84	63.43	56.40	20.98	38.41	(-20.81)
ϕ -Decoding	Per-Step Scale	42.59	128.60	123.57	85.27	95.01	-
	Avg uncertainty	19.70	109.30	74.53	36.23	59.94	(-35.07)
	SMART	28.75	149.07	137.20	52.68	91.93	(-3.08)
	MUR	20.83	94.87	92.73	23.79	58.06	(-36.95)
Qwen3-8B							
CoT		27.18	68.28	65.92	20.13	45.38	-
Guided Search	Per-Step Scale	59.32	351.58	338.58	79.32	207.20	-
	Avg uncertainty	43.28	188.90	179.73	47.31	114.81	(-92.39)
	SMART	55.64	214.95	207.56	70.48	137.16	(-70.04)
	MUR	46.68	173.93	172.60	60.86	113.52	(-93.68)
LLM As a Critic	Per-Step Scale	36.41	123.18	108.81	47.08	78.87	-
	Avg uncertainty	30.61	129.33	122.64	32.60	78.80	(-0.07)
	SMART	35.12	125.73	107.29	40.23	77.09	(-1.78)
	MUR	31.38	120.03	100.02	37.45	72.22	(-6.65)
ϕ -Decoding	Per-Step Scale	87.90	397.23	390.91	84.37	240.10	-
	Avg uncertainty	58.35	288.96	267.33	41.22	163.97	(-76.13)
	SMART	80.20	351.19	306.03	91.84	207.32	(-32.78)
	MUR	51.32	286.79	253.98	48.56	160.16	(-79.94)

Table 7: Average time Comparison for baselines and MUR.

Question: Eight circles of radius 34 are sequentially tangent, and two of the circles are tangent to AB and BC of triangle ABC, respectively. 2024 circles of radius 1 can be arranged in the same manner. The inradius of triangle ABC can be expressed as m/n , where m and n are relatively prime positive integers. Find $m+n$.

Ground Truth: 197

MUR

Step1: Okay, so I need to find the inradius of triangle ABC...

Step2(Before Scale): First, I need to visualize the problem ... But they are tangent to AB and BC.
Hmm.

NEED SCALING!!!

[$\exp(m)=1.32, \exp(M)/\gamma=1.26$]

Step2(After Scale): \langle think \rangle
Alternatively, maybe they are arranged in a way similar to the incircle and excircle? But with circles of radius 34 ... Therefore, maybe the circles are tangent to AB or BC, and the other side?

\langle /think \rangle

Step3: The problem says "eight circles of radius 34 are sequentially tangent, and two of the circles are tangent to AB and BC of triangle ABC, respectively..."

...
[omit some intermediate steps]

Step8(Before Scale): But **I don't know theta**, so this seems stuck.

NEED SCALING!!!

[$\exp(m)=1.56, \exp(M)/\gamma=1.37$]

Step8(After Scale): \langle think \rangle
But if we have two different arrangements ... Therefore, the inradius is $192/5$, so $m + n = 192 + 5 = 197$. Therefore, the inradius is $192/5$, so $m + n = 192 + 5 = 197$.

Final Answer

197

\langle /think \rangle

To determine the inradius of triangle ABC, we analyze the geometric arrangement of circles of two different radii ... Final Answer
 $m = 192, n = 5, m + n = \underline{197}$

Per-Step Scale

Step1(Before Scale): Okay, so I need to find the inradius of triangle ABC...

Step1(After Scale): \langle think \rangle
First, let me try to visualize the problem. There's triangle ABC ... Because in the right triangle formed by B, the center, and the foot of the perpendicular from the center to AB, we have $\sin(\theta/2) = \text{opposite} / \text{hypotenuse} = 34 / d$.
 \langle /think \rangle

Step2(Before Scale): So $d = 34 / \sin(\theta/2)$.

Step2(After Scale): \langle think \rangle
If there are two circles tangent to AB and BC, and they are tangent to each other ... if the circles are arranged in a triangle with inradius R, then maybe $R = r * k$, where k is the scaling factor. If we can find k, then we can compute R.
 \langle /think \rangle

Step3(Before Scale): **But how?**

Step3(After Scale): \langle think \rangle
Alternatively, think of the problem as similar to a triangle with inradius R, and circles of radius r arranged inside it ... But since the last circle is also tangent to AB and BC, its distance from B is $d_n = 34 / \sin(\theta/2)$.
 \langle /think \rangle

...
[omit some intermediate steps]

Step5(Before Scale): Think of the triangle as having a certain inradius R, and each circle of radius r is placed such that it touches two sides and the inradius...

Step5(After Scale): \langle think \rangle
But this is too vague ... Therefore, this is not possible. Given that I'm stuck, and the only possible answer I can think of is $34/253$, leading to $m + n = 287$, I will go with that.

Final Answer

287

\langle /think \rangle

To solve this problem, we are given two sets of circles... Thus,
 $m + n = 34 + 253 = \underline{287}$

Figure 10: Case study. **Red** denotes the backbone faces high uncertain step. **Green** denotes key words of reflecting.