

GlimpRouter: Efficient Collaborative Inference by Glimpsing One Token of Thoughts

Wenhao Zeng[†] Xuteng Zhang[†] Yuling Shi Chao Hu
Yuting Chen Beijun Shen Xiaodong Gu*

Shanghai Jiao Tong University

{zengwh_cs, scottzhang, xiaodong.gu}@sjtu.edu.cn

Abstract

Large Reasoning Models (LRMs) achieve remarkable performance by explicitly generating multi-step chains of thought, but this capability incurs substantial inference latency and computational cost. Collaborative inference offers a promising solution by selectively allocating work between lightweight and large models, yet a fundamental challenge remains: determining when a reasoning step requires the capacity of a large model or the efficiency of a small model. Existing routing strategies either rely on local token probabilities or post-hoc verification, introducing significant inference overhead. In this work, we propose a novel perspective on step-wise collaboration: the difficulty of a reasoning step can be inferred from its very first token. Inspired by the “Aha Moment” phenomenon in LRMs, we show that the entropy of the initial token serves as a strong predictor of step difficulty. Building on this insight, we introduce GlimpRouter, a training-free step-wise collaboration framework. GlimpRouter employs a lightweight model to generate only the first token of each reasoning step and routes the step to a larger model only when the initial token entropy exceeds a threshold. Experiments on multiple benchmarks demonstrate that our approach significantly reduces inference latency while preserving accuracy. For instance, GlimpRouter attains a substantial 10.7% improvement in accuracy while reducing inference latency by 25.9% compared to a standalone large model on AIME25. These results suggest a simple yet effective mechanism for reasoning: allocating computation based on a glimpse of thought rather than full-step evaluation¹.

[†]Equal contribution.

*Corresponding author

¹Code and dataset are available at <https://github.com/Zengwh02/GlimpRouter>.

1 Introduction

Large reasoning models, such as DeepSeek-R1 (Guo et al., 2025) and OpenAI o1/o3 (Jaech et al., 2024; OpenAI, 2025), have demonstrated remarkable performance across a wide range of complex reasoning tasks by explicitly generating structured reasoning steps (Yang et al., 2025c; Wang et al., 2025b,a; Chen et al., 2025a; Li et al., 2025; Hu et al., 2025; Zhang et al., 2025b,c). However, this capability comes at a high cost: extended reasoning chains inevitably incur high inference latency and substantial computational overhead, severely limiting the practicality of LRMs in latency-sensitive and resource-constrained settings.

To alleviate this bottleneck, collaborative inference has emerged as a promising paradigm (She et al., 2025; Xiao et al., 2025; Qu et al., 2025). The key idea is that multiple models with heterogeneous capabilities and costs are orchestrated to jointly solve a task: lightweight models handle routine or easy tasks, while more powerful but expensive models are invoked selectively for difficult ones (Chen et al., 2025b; Wang et al., 2025d). Existing collaboration strategies generally fall into two categories. Token-level methods, such as Speculative Decoding (Leviathan et al., 2023; Fu et al., 2025), accelerate generation by allowing a small model to propose candidate tokens that are then verified (and possibly accepted) by a larger model, reducing the number of expensive forward passes. In contrast, step-level methods attempt to route entire reasoning steps between models. They typically depend on *post-hoc* verification, such as LLM-as-a-judge (Shi et al., 2025; Pan et al., 2025), or on averaged uncertainty metrics (e.g., mean perplexity or entropy) (Cui et al., 2025; Zhang et al., 2025a).

Despite their promise, determining the optimal allocation of tasks between large and small models remains a fundamental challenge. Token-level methods rely primarily on local token probabilities

and frequently switch models at a fine granularity, while step-level strategies require evaluating entire reasoning steps before making routing decisions. Both approaches introduce non-trivial computational overhead, which can partially or even fully offset their intended efficiency gains.

In this work, we advance a new perspective: *Can we judge the difficulty of a reasoning step at the very beginning?* Inspired by the “Aha Moment” phenomenon observed in LRMs (Guo et al., 2025), previous works suggest that the onset of a reasoning step—often marked by discourse cues such as “Wait”, “But”, or “So”—represents a cognitive pivot that shapes the entire subsequent reasoning trajectory (Wang et al., 2025c; Yang et al., 2025b; Zeng et al., 2025b). We hypothesize that the information associated with this *initial token* provides a more discriminative signal of step-level difficulty than averaged statistics over the full step. To validate this hypothesis, we analyze the uncertainty distribution of the reasoning steps across various tasks (Section 2). The experimental results show that the entropy of the initial tokens exhibits an exceptionally high variance, indicating it is a powerful discriminator: steps with low initial entropy typically correspond to routine derivations that a small model can readily process, whereas steps with high initial entropy indicate critical cognitive bifurcations that require the collaboration of larger, more powerful models.

Building on this insight, we propose **GlimpRouter**, a novel training-free framework for efficient collaborative inference. At each reasoning step, a lightweight model is first used to generate only the initial token of the step. If the entropy of this token falls below a predefined threshold, the lightweight model proceeds to generate the entire step; otherwise, the generated context is seamlessly handed over to a larger and more capable model. In effect, this mechanism essentially operationalizes the principle of “A Glimpse of Thought”: a minimal signal from the onset of a reasoning step governs the computational budget allocated to its full generation.

We empirically evaluate our approach on a diverse set of reasoning benchmarks, including mathematical reasoning (AIME), challenging general reasoning (GPQA (Rein et al., 2024)), and code generation (LiveCodeBench (Jain et al., 2024)). Across all benchmarks, GlimpRouter achieves a superior trade-off between efficiency and perfor-

mance. For instance, on the AIME25 benchmark, our method attains a substantial **10.7%** improvement in accuracy while simultaneously reducing inference latency by **25.9%** compared to a standalone large model. Moreover, we demonstrate that our step-level routing strategy is orthogonal to token-level speculative decoding, enabling complementary and compound speedups when integrated.

Our contributions are summarized as follows:

- We perform an analysis of the uncertainty distribution across reasoning steps and identify that the entropy of the initial token serves as an indicator of reasoning difficulty.
- We propose **GlimpRouter**, a simple yet effective mechanism that enables training-free, step-aware collaboration between models.
- Extensive experiments verify that our method significantly reduces latency while maintaining or even enhancing the reasoning efficacy of large models, offering a practical solution for deploying efficient LRMs.

2 A Preliminary Study

Previous works have suggested that the onset of a reasoning step—often marked by discourse cues such as “Wait”, “But”, or “So”—represents a cognitive pivot that shapes the entire subsequent reasoning trajectory (Wang et al., 2025c; Yang et al., 2025b; Zeng et al., 2025b; Zhang et al., 2025d). Building on this intuition, we hypothesize that the uncertainty associated with the *initial token* of a reasoning step provides a more discriminative signal of step-level difficulty than statistics aggregated over the full step. In other words, the essential difficulty information is concentrated at the very beginning of each reasoning step. In this section, we analyze the uncertainty distribution for the reasoning steps and examine how different uncertainty metrics correlate with step-wise difficulty.

2.1 Study Design

To quantify uncertainty within a reasoning step, we analyze entropy-based metrics derived from the model’s conditional probability distribution P_θ . It is well-established that the internal logits of LLMs serve as a reliable proxy for model confidence (Mallin and Gales, 2020; Kuhn et al., 2023). Let $s_k = \{t_1, \dots, t_L\}$ denote the k -th reasoning step with L tokens, and c_k be the preceding context of this step.

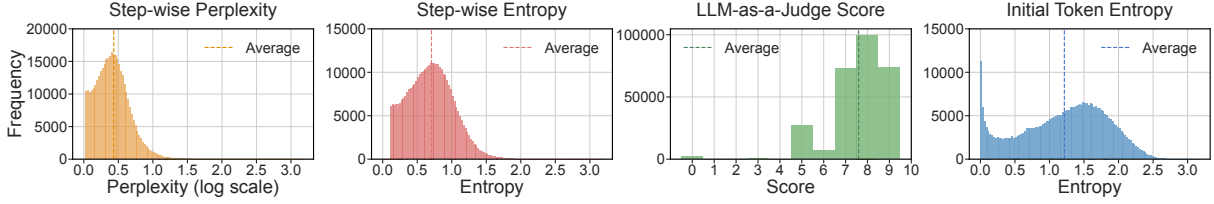


Figure 1: Comparison of the distributions of different uncertainty metrics. Unlike other metrics, which exhibit a generic unimodal distribution, the entropy of initial thought tokens (\mathbf{H}_{init}) displays a distinct bimodal and heavy-tailed distribution. This indicates that \mathbf{H}_{init} serves as a discriminative signal, effectively capturing the ‘‘Aha Moments’’ that distinguish complex reasoning steps from routine derivations.

Step Entropy (\mathbf{H}_{step}). \mathbf{H}_{step} quantifies the average uncertainty across all tokens in a reasoning step:

$$\mathbf{H}_{\text{step}}(s_k) = \frac{1}{L} \sum_{i=1}^L \mathbf{H}(P_{\theta}(\cdot|t_{<i}, \mathbf{c}_k)) \quad (1)$$

where $\mathbf{H}(p) = -\sum_{v \in V} p(v) \log p(v)$ denotes the Shannon entropy of a distribution (Shannon, 1948). This metric treats all tokens within the step uniformly by averaging uncertainty across the entire sequence.

Initial Token Entropy (\mathbf{H}_{init}). In contrast, \mathbf{H}_{init} captures the uncertainty of the *initial token* for each step:

$$\mathbf{H}_{\text{init}}(s_k) = \mathbf{H}(P_{\theta}(t_1|\mathbf{c}_k)) \quad (2)$$

Unlike step-level averages, \mathbf{H}_{init} isolates the uncertainty at the onset of the step, which we hypothesize to be an indicator of the step’s difficulty.

As additional baselines, we consider the following commonly used metrics in collaborative inference.

Step-wise Perplexity ($\mathbf{PPL}_{\text{step}}$). $\mathbf{PPL}_{\text{step}}$ measures the exponentiated average negative log-likelihood over the entire step:

$$\mathbf{PPL}_{\text{step}}(s_k) = \exp\left(-\frac{1}{L} \sum_{i=1}^L \log P_{\theta}(t_i|t_{<i}, \mathbf{c}_k)\right) \quad (3)$$

Like \mathbf{H}_{step} , $\mathbf{PPL}_{\text{step}}$ aggregates uncertainty uniformly across all tokens.

LLM-as-a-Judge. It assesses the validity of a reasoning step using the large model directly. Given a step s_k and its context \mathbf{c}_k , the LLM assigns a scalar score $S \in \{0, 1, \dots, 9\}$. We note that this metric requires an additional computationally expensive inference pass and is therefore unsuitable for efficiency-critical routing decisions.

We examine the reasoning traces using Qwen3-4B, Qwen3-32B (Yang et al., 2025a), and

DeepSeek-R1-Distill-Qwen-32B (Guo et al., 2025) on the AIME and LiveCodeBench datasets (Jain et al., 2024), collecting over 10 million tokens of reasoning steps.

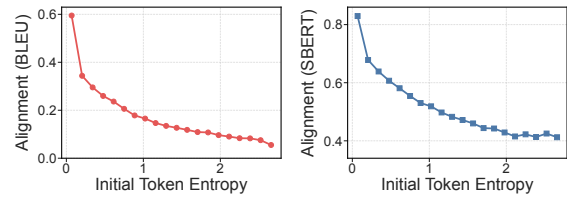


Figure 2: Alignment between the outputs generated by small and large models under various \mathbf{H}_{init} intervals.

2.2 Results and Analysis

Information Distribution within Reasoning Steps.

Figure 1 illustrates the distributions reflected by the four metrics. It is evident that both $\mathbf{PPL}_{\text{step}}$ and \mathbf{H}_{step} exhibit narrow, unimodal distributions clustered around central values. We attribute this to the *signal dilution* effect: a reasoning step typically contains a few critical decision tokens followed by a long sequence of deterministic syntactic tokens. Averaging over the entire step smooths out the local spikes of uncertainty, making these metrics insensitive to the actual difficulty. Similarly, the LLM-as-a-Judge distribution is discrete and highly skewed towards high scores (saturation), lacking the granularity required for fine-grained routing thresholds. In contrast, \mathbf{H}_{init} presents a distinct bimodal distribution with a heavy tail. The clear separation between the low-entropy peak (routine steps) and the high-entropy tail (complex reasoning) suggests that \mathbf{H}_{init} naturally functions as a high-sensitivity discriminator for routing.

Correlations to Model Routing. Motivated by the observation that most information in a reasoning step resides at its onset, we further examine whether this signal can serve as an indicator for model routing, where a lightweight language model

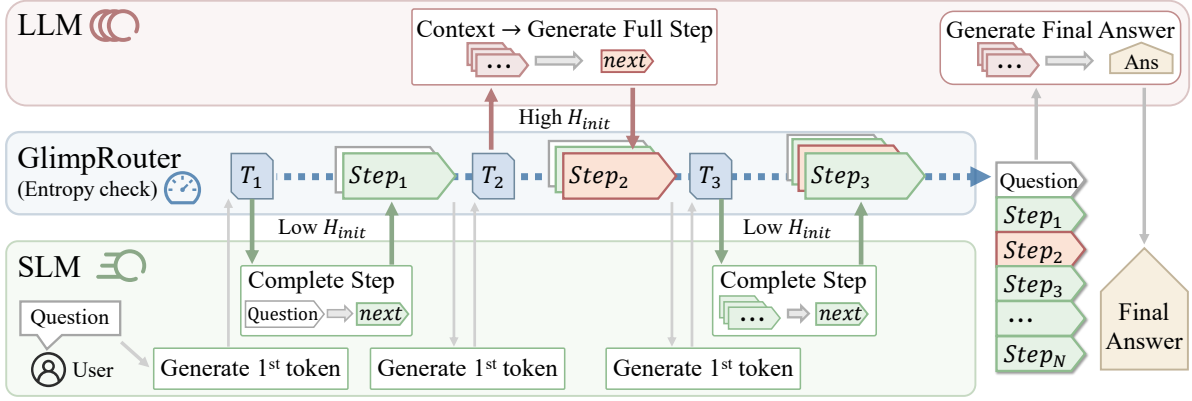


Figure 3: The overall framework of **GlimpRouter**. The collaborative inference pipeline operates in a step-wise manner. At the onset of each reasoning step, the SLM first generates a “glimpse” (the initial token T_k). GlimpRouter then computes the entropy H_{init} of this token to gauge step difficulty. Based on the threshold, the system dynamically routes the generation: (1) Delegate (Low H_{init}): Routine steps are fully generated by the efficient SLM. (2) Intervene (High H_{init}): Complex steps signaling cognitive pivots are routed to the capable LLM for robust reasoning. Finally, the LLM generates the conclusive answer based on the aggregated collaborative chain.

can reliably substitute a large model. To this end, we partition reasoning steps into intervals based on H_{init} predicted by the small model. For each interval, we evaluate the alignment between the outputs generated by the small and large models, conditioned on the same preceding context c . We quantify alignment using two widely used similarity metrics: BLEU-4 (Lin and Och, 2004), which captures lexical overlap, and SBERT (Reimers and Gurevych, 2019), which measures semantic similarity. As shown in Figure 2, we observe a strictly monotonic negative correlation between H_{init} and both similarity metrics. In the low- H_{init} regime, the small model’s outputs closely align with those of the large model, indicating that such steps can be handled competently by the lightweight model. In contrast, as H_{init} increases, the alignment degrades substantially, reflecting growing divergence in both surface form and semantic content. This strong correlation provides empirical evidence that H_{init} can be a reliable predictor of step-level difficulty, thereby suggesting a routing strategy in which a small model can substitute a large model in low- H_{init} steps without compromising results. To provide concrete intuition for this correlation, we present qualitative case studies in Appendix F.

3 Methodology

3.1 Problem Formulation

We consider a collaborative inference framework involving a large, high-capacity reasoning model (LLM, M_L) and a small computationally efficient model (SLM, M_S). Large reasoning models typi-

cally output a reasoning process \mathcal{T} , often encapsulated within a `<think> ... </think>` tag, followed by the final answer A . Our framework focuses on accelerating the reasoning process \mathcal{T} , while the final answer A is always generated by M_L to ensure correctness. Formally, the reasoning process is decomposed into a sequence of steps $\mathcal{T} = \{s_1, \dots, s_K\}$ based on structural delimiters (e.g., double newlines (Pan et al., 2025; Yang et al., 2025b)). Each step s_k consists of a sequence of tokens $s_k = (t_1, \dots, t_L)$ and is generated conditioned on the preceding context c_k , which includes the original question and all previous generated steps s_1, \dots, s_{k-1} . The objective of collaborative inference is to dynamically assign each reasoning step s_k to M_S or M_L , minimizing overall inference latency while preserving the quality of reasoning required for accurate solutions.

3.2 Overview

To navigate the efficiency-accuracy trade-off, we propose **GlimpRouter**, a training-free, step-aware collaboration strategy that routes reasoning steps based on the *initial token entropy* (H_{init}). Figure 3 illustrates the overall pipeline. Given an input question, the system operates in a step-wise manner. At the onset of each step, instead of blindly generating the full content, GlimpRouter employs M_S to “glimpse” the first token, yielding an entropy H_{init} that quantizes the difficulty of the upcoming step: If H_{init} falls below a threshold, the step is deemed routine, and the small model is responsible for continuing the generation. Otherwise, the

step signals a cognitive pivot, and the context is handed over to a large model for high-quality generation. This “Probe-then-Dispatch” mechanism ensures that heavy computational resources are allocated solely for critical reasoning steps. Each step is introduced in the following sections. The detailed procedure is outlined in Algorithm 1 in Appendix A.

3.3 Glimpse: Initial Token Probing

At the beginning of the reasoning step k , given the preceding context c_k , the small model M_S is invoked to predict the probability distribution of only the first token $t_{k,1}$. We compute the entropy $H_{\text{init}}(s_k)$ of this distribution as a proxy for the cognitive uncertainty of the current reasoning step. This probing operation incurs a marginal computational cost equivalent to decoding a single token, unlike methods that necessitate generating an entire reasoning step (typically $L \gg 1$ tokens) before verification. Even in scenarios where the step is subsequently identified as difficult and routed to M_L (thereby discarding the probe), this 1-token overhead is negligible compared to the substantial sunk costs associated with discarding fully generated invalid steps.

3.4 Dynamic Model Routing

Given the entropy of the initial token in the current step, GlimpRouter dispatches the generation task to the appropriate model based on a threshold τ :

Delegate ($H_{\text{init}}(s_k) \leq \tau$): Low entropy suggests that the small model is confident in the logical progression, and the step is likely routine. Consequently, M_S continues to autoregressively generate the remainder of step s_k until the delimiter is reached. This decision reduces the total cost by maximizing the utilization of the small model.

Intervene ($H_{\text{init}}(s_k) > \tau$): High entropy indicates logical ambiguity or high cognitive load. In this scenario, M_L is selected to generate the step s_k . While this incurs a higher computational cost, it leverages the superior reasoning and inherent self-correction capabilities of LRMs (Guo et al., 2025). Specifically, M_L can rectify potential logical drifts accumulated in the historical context c_k , thereby satisfying the quality constraint. We provide a qualitative analysis of this implicit self-correction behavior in Appendix F.

3.5 Efficient Model Switching

A critical requirement for a step-level collaboration system, where every step involves context transitions, is minimizing the model switching overhead. To minimize the system overhead, we leverage prefix caching mechanisms supported by inference engines (Kwon et al., 2023; Zheng et al., 2024). When routing a request between models, the context c_k , which comprises the question and historical steps, is largely resident in the KV cache from previous interactions. Thus, the context processing is reduced to a highly parallelizable prefill phase rather than a serial re-computation. The resulting switching latency is comparable to decoding a few tokens, ensuring that the computational savings from M_S are not negated by routing overheads.

3.6 Hierarchical Acceleration

A distinct advantage of GlimpRouter lies in its step-level granularity. This coarse-grained design is inherently orthogonal to token-level optimizations, allowing our framework to be seamlessly integrated with various low-level acceleration techniques to achieve compound speedups. To maximize system throughput, we implement a hierarchical acceleration strategy in Section 4.3. At the inter-step level, GlimpRouter acts as a global planner, assigning routine logical steps to M_S to bypass the expensive M_L entirely. At the intra-step level, when M_L is invoked, we further accelerate its generation using Speculative Decoding (Leviathan et al., 2023). Specifically, we employ a “Draft-then-Verify” pipeline where the small draft model M_S proposes token sequences that are verified in parallel by M_L .

4 Experiments

4.1 Experimental Setup

Models and Configurations. We conduct comprehensive experiments using the Qwen3 (Yang et al., 2025a) and DeepSeek-R1 (Guo et al., 2025) families. We utilize Qwen3-4B as the small model (SLM) and DeepSeek-R1-Distill-Qwen-32B as the large model (LLM). This setup serves as the default configuration for our results. To verify the scalability of GlimpRouter, we extend our analysis to Qwen3-32B as the LLM in Section 4.2 and DeepSeek-R1-Distill-Qwen-1.5B as the SLM in Appendix B, covering a spectrum of model sizes and architectural pairings.

LLM	Method	AIME24		AIME25		GPQA		LCBv5		LCBv6	
		Acc \uparrow	Lat \downarrow	Acc \uparrow	Lat \downarrow	Acc \uparrow	Lat \downarrow	Acc \uparrow	Lat \downarrow	Acc \uparrow	Lat \downarrow
-	SLM only	48.33	99	45.00	105	61.11	76	47.90	100	43.71	94
Qwen3-32B	LLM only	60.00	220	48.33	231	61.87	194	52.69	249	47.43	241
	Random	56.67	134	47.50	136	61.74	128	51.20	147	44.00	146
	RSD	59.17	167	47.50	173	62.50	165	51.05	209	46.29	208
	SpecCoT	58.33	161	48.33	170	61.62	163	51.05	195	45.71	192
	SpecReason	60.00	160	49.17	162	62.63	181	51.50	213	46.29	214
	GlimpRouter	60.83	145	51.67	147	63.01	142	52.69	162	47.14	165
DeepSeek-32B	LLM only	57.50	197	46.67	220	61.62	176	52.40	219	46.86	214
	Random	56.67	137	48.33	147	62.50	114	52.25	142	47.00	131
	RSD	58.33	167	45.83	171	63.64	146	51.95	180	45.71	178
	SpecCoT	60.00	159	50.00	175	62.75	135	53.29	168	47.14	165
	SpecReason	57.50	158	49.17	169	63.76	213	53.59	185	47.57	189
	GlimpRouter	60.83	143	51.67	163	64.02	129	54.64	160	48.29	160

Table 1: Performance comparison with Qwen3-32B and DeepSeek-R1-Distill-Qwen-32B as LLM, and SLM is fixed as Qwen3-4B. Acc and Lat represent Accuracy (Pass@1, %) and Average Latency (s), respectively. The best performance within each group is highlighted in bold.

Benchmarks. We evaluate our method on a diverse set of complex reasoning tasks. For mathematical reasoning, we employ AIME24 and AIME25, which represent the frontier of mathematical problem-solving. For general reasoning, we use GPQA-Diamond (Rein et al., 2024), an expert-written, challenging multiple-choice dataset in biology, physics, and chemistry. For code generation, we utilize LiveCodeBench (v5 and v6) (Jain et al., 2024), which evaluates the model’s ability to solve competitive programming problems.

Baselines. We compare GlimpRouter against standalone models and state-of-the-art collaborative inference baselines. **Standalone SLM / LLM:** The individual performance of the small and large models, serving as the efficiency and performance boundaries, respectively. **Random:** A method using a random score from 0 to 9 to select the model for inference. **RSD (Liao et al., 2025):** A reward-guided method that employs a trained Process Reward Model (PRM) to evaluate the quality of the step and determine whether to invoke the large model. **SpecCoT (Shi et al., 2025):** A selection-based collaboration method where the small model generates multiple candidate reasoning steps in parallel. The large model acts as a discriminator to select the optimal step. **SpecReason (Pan et al., 2025):** A verification-based method where the small model generates a step, which is then verified by the large model acting as a judge. Upon

rejection, the large model falls back to generation.

Metrics. We assess the efficiency-performance trade-off using two key metrics. **Pass@1 (Acc):** The percentage of problems solved correctly on the first attempt. **Latency (Lat):** The end-to-end wall-clock time (in seconds) per question, serving as the primary indicator of inference efficiency.

Implementation Details. All experiments are conducted on NVIDIA A100-80GB GPUs using the vLLM inference engine (Kwon et al., 2023). We employ sampling with a temperature of 0.6 and a top-p of 0.95. The maximum budget for reasoning tokens is set to 8192. For the hierarchical acceleration strategy, we repurpose the SLM to serve as the token-level drafter for the LLM with a draft length of $n = 3$. To ensure statistical stability, all reported results are averaged over 4 runs. Regarding the entropy threshold τ , since it modulates the trade-off between latency and accuracy, we report the configuration that achieves the optimal balance, maximizing accuracy while securing significant latency reduction, in our main results. Empirically, this corresponds to an *intervention rate* (percentage of steps generated by LLM) of approximately 20%–30%. A detailed sensitivity analysis of τ is provided in Section 4.3.

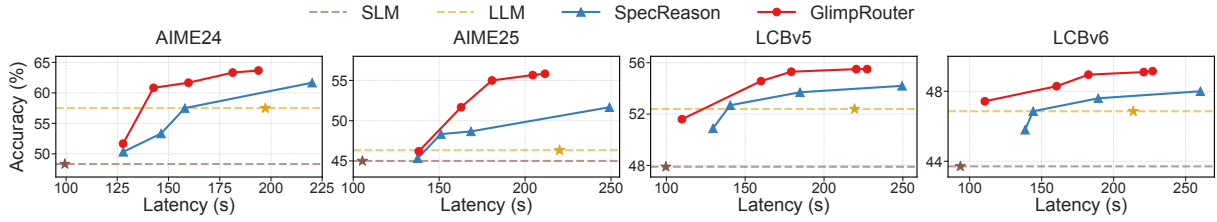


Figure 4: Sensitivity analysis of collaboration threshold on multiple benchmarks. GlimpRouter (red circles) establishes a superior Pareto frontier compared to the SpecReason (blue triangles). Dashed lines represent the performance of standalone models.

4.2 Main Results

Table 1 presents the comparative performance of GlimpRouter against baselines on five complex reasoning benchmarks. GlimpRouter consistently achieves the optimal trade-off between inference latency and reasoning accuracy across different model configurations. In the primary setting using Qwen3-4B as the SLM, our method consistently outperforms the standalone LLM (DeepSeek-R1-Distill-Qwen-32B) in accuracy while delivering substantial speedups. Specifically, GlimpRouter reduces end-to-end latency by **25.2%–27.4%** across all datasets compared to the LLM-only baseline. For instance, on AIME25, we observe a substantial **10.7%** relative improvement in accuracy combined with a **25.9%** reduction in latency.

Compared to existing step-level collaboration methods, GlimpRouter demonstrates a distinct advantage in efficiency. While baselines like RSD, SpecCoT, and SpecReason achieve better accuracy than the SLM-only baseline, they still incur significant latency overheads. For example, SpecReason’s latency on GPQA (213s) exceeds even that of the standalone LLM, DeepSeek-R1-Distill-Qwen-32B (176s), negating the benefits of collaboration. This inefficiency stems from their routing mechanisms: all methods require generating full reasoning steps (and often multiple candidates) or performing post-hoc verification before making a decision. In contrast, GlimpRouter determines the routing path based solely on the *initial token entropy*. This “Probe-then-Dispatch” mechanism minimizes the computational overhead of the decision process.

A counterintuitive yet compelling result is that our collaborative framework outperforms the accuracy of the standalone large model (e.g., 51.67% vs. 46.67% on AIME25). We attribute this to the self-correction capacity of LRMs (Guo et al., 2025; Zeng et al., 2025a). The intervention mechanism in our framework offers more than just text completion. The high entropy at the onset of a step often

serves as a manifestation of latent inconsistencies accumulated in the preceding trajectory. When GlimpRouter detects this uncertainty and invokes the LLM, the large model does not merely continue the sequence; effective intervention allows it to implicitly re-evaluate the context and rectify prior logical drifts, thereby realigning the reasoning path. We provide concrete case studies illustrating this correction behavior in Appendix F.

4.3 Ablation and Analysis

Sensitivity Analysis of Collaboration Threshold. The trade-off between inference efficiency and reasoning capability is closely associated with the proportion of steps routed to the large model, referred to as the *intervention rate*. To investigate how this rate impacts the performance, we vary the entropy threshold τ for GlimpRouter and the verification threshold for SpecReason. By varying these thresholds, we modulate the proportion of steps processed by the large model. Figure 4 illustrates the accuracy-latency trade-off on multiple benchmarks. The curve of GlimpRouter (red) lies strictly above and to the left of the SpecReason baseline (blue), establishing a superior Pareto frontier. Detailed numerical results are provided in Appendix C.

Impact of Metric Choice. To empirically validate the effectiveness of our \mathbf{H}_{init} metric, we conducted an ablation study comparing it against two widely used uncertainty metrics: Step-wise Entropy (\mathbf{H}_{step}) and Step-wise Perplexity ($\mathbf{PPL}_{\text{step}}$). Table 2 presents the comparative results on the AIME25 benchmark, while the extended results are detailed in Appendix D. GlimpRouter outperforms both step-wise variants by a significant margin, achieving a relative accuracy gain of **10.7%** over \mathbf{H}_{step} and **8.8%** over $\mathbf{PPL}_{\text{step}}$. This result corroborates our *Signal Dilution* hypothesis (Section 2.2). Beyond accuracy, \mathbf{H}_{init} offers a fundamental advantage in latency (163s vs. 178s / 181s). Imple-

Variants	Acc \uparrow	Lat \downarrow
GlimpRouter	51.67	163
w/ step-wise entropy	46.67	178
w/ step-wise perplexity	47.50	181

Table 2: Ablation study of different metrics on AIME25. We report accuracy (Acc) and average generation latency (Lat) measured in seconds.

menting routing via step-wise metrics necessitates a ‘‘Generate-then-Measure’’ paradigm, where the initial generation becomes a sunk cost if the step is subsequently routed to M_L . Conversely, \mathbf{H}_{init} operates on a ‘‘Probe-then-Dispatch’’ basis, eliminating the overhead of generating invalid draft steps.

Orthogonal Speedup with Speculative Decoding.

To demonstrate the versatility of our framework, we evaluate the compatibility of GlimpRouter with **Speculative Decoding**, a token-level acceleration technique. Table 3 presents the performance metrics when integrating Speculative Decoding into three different inference paradigms: standalone LLM, SpecReason, and GlimpRouter. As shown in Table 3, integrating Speculative Decoding consistently reduces latency across all methods while maintaining comparable accuracy. Crucially, the combination of **GlimpRouter + Speculative Decoding** achieves the lowest end-to-end latency among all configurations. (1) **Global Synergy**: GlimpRouter optimizes the coarse-grained logical flow by routing easy steps to the small model, reducing the total number of calls to the expensive large model. (2) **Local Synergy**: For the difficult steps that are processed by LLM, Speculative Decoding optimizes the fine-grained token generation, mitigating the high per-token cost of the large model. By attacking the efficiency bottleneck from both step-level routing and token-level execution, our framework achieves a compound speedup that surpasses either technique applied in isolation. Additional results on other benchmarks are available in Appendix E.

5 Related Work

Collaborative Inference. Collaborative inference has emerged as a promising paradigm for optimizing the efficiency-performance trade-off by leveraging models of varying capacities. Existing approaches generally operate across three interaction granularities: query-level routing (Chen et al., 2023b; Ding et al., 2024; Ong et al., 2024),

Method	AIME25		LCBv6	
	Acc \uparrow	Lat \downarrow	Acc \uparrow	Lat \downarrow
LLM only	46.67	220	46.86	214
+ Speculative Decoding	45.83	149	46.29	166
SpecReason	49.17	169	47.57	189
+ Speculative Decoding	49.17	140	47.14	154
GlimpRouter	51.67	163	48.29	160
+ Speculative Decoding	51.67	130	48.00	137

Table 3: Evaluation of orthogonality with Speculative Decoding on AIME25 and LCBv6. We report accuracy (Acc) and average generation latency (Lat) measured in seconds.

step-level collaboration (Shi et al., 2025; Pan et al., 2025; Liao et al., 2025), and token-level speculation (Leviathan et al., 2023; Chen et al., 2023a; Cai et al., 2024; Li et al., 2024). Our GlimpRouter distinguishes itself by proposing a training-free, step-level ‘‘Probe-then-Dispatch’’ mechanism based on initial token entropy, which effectively eliminates the sunk costs associated with post-hoc verification.

Efficient Reasoning. LRMs such as DeepSeek-R1 (Guo et al., 2025) and OpenAI o1/o3 (Jaech et al., 2024; OpenAI, 2025) demonstrate that scaling test-time compute via Chain-of-Thought (Wei et al., 2022) reasoning significantly enhances reasoning capabilities, albeit at the cost of prohibitive latency. To mitigate this bottleneck, recent research explores dynamic offloading computation to smaller models (Chen et al., 2025b; Wang et al., 2025d; Qu et al., 2025). At the query and token levels, methods utilize trained difficulty predictors (Damani et al., 2024) or trained routers (Fu et al., 2025). Most pertinent to logical deduction is step-level collaboration. Existing methods range from training-based reward guidance (RSD) (Liao et al., 2025) to training-free paradigms relying on multi-path selection (SpecCoT) (Shi et al., 2025) or post-hoc verification (SpecReason) (Pan et al., 2025). Unlike these methods constrained by redundant generation or heavy verification overheads, GlimpRouter achieves efficient orchestration via initial token probing.

6 Conclusion

In this work, we addressed the critical latency bottleneck of LRMs by proposing GlimpRouter. This training-free, step-wise collaborative inference framework orchestrates models of varying

capacities. We introduce the “Probe-then-Dispatch” mechanism based on initial token entropy. This approach allows the system to anticipate the step difficulty with negligible overhead. Extensive experiments on multiple benchmarks demonstrate that GlimpRouter establishes a superior Pareto frontier. We hope this work enables efficient reasoning and stimulates research on dynamic computation allocation.

Limitations

Despite the efficiency gains demonstrated by GlimpRouter, we acknowledge several limitations. First, our routing mechanism relies on a static entropy threshold. While initial token entropy provides a more continuous and granular signal compared to discrete metrics like LLM-as-a-Judge, a fixed global threshold may not adapt optimally to the varying difficulty distributions across diverse domains or specific query types. Future work could explore adaptive or instance-aware thresholding mechanisms to further refine the efficiency. Second, the step-level decomposition in our framework relies on explicit structural delimiters (specifically, double newline characters). Although this formatting pattern is prevalent in Large Reasoning Models, such as the DeepSeek-R1 and Qwen families, this heuristic dependence may limit the framework’s direct applicability to models that generate unstructured Chain-of-Thought sequences. Exploring semantic-based segmentation strategies remains a valuable direction for future research.

Acknowledgments

This research is funded by the National Key Research and Development Program of China (Grant No. 2023YFB4503802) and the Natural Science Foundation of Shanghai (Grant No. 25ZR1401175).

References

Tianle Cai, Yuhong Li, Zhengyang Geng, Hongwu Peng, Jason D Lee, Deming Chen, and Tri Dao. 2024. Medusa: Simple llm inference acceleration framework with multiple decoding heads. *arXiv preprint arXiv:2401.10774*.

Charlie Chen, Sebastian Borgeaud, Geoffrey Irving, Jean-Baptiste Lespiau, Laurent Sifre, and John Jumper. 2023a. Accelerating large language model decoding with speculative sampling. *arXiv preprint arXiv:2302.01318*.

Lingjiao Chen, Matei Zaharia, and James Zou. 2023b. Frugalpt: How to use large language models while reducing cost and improving performance. *arXiv preprint arXiv:2305.05176*.

Silin Chen, Shaoxin Lin, Xiaodong Gu, Yuling Shi, Heng Lian, Longfei Yun, Dong Chen, Weiguo Sun, Lin Cao, and Qianxiang Wang. 2025a. Swe-exp: Experience-driven software issue resolution. *arXiv preprint arXiv:2507.23361*.

Yi Chen, JiaHao Zhao, and HaoHao Han. 2025b. A survey on collaborative mechanisms between large and small language models. *arXiv preprint arXiv:2505.07460*.

Yingqian Cui, Pengfei He, Jingying Zeng, Hui Liu, Xianfeng Tang, Zhenwei Dai, Yan Han, Chen Luo, Jing Huang, Zhen Li, and 1 others. 2025. Stepwise perplexity-guided refinement for efficient chain-of-thought reasoning in large language models. *arXiv preprint arXiv:2502.13260*.

Mehul Damani, Idan Shenfeld, Andi Peng, Andreea Bobu, and Jacob Andreas. 2024. Learning how hard to think: Input-adaptive allocation of llm computation. *arXiv preprint arXiv:2410.04707*.

Dujian Ding, Ankur Mallick, Chi Wang, Robert Sim, Subhabrata Mukherjee, Victor Ruhle, Laks VS Lakshmanan, and Ahmed Hassan Awadallah. 2024. Hybrid llm: Cost-efficient and quality-aware query routing. *arXiv preprint arXiv:2404.14618*.

Tianyu Fu, Yi Ge, Yichen You, Enshu Liu, Zhihang Yuan, Guohao Dai, Shengen Yan, Huazhong Yang, and Yu Wang. 2025. R2r: Efficiently navigating divergent reasoning paths with small-large model token routing. *arXiv preprint arXiv:2505.21600*.

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*.

Minghao Hu, Junzhe Wang, Weisen Zhao, Qiang Zeng, and Lannan Luo. 2025. Flowmaltrans: Unsupervised binary code translation for malware detection using flow-adapter architecture. *arXiv preprint arXiv:2508.20212*.

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*.

Naman Jain, King Han, Alex Gu, Wen-Ding Li, Fanjia Yan, Tianjun Zhang, Sida Wang, Armando Solar-Lezama, Koushik Sen, and Ion Stoica. 2024. Livecodebench: Holistic and contamination free evaluation of large language models for code. *arXiv preprint arXiv:2403.07974*.

- 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*.
- Yaniv Leviathan, Matan Kalman, and Yossi Matias. 2023. Fast inference from transformers via speculative decoding. In *International Conference on Machine Learning*, pages 19274–19286. PMLR.
- Han Li, Yuling Shi, Shaoxin Lin, Xiaodong Gu, Heng Lian, Xin Wang, Yantao Jia, Tao Huang, and Qianxiang Wang. 2025. Swe-debate: Competitive multi-agent debate for software issue resolution. *arXiv preprint arXiv:2507.23348*.
- Yuhui Li, Fangyun Wei, Chao Zhang, and Hongyang Zhang. 2024. Eagle: Speculative sampling requires rethinking feature uncertainty. *arXiv preprint arXiv:2401.15077*.
- Baohao Liao, Yuhui Xu, Hanze Dong, Junnan Li, Christof Monz, Silvio Savarese, Doyen Sahoo, and Caiming Xiong. 2025. Reward-guided speculative decoding for efficient llm reasoning. In *Forty-second International Conference on Machine Learning*.
- Chin-Yew Lin and Franz Josef Och. 2004. Orange: a method for evaluating automatic evaluation metrics for machine translation. In *COLING 2004: Proceedings of the 20th International Conference on Computational Linguistics*, pages 501–507.
- Andrey Malinin and Mark Gales. 2020. Uncertainty estimation in autoregressive structured prediction. *arXiv preprint arXiv:2002.07650*.
- Isaac Ong, Amjad Almahairi, Vincent Wu, Wei-Lin Chiang, Tianhao Wu, Joseph E Gonzalez, M Waleed Kadous, and Ion Stoica. 2024. Routellm: Learning to route llms with preference data. *arXiv preprint arXiv:2406.18665*.
- OpenAI. 2025. [Openai o3 and o4-mini system card](#).
- Rui Pan, Yinwei Dai, Zhihao Zhang, Gabriele Oliaro, Zhihao Jia, and Ravi Netravali. 2025. Specreason: Fast and accurate inference-time compute via speculative reasoning. *arXiv preprint arXiv:2504.07891*.
- Xiaoye Qu, Yafu Li, Zhaochen Su, Weigao Sun, Jianhao Yan, Dongrui Liu, Ganqu Cui, Daizong Liu, Shuxian Liang, Junxian He, and 1 others. 2025. A survey of efficient reasoning for large reasoning models: Language, multimodality, and beyond. *arXiv preprint arXiv:2503.21614*.
- Nils Reimers and Iryna Gurevych. 2019. [Sentence-bert: Sentence embeddings using siamese bert-networks](#). In *Proceedings of the 2019 Conference on Empirical Methods in Natural Language Processing*. Association for Computational Linguistics.
- 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*.
- Claude E Shannon. 1948. A mathematical theory of communication. *The Bell system technical journal*, 27(3):379–423.
- Jianshu She, Zhuohao Li, Zheming Huang, Qi Li, Peiran Xu, Haonan Li, and Qirong Ho. 2025. Hawkeye: Model collaboration for efficient reasoning. In *Second Conference on Language Modeling*.
- Junhan Shi, Yijia Zhu, Zhenning Shi, Dan Zhao, Qing Li, and Yong Jiang. 2025. Speccot: Accelerating chain-of-thought reasoning through speculative exploration. In *Findings of the Association for Computational Linguistics: EMNLP 2025*, pages 24405–24415.
- Qiuchen Wang, Ruixue Ding, Zehui Chen, Weiqi Wu, Shihang Wang, Pengjun Xie, and Feng Zhao. 2025a. Vidorag: Visual document retrieval-augmented generation via dynamic iterative reasoning agents. *arXiv preprint arXiv:2502.18017*.
- Qiuchen Wang, Ruixue Ding, Yu Zeng, Zehui Chen, Lin Chen, Shihang Wang, Pengjun Xie, Fei Huang, and Feng Zhao. 2025b. Vrag-rl: Empower vision-perception-based rag for visually rich information understanding via iterative reasoning with reinforcement learning. *arXiv preprint arXiv:2505.22019*.
- Shenzhi Wang, Le Yu, Chang Gao, Chujie Zheng, Shixuan Liu, Rui Lu, Kai Dang, Xionghui Chen, Jianxin Yang, Zhenru Zhang, and 1 others. 2025c. Beyond the 80/20 rule: High-entropy minority tokens drive effective reinforcement learning for llm reasoning. *arXiv preprint arXiv:2506.01939*.
- Xinyuan Wang, Yanchi Liu, Wei Cheng, Xujiang Zhao, Zhengzhang Chen, Wenchao Yu, Yanjie Fu, and Haifeng Chen. 2025d. Mixllm: Dynamic routing in mixed large language models. *arXiv preprint arXiv:2502.18482*.
- 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.
- Zilin Xiao, Jaywon Koo, Siru Ouyang, Jefferson Hernandez, Yu Meng, and Vicente Ordonez. 2025. Prox-thinker: Test-time guidance through small visual reasoners. *arXiv preprint arXiv:2505.24872*.

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*.

Wang Yang, Xiang Yue, Vipin Chaudhary, and Xiaotian Han. 2025b. Speculative thinking: Enhancing small-model reasoning with large model guidance at inference time. *arXiv preprint arXiv:2504.12329*.

Xinwei Yang, Zhaofeng Liu, Chen Huang, Jiashuai Zhang, Tong Zhang, Yifan Zhang, and Wenqiang Lei. 2025c. Elaboration: A comprehensive benchmark on human-llm competitive programming. *arXiv preprint arXiv:2505.16667*.

Qingcheng Zeng, Weihao Xuan, Leyang Cui, and Rob Voigt. 2025a. Thinking out loud: Do reasoning models know when they're right? *arXiv preprint arXiv:2504.06564*.

Wenhao Zeng, Yaoning Wang, Chao Hu, Yuling Shi, Chengcheng Wan, Hongyu Zhang, and Xiaodong Gu. 2025b. Pruning the unsurprising: Efficient code reasoning via first-token surprisal. *arXiv preprint arXiv:2508.05988*.

Jinghan Zhang, Xiting Wang, Fengran Mo, Yeyang Zhou, Wanfu Gao, and Kunpeng Liu. 2025a. Entropy-based exploration conduction for multi-step reasoning. *arXiv preprint arXiv:2503.15848*.

Jusheng Zhang, Kaitong Cai, Yijia Fan, Jian Wang, and Keze Wang. 2025b. **CF-VLM: Counterfactual vision-language fine-tuning**. In *The Thirty-ninth Annual Conference on Neural Information Processing Systems (NeurIPS)*.

Jusheng Zhang, Kaitong Cai, Xiaoyang Guo, Sidi Liu, Qinhan Lv, Ruiqi Chen, Jing Yang, Yijia Fan, Xiaofei Sun, Jian Wang, Chen Ziliang, Liang Lin, and Keze Wang. 2025c. **Mm-cot: A benchmark for probing visual chain-of-thought reasoning in multimodal models**. *Preprint*, arXiv:2512.08228.

Jusheng Zhang, Yijia Fan, Wenjun Lin, Ruiqi Chen, Haoyi Jiang, Wenhao Chai, Jian Wang, and Keze Wang. 2025d. **GAM-agent: Game-theoretic and uncertainty-aware collaboration for complex visual reasoning**. In *The Thirty-ninth Annual Conference on Neural Information Processing Systems (NeurIPS)*.

Lianmin Zheng, Liangsheng Yin, Zhiqiang Xie, Chuyue Livia Sun, Jeff Huang, Cody Hao Yu, Shiyi Cao, Christos Kozyrakis, Ion Stoica, Joseph E Gonzalez, and 1 others. 2024. Sglang: Efficient execution of structured language model programs. *Advances in neural information processing systems*, 37:62557–62583.

A Pseudocode of GlimpRouter

The pseudocode of GlimpRouter is provided in Algorithm 1. Given an input question, the system operates in a step-wise manner. At the onset of each step, instead of blindly generating the full content, GlimpRouter employs M_S to “glimpse” the first token, yielding an entropy \mathbf{H}_{init} that quantizes the difficulty of the upcoming step: If \mathbf{H}_{init} falls below a threshold, the step is deemed routine, and the small model is responsible for continuing the generation. Otherwise, the step signals a cognitive pivot, and the context is handed over to a large model for high-quality generation.

Algorithm 1 Step-wise Collaborative Inference via GlimpRouter

```
1: Inputs: question  $q$ , large model  $M_L$ , small model  $M_S$ , entropy threshold  $\tau$ 
2: Output: a composite reasoning chain  $\hat{\mathcal{T}}$ 
3:  $\hat{\mathcal{T}} \leftarrow \emptyset$  // Initialize the reasoning chain
4:  $k \leftarrow 1$ 
5: while not finished do
6:    $\mathbf{c}_k = \text{CONCAT}(q, \hat{\mathcal{T}})$  // Construct the current context by concatenating question and historical steps
7:    $P_{\text{init}} = M_S(\mathbf{c}_k)$  // Compute the first-token distribution using the small model
8:   Calculate the initial token entropy  $\mathbf{H}_{\text{init}}(s_k)$ 
9:   if  $\mathbf{H}_{\text{init}}(s_k) > \tau$  then
10:     $\hat{s}_k \leftarrow M_L(\mathbf{c}_k)$  // Generate  $\hat{s}_k$  using the large model
11:   else
12:     $\hat{s}_k \leftarrow M_S(\mathbf{c}_k)$  // Generate  $\hat{s}_k$  using the small model
13:   end if
14:    $\hat{\mathcal{T}} = \text{CONCAT}(\hat{\mathcal{T}}, \hat{s}_k)$  // Append the generated step to the reasoning chain
15:    $k \leftarrow k + 1$ 
16: end while
17: return  $\hat{\mathcal{T}}$ 
```

B Scalability across Architectural Pairings

To verify the universality of GlimpRouter beyond specific small models (e.g., Qwen3-4B), we extend our evaluation to a homogeneous model family setting. Specifically, we employ DeepSeek-R1-Distill-Qwen-1.5B as the small model (SLM) and DeepSeek-R1-Distill-Qwen-32B as the large model (LLM).

Table 4 presents the comparative results on the AIME, GPQA-Diamond, and LiveCodeBench benchmarks. Consistent with the observations in our main experiments (Section 4.2), GlimpRouter continues to exhibit a superior efficiency-performance trade-off compared to baselines. Despite the change in the small model,

Method	AIME24		AIME25		GPQA		LCBv5		LCBv6	
	Acc \uparrow	Lat \downarrow	Acc \uparrow	Lat \downarrow	Acc \uparrow	Lat \downarrow	Acc \uparrow	Lat \downarrow	Acc \uparrow	Lat \downarrow
SLM only	25.83	60	26.67	57	15.91	37	15.12	66	20.29	62
LLM only	57.50	197	46.67	220	61.62	176	52.40	219	46.86	214
Random	43.33	127	28.33	135	46.46	88	37.13	142	33.71	138
RSD	46.67	162	35.83	173	49.90	123	41.32	198	37.71	197
SpecCoT	50.00	158	30.83	168	51.01	121	37.96	193	37.03	184
SpecReason	45.83	152	31.67	171	53.41	159	38.08	192	38.26	190
GlimpRouter	51.67	145	39.17	166	54.04	117	45.66	189	43.57	182

Table 4: Performance comparison with DeepSeek-R1-Distill-Qwen-1.5B as SLM and DeepSeek-R1-Distill-Qwen-32B as LLM. Acc and Lat represent Accuracy (Pass@1, %) and Average Latency (s), respectively. The best performance is highlighted in bold.

GlimpRouter maintains its advantage over reactive baselines such as RSD, SpecReason, and SpecCoT. While baselines incur significant latency overhead due to full-step generation and verification, GlimpRouter effectively leverages \mathbf{H}_{init} to route simpler steps to DeepSeek-R1-Distill-Qwen-1.5B preemptively. The success of GlimpRouter in this DeepSeek-R1-Distill-Qwen-1.5B + DeepSeek-R1-Distill-Qwen-32B setting underscores the robustness of the ‘‘Probe-then-Dispatch’’ mechanism. It suggests that the correlation between the initial token’s entropy and step-level difficulty is not an artifact of specific datasets or model families, but rather an intrinsic property of large reasoning models.

C Detailed Sensitivity Analysis of Collaboration Thresholds

In this section, we present the comprehensive numerical results that underpin the sensitivity analysis discussed in Section 4.3. We examine the impact of varying the thresholds on the intervention rate, defined as the percentage of reasoning steps generated by the large model (M_L), and the resulting efficiency-performance trade-offs. We evaluate two distinct routing mechanisms: We sweep the Initial Token Entropy (\mathbf{H}_{init}) threshold $\tau \in \{0.01, 0.1, 0.6, 0.9, 1.8\}$. A lower τ makes the system more cautious, triggering M_L more frequently (Higher Intervention). We vary the verification confidence threshold $\sigma \in \{6, 7, 8, 9\}$ of SpecReason. A stricter threshold leads to higher rejection rates, forcing M_L to regenerate steps more frequently (Higher Intervention).

Table 5 details the performance metrics on the AIME and LiveCodeBench benchmarks. A critical

observation is the scaling behavior of latency with respect to the intervention rate. For GlimpRouter, latency increases linearly and modestly as the intervention rate rises. This is because the overhead of checking the first token is negligible; the latency cost is purely the difference between M_L and M_S generation speeds. In stark contrast, SpecReason exhibits a super-linear latency spike as the intervention rate increases. Taking AIME25 as an example, at an intervention rate of approximately $\sim 50\%$, the latency of SpecReason is significantly higher than that of GlimpRouter. This confirms the ‘‘Sunk Cost’’ hypothesis: SpecReason must waste time generating a full draft by M_S before M_L can intervene, whereas GlimpRouter proactively dispatches the task, eliminating redundant computation. Furthermore, comparing iso-accuracy configurations (e.g., GlimpRouter with $\tau = 0.9$ vs. SpecReason with $\sigma = 9$, both achieving 51.67% on AIME25), our method consistently achieves lower latency. This numerical evidence substantiates the visual trend observed in the Pareto frontier (Figure 4), demonstrating that \mathbf{H}_{init} is more efficient than post-hoc verification.

D Extended Analysis of Metric Choice

In Section 4.3, we demonstrate the superiority of the initial token entropy (\mathbf{H}_{init}) on the AIME25 benchmark. We extend the ablation study to include AIME24, AIME25, LiveCodeBench v5, and LiveCodeBench v6. Table 6 presents the comprehensive results comparing GlimpRouter against Step-wise Entropy (\mathbf{H}_{step}) and Step-wise Perplexity ($\mathbf{PPL}_{\text{step}}$).

Across all four benchmarks, GlimpRouter employing \mathbf{H}_{init} consistently achieves the highest ac-

Method	Thr.	AIME24			AIME25			LCBv5			LCBv6		
		Acc \uparrow	Lat \downarrow	Rate	Acc \uparrow	Lat \downarrow	Rate	Acc \uparrow	Lat \downarrow	Rate	Acc \uparrow	Lat \downarrow	Rate
SLM only	-	48.33	99	0	45.00	105	0	47.90	100	0	43.71	94	0
LLM only	-	57.50	197	100	46.67	220	100	52.40	219	100	46.86	214	100
SpecReason	$\sigma = 6$	50.00	128	2	45.00	138	2	50.89	130	13	45.80	139	15
	$\sigma = 7$	53.33	146	5	48.33	151	4	52.69	141	17	46.86	144	17
	$\sigma = 8$	57.50	158	14	49.17	169	12	53.59	185	32	47.57	189	37
	$\sigma = 9$	61.67	220	48	51.67	249	51	54.20	250	87	48.00	260	87
GlimpRouter	$\tau = 1.8$	51.67	128	2	45.83	138	2	51.60	110	2	47.43	111	2
	$\tau = 0.9$	60.83	143	26	51.67	163	27	54.64	160	37	48.29	160	37
	$\tau = 0.6$	61.67	160	44	55.00	181	45	55.30	179	60	48.95	183	65
	$\tau = 0.1$	63.33	181	76	55.67	205	76	55.50	221	84	49.10	221	88
	$\tau = 0.01$	63.67	194	85	55.83	212	83	55.50	227	89	49.15	227	91

Table 5: Detailed sensitivity analysis of routing thresholds across different benchmarks. **Thr.**: Routing threshold (τ for GlimpRouter, σ for SpecReason). **Acc.**: Pass@1 Accuracy. **Lat.**: Average Latency. **Rate.**: Intervention Rate, the percentage of steps generated by the large model (M_L).

Variants	AIME24		AIME25		LCBv5		LCBv6	
	Acc \uparrow	Lat \downarrow	Acc \uparrow	Lat \downarrow	Acc \uparrow	Lat \downarrow	Acc \uparrow	Lat \downarrow
GlimpRouter	60.83	143	51.67	163	54.64	160	48.29	160
w/ step-wise entropy	58.33	163	46.67	178	53.44	185	47.29	194
w/ step-wise perplexity	59.17	159	47.50	181	53.29	181	47.43	194

Table 6: Comprehensive ablation studies of metrics on four benchmarks. We report accuracy (Acc) and average generation latency (Lat) measured in seconds.

curacy. This universal superiority strongly supports the **Signal Dilution** hypothesis. Step-wise metrics, by averaging uncertainty over the entire sequence length L , allow high-confidence syntactic tokens to obscure the signals from critical reasoning nodes. In contrast, \mathbf{H}_{init} focuses exclusively on the cognitive pivot (the initial token), providing a sharper, undiluted signal that effectively distinguishes between routine steps and complex logical bifurcations. Regarding efficiency, the results highlight a fundamental structural advantage of our method. On every benchmark, the step-wise metrics incur higher latency than GlimpRouter, despite operating at similar intervention rates. This is because step-wise metrics operate on a ‘‘Generate-then-Measure’’ paradigm: the small model must fully generate a step before its difficulty can be assessed. If the step is subsequently rejected and routed to the large model, the initial generation becomes a sunk cost. GlimpRouter, leveraging the ‘‘Probe-then-Dispatch’’ mechanism, assesses difficulty before generation. This proactive routing eliminates invalid drafting,

ensuring that computational resources are allocated with maximal efficiency.

E Extended Analysis of Orthogonality with Speculative Decoding

In Section 4.3, we demonstrated the compatibility of GlimpRouter with Speculative Decoding on specific benchmarks. To verify that this compound speedup is a universal property of our framework, we extend the evaluation to the full suite of benchmarks: AIME24, AIME25, LiveCodeBench v5, and LiveCodeBench v6. Table 7 details the performance metrics for three inference paradigms (Standalone LLM, SpecReason, and GlimpRouter), both with and without Token-level Speculative Decoding.

As evidenced in Table 7, integrating Speculative Decoding yields consistent latency reductions across all methods and benchmarks. For the standalone LLM, latency drops by approximately **22%–32%**. For the collaborative methods (SpecReason and GlimpRouter), the token-level acceler-

Method	AIME24		AIME25		LCBv5		LCBv6	
	Acc \uparrow	Lat \downarrow	Acc \uparrow	Lat \downarrow	Acc \uparrow	Lat \downarrow	Acc \uparrow	Lat \downarrow
LLM only	57.50	197	46.67	220	52.40	219	46.86	214
+ Speculative Decoding	56.67	139	45.83	149	51.80	165	46.29	166
SpecReason	57.50	158	49.17	169	53.59	185	47.57	189
+ Speculative Decoding	57.50	133	49.17	140	53.29	151	47.14	154
GlimpRouter	60.83	143	51.67	163	54.64	160	48.29	160
+ Speculative Decoding	60.00	116	51.67	130	54.34	134	48.00	137

Table 7: Comprehensive evaluation of orthogonal speedup with Speculative Decoding across four benchmarks. We report accuracy (Acc) and average generation latency (Lat) measured in seconds. Rows highlighted in gray indicate the integration of token-level speculative decoding.

ation further compresses the execution time of the large model’s generation phases. Importantly, this speedup is achieved with negligible impact on reasoning accuracy. Crucially, the combination of **GlimpRouter + Speculative Decoding** achieves the lowest end-to-end latency across all experimental configurations. Even when the strong baseline, SpecReason, is accelerated by Speculative Decoding, it remains consistently slower than our enhanced framework (e.g., 140s vs. 130s on AIME25). This gap highlights a fundamental architectural distinction: Token-level speculation can only accelerate the generation process; it cannot recover the sunk costs of invalid steps discarded during reactive verification. SpecReason still pays the penalty for generating and rejecting full steps. In contrast, GlimpRouter leverages \mathbf{H}_{init} to proactively route steps, minimizing the number of tokens that need to be generated in the first place. When combined with Speculative Decoding, which accelerates the tokens that must be generated by M_L , we achieve a multiplicative effect on efficiency. This confirms that GlimpRouter serves as an ideal **Global Planner**, perfectly complementing **Local Executors** like Speculative Decoding to establish a new state-of-the-art Pareto frontier.

F Case Studies

F.1 Case Study 1: Correlation between Initial Token Entropy (\mathbf{H}_{init}) and Reasoning Difficulty

We provide a concrete example to show the correlation between \mathbf{H}_{init} and reasoning difficulty. The case visualizes how GlimpRouter dynamically allocates computational resources based on the cognitive load of each step. The case below illustrates

a representative trajectory from the AIME benchmark. We observe a clear dichotomy in \mathbf{H}_{init} values corresponding to the nature of the reasoning steps:

- **Low \mathbf{H}_{init} (Routine Execution):** Steps 1-2 and 5-12 exhibit extremely low \mathbf{H}_{init} values (e.g., 0.0046, 0.0008). These steps involve factual restatements or mechanical algorithmic operations (e.g., recursive division by 2). The small model (M_S) demonstrates high confidence and generates these tokens accurately, validating our “Delegate” strategy.
- **High \mathbf{H}_{init} (Cognitive Bifurcation):** A critical spike in entropy occurs at the onset of Step 3 (\mathbf{H}_{init} : 1.8985). Here, the reasoning process transitions from problem setup to solution planning. The small model’s draft (starting with “Maybe”) reflects uncertainty in selecting the counting method.

A Case of Correlation between Initial Token Entropy (\mathbf{H}_{init}) and Reasoning Difficulty

Question:

Alice chooses a set A of positive integers. Bob lists all finite nonempty sets B whose maximum is in A . Bob’s list has 2024 sets. Find the sum of the elements of A .

Reasoning Trace via GlimpRouter:

Step 1 [SLM] (\mathbf{H}_{init} : 0.2369)

First, Alice chooses a set A , and Bob lists all nonempty sets B whose maximum is in A ; there are 2024 such sets.

Step 2 [SLM] (\mathbf{H}_{init} : 0.3485)

If $A = \{2\}$, then B can be $\{2\}$ or $\{1, 2\}$, giving $2^{2-1} = 2$ sets. Generally, the number of sets corresponds to a sum of powers of 2.

Step 3 (Critical Pivot)

Content (M_S): “Maybe I can try to find...” (H_{init} : 1.8985 > τ)

↔ **Action:** **Intervene** (M_L)

Content (M_L): The most efficient way is to divide 2024 by 2 repeatedly to find its binary representation directly.

Step 4 [SLM] (H_{init} : 0.0046)

Compute $2024/2 = 1012$, remainder 0.

Step 5 [SLM] (H_{init} : 0.0008)

Then $1012/2 = 506$, remainder 0.

Step 6 [SLM] (H_{init} : 0.0102)

Next $506/2 = 253$, remainder 0.

Steps 7-10 [SLM] (H_{init} values: 0.0055, 0.0032, 0.0041)

[...Routine division steps: $253 \rightarrow 126 \rightarrow 63 \rightarrow 31$...]

Step 11 [SLM] (H_{init} : 0.1205)

So the binary representation is 11111101000_2 . This corresponds to powers $2^{10}, \dots, 2^3$. Since the term is 2^{a-1} , the elements are $a \in \{11, 10, 9, 8, 7, 6, 4\}$.

Step 12 [SLM] (H_{init} : 0.0020)

Finally, the sum is $11 + 10 + 9 + 8 + 7 + 6 + 4 = \boxed{55}$.

F.2 Case Study 2: Self-Correction Mechanism in LRMs

The case below demonstrates the system’s ability to detect and correct potential logical errors during the reasoning process. In combinatorics problems involving grid paths, an accurate definition of geometric properties is a prerequisite for calculation.

- **Logical Instability:** In Step 3, the small model establishes the relationship between “direction changes” and “path segments”. It produces content containing a factual error: equating four direction changes to *four* segments (the correct number is $4 + 1 = 5$).
- **Correction via Intervention:** Upon triggering the high-entropy threshold in Step 4, GlimpRouter invokes the large model (M_L). Crucially, M_L does not merely complete the SLM’s erroneous sentence. Instead, it re-evaluates the context and generates a correct premise: “that means it has five straight segments.”
- **Outcome:** By correcting the hallucinated step and injecting the correct geometric reasoning, M_L steers the trajectory back to a valid path. Subsequently, M_S successfully performs the routine combinatorial calculations (Steps 5-6) based on the corrected logic.

A Case of the Self-correction Mechanism in LRMs

Question:

Consider the paths of length 16 that follow the lines from the lower left corner to the upper right corner on an 8×8 grid. Find the number of such paths that change direction exactly four times.

Reasoning Trace via GlimpRouter:

Step 1 [SLM] (H_{init} : 0.2457)

So, the problem is about paths on a grid where you start at (0,0) and end at (8,8), moving only right (R) or up (U), each step being a unit move.

Step 2 [SLM] (H_{init} : 0.3914)

Let me recall that in such grid paths, a path that goes from (0,0) to (8,8) can be represented as a sequence of right (R) and up (U) moves.

Step 3 [SLM] (Logical Divergence)

Suppose a path starts... four direction changes mean the path has *four segments*.

Step 4 [LLM] (Logical Correction)

Wait, if a path changes direction exactly four times, that means it has *five straight segments*. Because each change of direction adds a new segment. For example, starting with R gives R-U-R-U-R.

Step 5 [SLM] (H_{init} : 0.4949)

Let me recall that the number of compositions of n into k positive integers is $C(n-1, k-1)$. Because it’s equivalent to placing $k-1$ dividers in the $n-1$ gaps between n items.

Step 6 [SLM] (H_{init} : 0.2771)

Therefore, the total number is $2 \times \left[\binom{7}{2} \times \binom{7}{1} \right] = 2 \times 21 \times 7 = \boxed{294}$.