

Solve-Detect-Verify : Inference-Time Scaling with Flexible Generative Verifier

Anonymous ACL submission

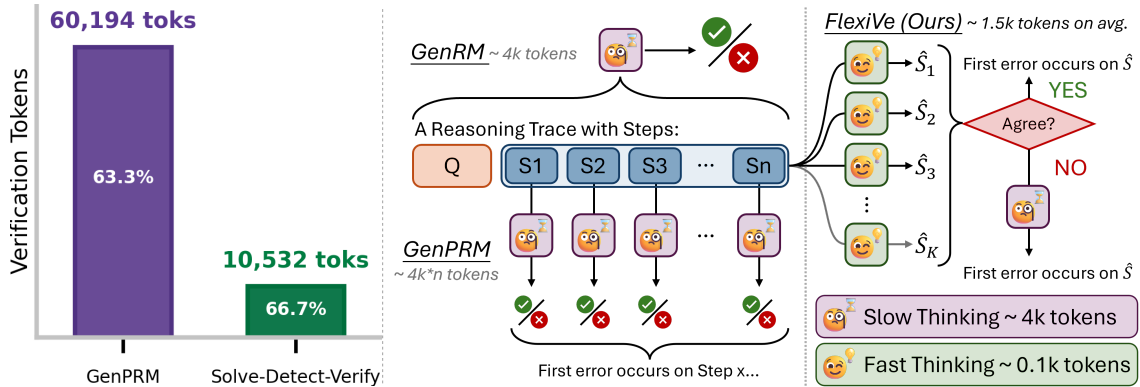


Figure 1: The *Solve-Detect-Verify* (SDV) pipeline transforms linguistic signals into efficiency. **Left:** On AIME 2024, SDV achieves 83.3% accuracy (vs. 63.3% for GenPRM) while using **6x fewer** verification tokens by pruning redundant reasoning. **Right:** The pipeline is powered by *FlexiVe*, a unified verifier. Unlike process-based verifiers that incur per-step overhead, *FlexiVe* analyzes traces **holistically**. It employs a “pragmatic” consensus strategy: parallel “Fast Thinking” checks ($\sim 0.1k$ tokens) provide an initial semantic intuition, escalating to deliberative “Slow Thinking” ($\sim 4k$ tokens) only when the model exhibits verbalized uncertainty.

Abstract

Complex reasoning with Large Language Models (LLMs) demands a careful balance between accuracy and computational cost. Verification is crucial for reliability but faces trade-off: robust process-based verifiers are computationally prohibitive, while fast verifiers lack precision. We introduce *FlexiVe*, a unified generative verifier designed to navigate this trade-off by dynamically allocating compute between “rapid fast thinking” and “deliberative slow thinking.” A key innovation is our training strategy: we use Group Relative Policy Optimization (GRPO) to specifically enhance the reliability of the fast mode. This targeted training generalizes effectively, elevating the slow mode to state-of-the-art open-source performance. To deploy *FlexiVe*, we propose the *Solve-Detect-Verify* (SDV) pipeline. Moving beyond static Best-of-N ranking, SDV employs an iterative refinement process that utilizes likelihood-based probing to detect solution completion, curtailing “overthinking”, and leverages *FlexiVe*’s feedback for targeted correction. *FlexiVe* establishes a new open-source state-of-the-art on ProcessBench, outperform-

ing GenPRM-32B while requiring $\sim 2.3x$ fewer TFLOPS and 15x less training data. On AIME 2024, the full SDV pipeline achieves 83.3% accuracy, surpassing strong baselines while using significantly fewer tokens.

1 Introduction

Recent advances in Large Language Models (LLMs) have enhanced capabilities in complex reasoning through step-by-step trace generation (Wei et al., 2022; Kojima et al., 2022). However, this shift toward deeper “System 2” processes (Shao et al., 2024a) introduces a fundamental trade-off between accuracy and efficiency.

This challenge is twofold. First, models often exhibit “overthinking” (Chen et al., 2024), generating redundant self-correction steps and hesitation markers even after implicitly reaching a correct solution. Second, ensuring reliability requires verification, which adds significant overhead. Sophisticated Generative Reward Models (GenRMs) (Liu et al., 2025) are often computationally prohibitive, while highly token-efficient mechanisms like “NoThinking” (Ma et al., 2025) suffer from severe drops in

050 precision when adapted for verification (see Fig- 099
051 ure 2). 100

052 To address this methodological gap, we require 101
053 a verifier that adapts computational effort to task 102
054 complexity and an inference pipeline that stream- 103
055 lines reasoning. We introduce *FlexiVe*, a uni- 104
056 fied generative verifier, and the *Solve-Detect-Verify* 105
057 (SDV) pipeline. Our contributions are: 106

- 058 • ***FlexiVe* : A Flexible, RL-Trained Genera- 107
059 tive Verifier.** We introduce a unified model 108
060 operating across a cost-performance spec- 109
061 trum, featuring a fast thinking mode, a de- 110
062 liberative slow thinking mode, and a dynamic 111
063 consensus-based “flexible” mode. Uniquely, 112
064 we employ Group Relative Policy Optimiza- 113
065 tion (GRPO) (Shao et al., 2024b) to specifi- 114
066 cally train the fast mode. We find this targeted 115
067 RL training fixes the low precision of fast ver- 116
068 ifiers and generalizes to elevate the slow mode 117
069 to state-of-the-art performance. 118
- 070 • ***Solve-Detect-Verify* (SDV) Pipeline.** We 119
071 propose an intelligent inference framework 120
072 that moves beyond standard Best-of-N (BoN) 121
073 paradigms. SDV employs an iterative refine- 122
074 ment process featuring a lightweight “De- 123
075 tect” module that uses **likelihood-based prob-** 124
076 **ing** (Kadavath et al., 2022; Lin et al., 2022) 125
077 to identify solution completion points and curtail 126
078 overthinking. It then triggers *FlexiVe* to pro- 127
079 vide targeted, generative feedback, guiding 128
080 the solver to refine its response. 129
- 081 • **State-of-the-Art Efficiency and Accuracy.** 130
082 *FlexiVe* sets a new open-source SOTA on Pro- 131
083 cessBench, outperforming larger models like 132
084 GenPRM-32B while requiring $\sim 2.3x$ fewer 133
085 TFLOPS and **15x less training data**. On the 134
086 AIME 2024 benchmark, the SDV pipeline 135
087 achieves 83.3% accuracy, outperforming com- 136
088 parable GenPRM BoN setups while using 137
089 only **1/6th** of the computational tokens. 138

090 2 Related Work

091 **Inference-Time Scaling Strategies** Inference-time 130
092 scaling strategies increase test-time compute to im- 131
093 prove reasoning accuracy (Welleck et al., 2024; 132
094 Wang et al., 2025), using methods from self- 133
095 consistency (Wang et al., 2023), verifier ranked 134
096 Best-of-N (BoN) (Ichihara et al., 2025), to tree- 135
097 based searches (Yao et al., 2023). While effec- 136
098 tive, these strategies are computationally intensive,

spurring work on optimized decoding (Sun et al., 2024) and compute trade-offs (Wu et al., 2025a). As scaling generations alone is insufficient (Chen et al., 2025) and verifier-guided search has known flaws (Wu et al., 2025b; Zhao et al., 2025a), intelligent frameworks like Solve-Detect-Verify (SDV) are needed. While building on established iterative refinement concepts (Madaan et al., 2023; Xie et al., 2023; Akyurek et al., 2023), SDV uniquely prioritizes efficiency through active detection and adaptive verification to avoid the computational redundancy (“overthinking”) typical of brute-force methods (Chen et al., 2024).

Verification Paradigms Verification, while crucial, adds computational cost. Generative and process-based verifiers like GenRMs and PRMs (Lightman et al., 2023; Liu et al., 2025; Zhang et al., 2025) offer detailed feedback but can be demanding (Singhi et al., 2025). Recent work reduces annotation reliance via bootstrapping (Zelikman et al., 2022) or label-free methods like Math-Shepherd (Wang et al., 2024b). Hybrid models like GenPRM (Zhao et al., 2025b) integrate code execution within a process-based framework, motivating its use as a key baseline. Alternative paradigms like code-based self-verification (Zhou et al., 2024a; Wang et al., 2024a) and autoformalization (Zhou et al., 2024b) use code for precision but may lack general applicability. In contrast, *FlexiVe* performs efficient, holistic trace analysis with dynamic budget allocation, targeting broader use cases.

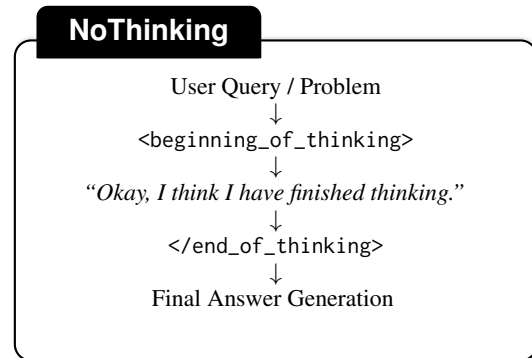


Figure 3: The *NoThinking* mechanism bypasses explicit thought generation, using a template to fill the thinking phase.

Adaptive Computation and Our Novelty Inspired by dual-process theory (Kahneman, 2011; Li et al., 2025), adaptive computation balances reasoning and efficiency (Graves, 2016). However, extreme efficiency methods like “NoThinking” (Ma et al., 2025) in Figure 3, when applied to verification, can yield low precision (Figure 2). The most

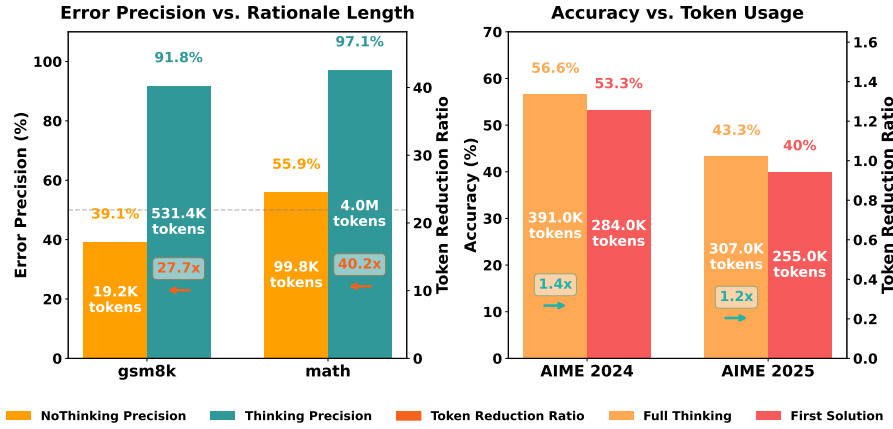


Figure 2: Empirical motivation for efficient verification and generation strategies. **(Left)** Comparison of error precision and token usage between *NoThinking* and *Thinking* verification on GSM8K and Math (ProcessBench). While *NoThinking* significantly reduces tokens, its error precision is substantially lower, suggesting high false positive rate. **(Right)** Accuracy and token usage comparison between generating a full solution (*Full Thinking*) and halting generation early upon detecting a complete intermediate solution (*First Solution*) on AIME 2024 and AIME 2025. Early detection offers significant token reduction with comparable accuracy.

related work, DyVe (Zhong et al., 2025), also uses “fast” and “slow” verification modes. However, its per-step approach incurs accumulating overhead. FlexiVe differs critically by (1) performing holistic, **consensus-based verification** on the entire reasoning trace to avoid iterative costs, and (2) optimizing its “fast” mode for reliable diagnosis via Reinforcement Learning (GRPO) (Shao et al., 2024a) for a more robust efficiency-accuracy balance.

3 Method

3.1 Problem Formulation

System Components Our inference-time scaling framework uses two primary Large Language Model (LLM) components: a solver LLM and *FlexiVe*, our specialized generative verifier. Both are reasoning-capable models. The solver, an off-the-shelf LLM, generates initial candidate solutions. *FlexiVe* is specifically trained for verification, detailed in Section 3.3.

Reasoning Trace Segmentation A reasoning trace S_{trace} is parsed into an ordered sequence of N_s steps, $S_{trace} = (step_1, \dots, step_{N_s})$. Each $step_i$ is a contiguous text segment delineated by predefined hesitation keywords (e.g., ‘Wait, double-check’, ‘Alternatively’, ‘Hmm’, ‘Let me check’ listed in Appendix A.1.3 Figure 8). This segmented trace forms the input for verification.

3.2 Verifier Architectures and Operation

Verifier architectures differ significantly in efficiency (Figure 1). **Process-Based Verifiers** (e.g., GenPRM) verify sequentially, incurring computational overhead by iteratively re-processing context

for every step. In contrast, **Holistic Verifiers** (e.g., *FlexiVe*) evaluate the trace in a single pass. Unlike standard GenRMs that use fixed budgets, *FlexiVe* dynamically modulates effort (Section 3.3).

Using a critic template (Zheng et al., 2024a), *FlexiVe* outputs $V_{out} = (F, idx_{pred})$, comprising a textual analysis F and the first error index idx_{pred} (or -1). We employ this generative approach (F) rather than scalar scores for two reasons: articulating reasoning improves verification accuracy by enforcing dependency checking (Liu et al., 2025; Zhang et al., 2025), and the resulting feedback serves as actionable diagnostics to guide solver refinement in our *Solve-Detect-Verify* pipeline.

3.3 FlexiVe : A Unified Generative Verifier

FlexiVe is a unified generative verifier designed to operate across the entire spectrum of cost-performance trade-offs by leveraging a single model with three distinct inference-time modes. At one extreme, its **Fast Thinking (NoThinking)** mode, inspired by the “NoThinking” mechanism (Ma et al., 2025), this mode prioritizes extreme efficiency. It utilizes a specific template (see Figure 3) to bypass explicit thought generation, filling the thinking phase with a placeholder before directly outputting the verification result. This approach results in responses that are approximately 40× shorter than the “Slow Thinking” mode (see Figure 2), enabling high-throughput, parallel sampling with minimal latency. At the other, the **Slow Thinking (Think)** mode generates a full, detailed reasoning trace to maximize verification accuracy. Our novel **Flexible Allocation (Flex)** mode dynam-

ically bridges these approaches, adaptively switching between Fast and Slow Thinking based on perceived task difficulty to optimally balance accuracy and cost.

Reinforcement Training for Reliable Fast Thinking A critical challenge for efficient verifiers is their low precision (Figure 2) under NoThinking mode (Figure 3). We address this through a targeted Reinforcement Learning strategy using Group Relative Policy Optimization (GRPO) (Shao et al., 2024a,b). Our goal is to maximize the reliability of the “fast thinking” mode while maintaining its efficiency.

To achieve this, we train *FlexiVe* specifically in the “fast thinking” configuration (activating the NoThinking template during training). The model predicts the index of the first error (e_{gt}) or -1 if correct. GRPO optimizes the policy by maximizing a composite reward $R_i = R_{\text{correct}} + R_{\text{length}}$.

The correctness reward R_{correct} is defined by:

$$R_{\text{correct}}(e_{\text{pred}}, e_{\text{gt}}) = \begin{cases} 1.0 & \text{if } e_{\text{pred}} = e_{\text{gt}} \\ 0.0 & \text{otherwise} \end{cases}. \quad (1)$$

To prevent the model from “reward hacking” with verbose outputs and ensure efficiency within the “fast thinking” constraint, we apply a length-based regularization term, R_{length} , proportional to the length L of the generated response:

$$R_{\text{length}}(L) = -\lambda \cdot L. \quad (2)$$

The hyperparameter λ (empirically set to 0.1) is crucial to ensure the “fast thinking” mode remains token-efficient, preventing the RL policy from converging to verbose outputs that violate the efficiency goal. Training involves sampling G outputs per prompt and calculating advantages relative to the group’s average (Shao et al., 2024a). A key finding, explored in Section 4.3, is that this targeted RL training not only substantially improves “fast thinking” precision but also generalizes remarkably, enhancing the accuracy of the “slow thinking” mode.

Flexible Allocation of Verification Budget (Flex@k) The dynamic “Flexible” mode utilizes a two-stage verification process to tailor computational effort to the difficulty of the trace.

At inference time, the process begins with an efficient, parallelizable probing stage. *FlexiVe* performs k independent “Fast Thinking” runs, utilizing the token-efficient NoThinking template, on the entire reasoning trace. The decision to escalate is determined dynamically by the consensus among these runs. Each run produces an outcome consisting of the predicted error index (or -1 if correct).

We measure consensus by the agreement ratio:

$$R_{\text{agreement}} = \frac{\max_i a_i}{k}, \quad (3)$$

where a_i is the count of the most frequent outcome.

If the consensus is high ($R_{\text{agreement}} \geq \tau$), it signals a straightforward case, and the “Fast Thinking” result, V_{fast} , is accepted efficiently. If consensus is low, it indicates ambiguity, and the framework escalates to the second stage: performing $\max(1, \lceil k/8 \rceil)$ resource-intensive “Slow Thinking” runs to produce a robust final outcome, V_{slow} . Methodologically, “Slow Thinking” re-processes the problem and solver responses without appending the template shown in Figure 3.

The overall verification result V is:

$$V = \begin{cases} V_{\text{fast}}, & \text{if } R_{\text{agreement}} \geq \tau, \\ V_{\text{slow}}, & \text{otherwise.} \end{cases} \quad (4)$$

The consensus threshold τ and sample count k are critical hyperparameters. We selected these based on a detailed sensitivity analysis (Appendix A.4.1) and Pareto frontier analysis (Section 4.3). We identify $\tau = 0.8$ and $k = 8$ (for Flex@8) as the optimal trade-off point (the “knee” of the performance curve), maximizing accuracy gains while minimizing computational overhead.

3.4 Solve-Detect-Verify

Solve-Detect-Verify is a framework designed to enhance LLM reasoning accuracy and efficiency through iterative refinement, moving beyond static Best-of-N ranking. It integrates three stages: Solve, Detect, and Verify/Refine. The full pipeline implementation is detailed in Appendix A.1.3.

Solve The ‘Solve’ stage initiates the process, wherein the solver LLM is tasked with generating an initial, step-by-step candidate solution (S_1) to a given problem. This stage forms the foundational attempt at problem-solving, producing a complete reasoning trace and a final answer for subsequent evaluation.

Detect The ‘Detect’ module, detailed in Algorithm 1 (see Appendix A.1.3), continuously monitors the output for hesitation keywords (Appendix Figure 8). Upon detection, generation pauses, and the LLM assesses solution completeness via a log-probability check ($\log p(\text{Yes})$ vs. $\log p(\text{No})$). This check efficiently reuses over 90% of the generation prefix (KV cache), minimizing overhead. If deemed complete, the pipeline advances; otherwise, generation resumes. This curtails “overthinking” and enables early verification.

Verify and Refine The candidate solution S_1 is assessed by *FlexiVe*. If correct, it is accepted. Oth-

erwise, diagnostic feedback (F_1) guides the solver to generate a new solution, S_2 . This feedback loop acts as an efficient context engineering strategy to refine the model’s reasoning path.

Iterative Refinement and Scalability The ‘Verify and Refine’ stages can be iterated to progressively improve the solution, as outlined in the full pipeline control flow in Algorithm 2 (Appendix A.1.3). The number of iterations, T , is a tunable parameter that creates a trade-off between computational cost and final accuracy. As shown in Figure 5 (top-right), each iteration yields monotonic accuracy gains, allowing the framework’s computational depth to be scaled according to specific performance and budget requirements.

4 Experiments

Our experiments address four primary questions: (1) How accurate and sample-efficient is *FlexiVe* compared to state-of-the-art (SOTA) verifiers? (2) Does *FlexiVe* offer a superior accuracy-efficiency trade-off (Pareto frontier) when measured in TFLOPS? (3) Does the *Solve-Detect-Verify* outperform standard inference-time scaling strategies like Best-of-N (BoN) ranking? (4) Are all components of the pipeline necessary and robust?

4.1 Experimental Setup

For detailed experimental configurations, including hyperparameter settings and full dataset statistics, please refer to Appendix A.1.1.

Evaluation Tasks and Datasets We assess *FlexiVe*’s step-level verification capability (F1 score) on the ProcessBench benchmark (Zheng et al., 2024a) (GSM8K, MATH, OlympiadBench, OmniMATH splits). For the full *Solve-Detect-Verify*, we evaluate end-to-end task accuracy and efficiency on challenging mathematical datasets: AIME 2024 and 2025 (Hugging Face, 2024; OpenCompass, 2025). Efficiency is measured using total generated tokens and, crucially, estimated TFLOPS¹ to account for architectural differences across baselines.

Baselines On ProcessBench, *FlexiVe* is compared against established Process Reward Models (PRMs) (Zheng et al., 2024a), including GenPRM (7B and 32B) (Zhao et al., 2025b). For evaluating the *Solve-Detect-Verify*, DeepSeek-R1 14B (DS14B) and 32B models (Shao et al., 2024a) serve as the base ‘worker’ LLMs. Performance is benchmarked against direct output, Self-Consistency (Majority Voting) (Wang et al., 2023), and BoN

ranking using external verifiers.

FlexiVe Training *FlexiVe* (14B) is initialized from DeepSeek-R1-Distill-Qwen-14B and trained using GRPO on the BIG-Bench Mistake dataset (Tyen et al., 2024). Notably, training utilized only 1,526 samples. The training focused specifically on the ‘fast mode’ (NoThinking mechanism activated) to optimize rapid, reliable error detection.

FlexiVe Configurations We evaluate *FlexiVe* in three distinct configurations, where k denotes the number of verification samples: (1) **Think@k**: Fixed “slow” budget. Performs k independent “slow thinking” (deliberative) runs with a majority vote. (2) **NoThinking@k**: Fixed “fast” budget. Performs k independent, token-efficient “fast thinking” runs with a majority vote. (3) **Flex@k**: Adaptive budget. Begins with k “fast thinking” runs and escalates to “slow thinking” only if initial consensus is below threshold τ . Provides a dynamic trade-off.

4.2 FlexiVe: A Unified, State-of-the-Art Verifier

We first evaluate the verification capabilities of the *FlexiVe* model on the ProcessBench benchmark and analyze the effectiveness of its novel RL training strategy.

State-of-the-Art Open-Source Verification Accuracy Table 1 details the F1 scores across various mathematical reasoning datasets. In the “High Compute” setting, *FlexiVe* (Think@64) establishes a new state-of-the-art for open-source models, achieving an average F1 score of 86.3%. This notably outperforms the compute-intensive GenPRM-32B (Maj@8) augmented with code execution (82.6% Avg F1). In the “Moderate Compute” setting, the adaptive *FlexiVe* (Flex@128) achieves a strong average F1 score of 80.8%, surpassing GenPRM-32B (Maj@8) without code execution (79.3% Avg F1).

4.3 Pareto Frontier Analysis: Accuracy and Efficiency

We analyze the accuracy-efficiency trade-off of *FlexiVe* against the process-based GenPRM (see Appendix A.3 for detailed data and extended analysis). Unlike GenPRM, which incurs non-linear costs by iteratively re-processing expanding contexts, *FlexiVe* operates holistically. This architectural advantage allows *FlexiVe* (**Think@k**) to define a superior Pareto frontier; specifically, Think@4 outperforms GenPRM-32B (~87% vs. ~84% F1) while consuming less than half the

¹Calculated as (input + output tokens) \times model parameters, normalized.

Model	# Samples	GSM8K	MATH	Olympiad Bench	Omni-MATH	Avg.
<i>Proprietary Models</i>						
GPT-4o-0806	unk	79.2	63.6	51.4	53.5	61.9
o1-mini	unk	93.2	88.9	87.2	82.4	87.9
<i>Open Source Models (7-8B)</i>						
Qwen2.5-Math-PRM-7B	~344K	82.4	77.6	67.5	66.3	73.5
RetrievalPRM-7B	404K	74.6	71.1	60.2	57.3	65.8
Universal-PRM-7B	unk	85.8	77.7	67.6	66.4	74.3
Direct Generative PRM-7B	23K	63.9	65.8	54.5	55.9	60.0
GenPRM-7B w/ Code Exec (Pass@1)	23K	78.7	80.3	72.2	69.8	75.2
GenPRM-7B w/ Code Exec (Maj@8)	23K	81.0	85.7	78.4	76.8	80.5
<i>Open Source Models (14-32B) w/ Moderate Compute</i>						
Dyve-14B	117K	68.5	58.3	49.0	47.2	55.8
GenPRM-32B w/o Code Exec (Maj@8)	23K	78.8	85.1	78.7	74.9	79.3
<i>FlexiVe</i> (Flex@32)	1526	<u>82.8</u>	83.3	<u>79.2</u>	73.4	<u>79.7</u>
<i>FlexiVe</i> (Flex@128)	1526	83.0	<u>85.0</u>	80.0	75.2	80.8
<i>Open Source Models (14-32B) w/ High Compute</i>						
GenPRM-32B (Pass@1) w/ Code Exec	23K	83.1	81.7	72.8	72.8	77.6
GenPRM-32B (Maj@8) w/ Code Exec	23K	85.1	86.3	78.9	80.1	82.6
<i>FlexiVe</i> (Think@64)	1526	88.1	90.1	86.7	80.4	86.3

Table 1: ProcessBench results reported with F1 scores. Results for *FlexiVe* are highlighted. **bold** indicates the best in the sub category. All *FlexiVe* variants are trained on only 1526 samples.

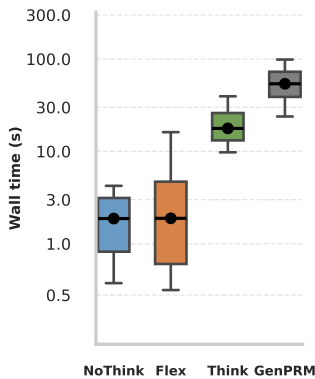


Figure 4: **Wall-clock Time Comparison.** Comparison of inference latency on ProcessBench (MATH). *FlexiVe*’s parallelizable **Flex** mode achieves a median latency of $\sim 2s$, matching the **NoThink** baseline and operating roughly $25\times$ faster than the sequential **GenPRM**.

TFLOPS (~ 12 vs. ~ 29). Crucially, theoretical TFLOPS metrics understate the practical latency benefits revealed in Figure 4. By executing low-cost “Fast Thinking” runs in parallel, the **Flex** mode achieves a median wall time of $\sim 2s$ —matching the **NoThink** baseline and approximately $25\times$ faster than the sequential **GenPRM** ($\sim 50s$). Based on this superior balance of TFLOPS efficiency and wall-clock latency, we select **Flex@8** as the preferred configuration for all subsequent scaling ex-

periments.

4.3.1 Sample Efficiency and Ablation

We conduct an ablation on ProcessBench (Table 2) to isolate the impact of our training strategy and base model.

RL Efficiency and Model Selection Using the identical DeepSeek-R1-14B base and small dataset (1,526 samples), standard Discriminative PRM (12.9% Avg F1) and Supervised Fine-Tuning (49.0% Avg F1) failed to generalize. Our RL strategy not only outperformed these baselines but also surpassed an SFT model trained on $6.5\times$ more synthetic data, confirming that gains stem from the GRPO objective rather than data scale. Furthermore, while DeepSeek-R1-14B provides a strong foundation (70.8% F1) compared to Llama-3 or QwQ, *FlexiVe* significantly elevates this performance to 75.6%, proving our alignment adds value beyond the base model’s capabilities.

Generalization Across Inference Modes Crucially, our RL training instills robustness across computational budgets. The unaligned base model collapses under the token-efficient “NoThink” template (37.1% F1), confirming that standard models lack inherent efficient verification skills. *FlexiVe* boosts this “fast-thinking” performance to 58.6% ($\sim 58\%$ relative improvement). This reliability pow-

Model / Configuration	Training Method	# Samples	GSM8K	MATH	Olym.	Omni.	Avg.
Base Model Selection (Think@1)							
Meta-Llama-3-8B-Instruct	None (Base)	-	26.8	13.2	12.3	13.2	16.4
QwQ-32B-Preview	None (Base)	-	75.5	59.2	35.7	35.3	51.4
DeepSeek-R1-14B	None (Base)	-	77.6	76.2	65.6	64.0	70.8
FlexiVe (Think@1)	RL (Ours)	1.5K	82.6	80.3	73.1	66.3	75.6
Training Strategy Ablation (Base: DeepSeek-R1-14B)							
Discriminative PRM	Math-Shepherd	445K	15.8	15.9	8.3	11.9	12.9
Discriminative PRM	SFT	1.5K	66.3	56.0	36.1	37.7	49.0
Generative Verifier	SFT	10K	71.9	69.0	59.7	47.9	62.1
FlexiVe (NoThink)	RL (Ours)	1.5K	82.6	80.3	73.1	66.3	75.6
RL Impact Across Inference Modes (Base: DeepSeek-R1-14B)							
Base Model (Flex@4)	None	-	57.9	62.8	59.6	59.5	60.0
FlexiVe (Flex@4)	RL (Ours)	1.5K	78.4	77.7	72.4	67.3	74.0
Base Model (NoThink@4)	None	-	39.5	36.0	33.9	39.0	37.1
FlexiVe (NoThink@4)	RL (Ours)	1.5K	66.8	61.3	53.8	52.5	58.6

Table 2: **Ablation Study on Base Models and Training Strategies.** **Top:** Comparison of base models (Think@1). **Middle:** Comparison of training methods on identical data (1.5K samples). **Bottom:** Comparison of RL impact across different inference modes. The base model fails to adapt to the efficient “NoThink” and “Flex” protocols, whereas our RL training yields massive gains (e.g., +21.5% in NoThink mode).

ers our adaptive “Flex” mode, where *FlexiVe* outperforms the base model by 14 percentage points (74.0% vs 60.0%), effectively unlocking a new, efficient inference capability.

4.4 Solve-Detect-Verify: An Efficient Alternative to BoN Ranking

We evaluate our *Solve-Detect-Verify* framework, demonstrating that its iterative refinement process is a more effective and efficient inference-time strategy than standard Best-of-N (BoN) ranking. The analysis is grounded in performance on the AIME 2024 benchmark.

Limitations of Standard BoN Ranking A common scaling strategy, BoN ranking, relies on an external verifier to select the best among N candidate solutions. However, our findings indicate this approach has significant limitations. As shown in Figure 5 (top-left), prominent verifiers like GenPRM-32B struggle to outperform even a simple majority vote baseline. We attribute this to ranking miscalibration, a known issue when verifiers evaluate the lengthy and complex reasoning traces typical of “thinking” models (Wu et al., 2025b). Unlike BoN, which depends on precise scalar scoring for ranking, our Solve-Detect-Verify (SDV) pipeline consistently achieves the highest accuracy across all sample sizes (N). At $N = 16$, SDV reaches 83.3% accuracy, surpassing the strong majority vote baseline (80.0%) and substantially outperforming GenPRM-32B BoN (66.7%). This suggests that *active* iterative self-correction is a more robust scaling mechanism than *passive* one-shot external

ranking.

The Advantage of Iterative Refinement The superior performance of SDV is attributable to its iterative refinement mechanism. Unlike BoN, which passively ranks static solutions, SDV actively improves upon them. Figure 5 (top-right) quantifies this benefit, showing a clear, monotonic increase in accuracy with each successive iteration on both the AIME 2024 and 2025 datasets. Unlike the parallel sampling (N) in the left panel of 5, this analysis tracks sequential refinement steps (T) on a single solution trajectory. For AIME 2024, accuracy improves from 60.0% after two iterations to over 70.0% after four, confirming that the refinement process is consistently productive.

Component-wise Token Efficiency The SDV pipeline is architected for efficiency, achieving superior accuracy without a corresponding increase in computational cost. The token breakdown in Figure 5 (bottom) provides a detailed analysis. The baseline ‘Solver LLM only’ approach uses an average of 12,788 tokens. **Detect** stage first prunes unnecessary generation paths, significantly reducing the average token count by over 35% to 8,204. **Verify** stage then applies targeted, corrective feedback, increasing the token count to 10,532 but yielding a substantial accuracy gain from 53.3% to 66.7%. Notably, we observe that the solver generates significantly fewer tokens during refinement compared to the initial phase. We hypothesize that while the base RL training encourages extensive exploration initially, the targeted feedback in the second pass constrains the search space, resulting in more con-

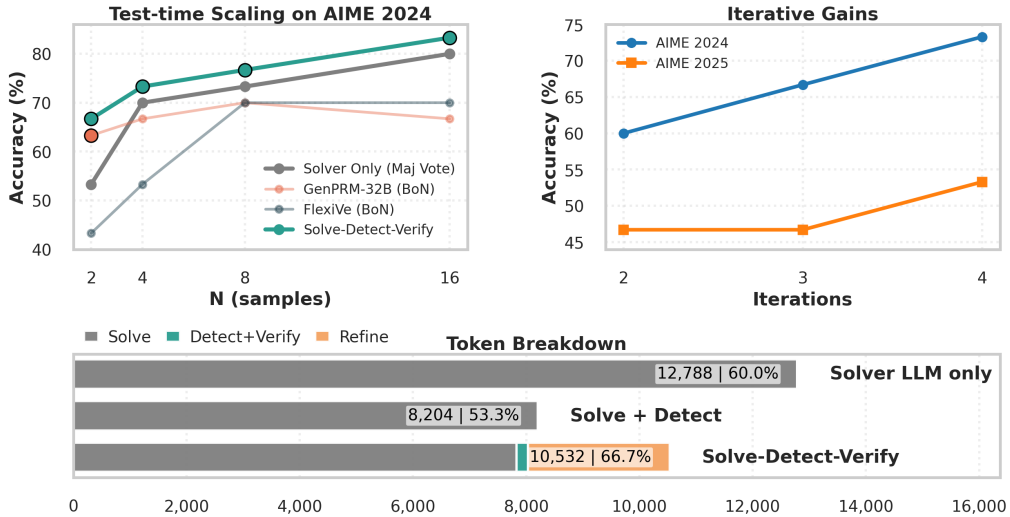


Figure 5: Performance and efficiency analysis of the Solve-Detect-Verify (SDV) pipeline on AIME 2024. **(Top-left)** SDV consistently outperforms standard Best-of-N (BoN) ranking methods in test-time accuracy scaling. **(Top-right)** The iterative nature of SDV yields monotonic accuracy improvements with each refinement step. **(Bottom)** A token breakdown for a single execution of the pipeline (Solve \rightarrow Detect \rightarrow Verify) reveals the pipeline’s efficiency: the ‘Detect’ stage reduces token usage, while the ‘Verify’ stage adds targeted computation to significantly boost accuracy, resulting in a net efficiency gain over the baseline solver.

cise corrections. The full SDV pipeline delivers a higher accuracy while consuming approximately 18% fewer tokens than the solver-only baseline, demonstrating a clear net gain in overall efficiency.

4.5 Discussions

Generalizability of Hesitation Detection We acknowledge that our hesitation keywords were derived empirically. To assess their generalizability, we evaluated the ‘Detect’ module on models with distinct training paradigms (Table 19). The results indicate that the mechanism’s effectiveness is tied to the training method. On RL-distilled models (e.g., Qwen3-8B), the detection behaves predictably, significantly reducing token usage (e.g., -3,576 tokens on AIME 2025) by pruning unproductive paths. Conversely, on SFT-trained models (e.g., S1-14B), the behavior is erratic, often increasing token usage (+2,374 tokens). This suggests that RL training instills a robust link between “verbalized hesitation” and model uncertainty, making our detection strategy a principled approach for the increasingly common class of RL-reasoning models.

Component Robustness and Qualitative Analysis Our extended analyses in the appendix validate the key design choices and robustness of our pipeline. The Flex@k verifier’s dynamic escalation is governed by a consensus threshold ($\tau = 0.8$) that optimally balances accuracy gains with a nearly 8x reduction in token usage compared to its full “slow thinking” mode (Appendix A.4.1, Table 18). Be-

Table 3: Sensitivity of Hesitation Keyword Detection Across Training Paradigms. RL-distilled models show consistent token reduction, whereas SFT models exhibit erratic behavior.

Model	Dataset	Baseline Acc.	S+D Acc.	Tok Δ
Qwen3-8B (RL-distilled)	AIME 24	83.3	60.9	-1,144
	AIME 25	73.3	66.7	-3,576
S1 14B (SFT-trained)	AIME 24	30.0	26.7	+2,206
	AIME 25	13.3	33.3	+2,374

sides, the iterative refinement loop demonstrates practical utility by successfully correcting 25% of incorrect initial solutions on AIME 2024 (Appendix A.4.3). Finally, We conduct a qualitative analysis of successful and failed feedback attempts to provide deeper insight into the correction process. Appendix A.4.4 illustrates two representative cases that demonstrate how the pipeline operates in practice.

5 Conclusion

We introduce *FlexiVe*, a dynamic verifier balancing computational cost and accuracy, integrated into the *Solve-Detect-Verify* pipeline for efficient LLM reasoning enhancement. Experiments confirm that our pipeline, leveraging *FlexiVe*, achieves significant gains in both accuracy and token efficiency over baselines. Our work demonstrates that the path to efficient and reliable LLM reasoning lies not only in developing flexible components but, critically, in designing intelligent pipelines that integrate them effectively.

555 **Limitation**

556 *FlexiVe* and *Solve-Detect-Verify* opens several ex-
557 citing avenues for future research. Our empirical
558 validation focuses on the challenging domain of
559 mathematical reasoning, a standard practice for
560 rigorously evaluating complex reasoning frame-
561 works (Zhong et al., 2025; Zhao et al., 2025b;
562 Zheng et al., 2024a; Wang et al., 2024b). A natural
563 and promising next step is to extend the demon-
564 strated benefits of *FlexiVe* to broader domains.
565 This presents a straightforward opportunity to adapt
566 the current “hesitation keywords”, an effective
567 heuristic for mathematical traces, to new linguistic
568 patterns. From a systems perspective, the pipeline’s
569 computational profile reflects a deliberate trade-
570 off for enhanced verification accuracy. We see a
571 clear path to optimizing this by integrating state-of-
572 the-art inference engines like vLLM (Kwon et al.,
573 2023) or SGLang (Zheng et al., 2024b). These fu-
574 ture steps represent a clear roadmap toward evol-
575 ving our framework into a more general-purpose,
576 highly efficient, and robust system for verified rea-
577 soning.

578 **Ethical Considerations**

579 Our work focuses on improving the reasoning capa-
580 bilities and inference efficiency of large language
581 models on publicly available mathematical bench-
582 mark datasets (gsm8k, math, olympiadbench, and
583 omnimath). We acknowledge the dual-use nature
584 of advanced AI problem-solvers; while they can
585 serve as valuable tools for education and research,
586 they could also be misused for academic dishon-
587 esty. The goal of our research is to contribute to the
588 scientific understanding of AI reasoning and create
589 more reliable and efficient systems, not to facilitate
590 misuse. Our method, *FlexiVe*, uses pre-trained
591 models without further fine-tuning, and we have
592 made no effort to remove existing safety guards.
593 We believe our work contributes to the transparent
594 and responsible development of AI systems.

References

- 595 Afra Feyza Akyurek, Ekin Akyurek, Ashwin Kalyan, Peter Clark, Derry Tanti Wijaya, and Niket Tandon. 2023. [RL4F: Generating natural language feedback with reinforcement learning for repairing model outputs](#). In *Proceedings of the 61st Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*, pages 7716–7733, Toronto, Canada. Association for Computational Linguistics. 649
- 596 [memory management for large language model serving with pagedattention](#). In *Proceedings of the ACM SIGOPS 29th Symposium on Operating Systems Principles (SOSP '23)*, page 1013–1029, New York, NY, USA. Association for Computing Machinery. 650
- 597 651
- 598 652
- 599 653
- 600
- 601 Zhong-Zhi Li, Haotian Wang, Kaiyan Zhang, Yancheng He, Yujia Xie, Yuxiang Huang, Zhengliang Shi, HongCheng Li, Wenxuan Wang, Zhiwei He, Dian Yu, Haitao Mi, Dong Yu, Jie Tang, and AnBo Zhang. 2025. From system 1 to system 2: A survey of reasoning large language models. *arXiv preprint arXiv:2502.17419*. 654
- 602 655
- 603 656
- 604 657
- 605 658
- 606 659
- 607 660
- 608 661
- 609 662
- 610 663
- 611 664
- 612 665
- 613 666
- 614 667
- 615 668
- 616 669
- 617 670
- 618 671
- 619 672
- 620 673
- 621 674
- 622 675
- 623 676
- 624 677
- 625 678
- 626 679
- 627 680
- 628 681
- 629 682
- 630 683
- 631 684
- 632 685
- 633 686
- 634 687
- 635 688
- 636 689
- 637 690
- 638 691
- 639 692
- 640 693
- 641 694
- 642 695
- 643 696
- 644 697
- 645 698
- 646 699
- 647 700
- 648 701
- 702

703	Deepseekmath: Pushing the limits of mathematical reasoning in open language models. <i>Preprint</i> , arXiv:2402.03300.	758
704		759
705		760
706	Nishad Singhi, Hritik Bansal, Arian Hosseini, Aditya Grover, Kai-Wei Chang, Marcus Rohrbach, and Anna Rohrbach. 2025. When to solve, when to verify: Compute-optimal problem solving and generative verification for llm reasoning. <i>Preprint</i> , arXiv:2504.01005.	761
707		762
708		763
709		764
710		765
711		766
712	Hanshi Sun, Momin Haider, Ruiqi Zhang, and 1 others. 2024. Fast best-of-n decoding via speculative rejection. In <i>Advances in Neural Information Processing Systems (NeurIPS)</i> .	767
713		768
714		769
715		770
716	Gladys Tyen, Hassan Mansoor, Victor Carbune, Peter Chen, and Tony Mak. 2024. LLMs cannot find reasoning errors, but can correct them given the error location. In <i>Findings of the Association for Computational Linguistics: ACL 2024</i> , pages 13894–13908, Bangkok, Thailand. Association for Computational Linguistics.	771
717		772
718		773
719		774
720		775
721		776
722		777
723	Leandro von Werra, Lewis Schmid, Thomas Wolf, and Lewis Tunstall. 2020-2024. Trl: Transformer reinforcement learning. https://github.com/huggingface/trl .	778
724		779
725		780
726		781
727	Jingxuan Wang, Yiming Ming, Zhengliang Shi, and 1 others. 2025. Inference-time scaling for complex tasks: Where we stand and what lies ahead. <i>arXiv preprint arXiv:2504.00294</i> .	782
728		783
729		784
730		785
731	Ke Wang, Houxing Ren, Aojun Zhou, and 1 others. 2024a. Mathcoder: Seamless code integration in llms for enhanced mathematical reasoning. In <i>International Conference on Learning Representations (ICLR)</i> .	786
732		787
733		788
734		789
735		790
736	Peiyi Wang, Lei Li, Zhihong Shao, Runxin Xu, and 1 others. 2024b. Math-Shepherd: Verify and reinforce LLMs step-by-step without human annotations. In <i>Proceedings of the 62nd Annual Meeting of the Association for Computational Linguistics (ACL)</i> .	791
737		792
738		793
739		794
740		795
741	Xuezhi Wang, Jason Wei, Dale Schuurmans, Quoc Le, Ed Chi, Sharan Narang, Aakanksha Chowdhery, and Denny Zhou. 2023. Self-consistency improves chain of thought reasoning in language models. In <i>International Conference on Learning Representations (ICLR)</i> .	796
742		797
743		798
744		799
745		800
746		801
747	Jason Wei, Xuezhi Wang, Dale Schuurmans, Maarten Bosma, Brian Ichter, Fei Xia, Ed Chi, Quoc Le, and Denny Zhou. 2022. Chain-of-thought prompting elicits reasoning in large language models. In <i>Advances in Neural Information Processing Systems</i> , volume 35, pages 24824–24837. Curran Associates, Inc.	802
748		803
749		804
750		805
751		806
752		807
753		808
754	Sean Welleck, Amanda Bertsch, Matthew Finlayson, and 1 others. 2024. From decoding to meta-generation: Inference-time algorithms for large language models.	809
755		810
756		811
757		812
		813
	Thomas Wolf, Lysandre Debut, Victor Sanh, Julien Chaumond, Clement Delangue, Anthony Moi, Pierric Cistac, Tim Rault, Remi Louf, Morgan Funtowicz, Joe Davison, Sam Shleifer, Patrick von Platen, Clara Ma, Yacine Jernite, Julien Plu, Canwen Xu, Teven Le Scao, Sylvain Gugger, and 3 others. 2020. Transformers: State-of-the-Art natural language processing. In <i>Proceedings of the 2020 Conference on Empirical Methods in Natural Language Processing: System Demonstrations</i> , pages 38–45, Online. Association for Computational Linguistics.	
	Yiming Wu, Zihan Sun, Sida Li, and 1 others. 2025a. Inference scaling laws: An empirical analysis of compute-optimal inference for llm problem-solving. In <i>International Conference on Learning Representations (ICLR)</i> .	
	Yiming Wu, Zihan Sun, Sida Li, and 1 others. 2025b. Scaling flaws of verifier-guided search in mathematical reasoning. <i>arXiv preprint arXiv:2502.00271</i> .	
	Noah Xie, AI AUTODIDAX, M Sarmad Parvez, Michael Song, Zhenqiao Zhang, Ziyu Chen, Shrimai Joshi, Robert Gmyr, Yufan Li, Siyuan Li, and 1 others. 2023. Reflexion: Language agents with verbal reinforcement learning. In <i>Advances in Neural Information Processing Systems (NeurIPS)</i> , volume 36.	
	Shunyu Yao, Dian Yu, Jeffrey Zhao, Izhak Sha, Thomas L Chen, Boyuan Rius, Yuxuan Du, Yang Liu, Zipeng Jiang, Tushar Han, and 1 others. 2023. Tree of thoughts: Deliberate problem solving with large language models. In <i>Advances in Neural Information Processing Systems (NeurIPS)</i> , volume 36.	
	Eric Zelikman, Yuhuai Wu, Jesse Mu, and Noah D. Goodman. 2022. Star: Bootstrapping reasoning with reasoning. In <i>Advances in Neural Information Processing Systems (NeurIPS)</i> .	
	Lunjun Zhang, Arian Hosseini, Hritik Bansal, Mehran Kazemi, Aviral Kumar, and Rishabh Agarwal. 2025. Generative verifiers: Reward modeling as next-token prediction. <i>Preprint</i> , arXiv:2408.15240.	
	Eric Zhao, Pranjal Awasthi, and Sreenivas Gollapudi. 2025a. Sample, scrutinize and scale: Effective inference-time search by scaling verification. <i>arXiv preprint arXiv:2502.00891</i> .	
	Jian Zhao, Runze Liu, Kaiyan Zhang, Zhimu Zhou, Junqi Gao, Dong Li, Jiafei Lyu, Zhouyi Qian, Biqing Qi, Xiu Li, and Bowen Zhou. 2025b. Genprm: Scaling test-time compute of process reward models via generative reasoning. <i>Preprint</i> , arXiv:2504.00891.	
	Chujie Zheng, Zhenru Zhang, Beichen Zhang, Runji Lin, Keming Lu, Bowen Yu, Dayiheng Liu, Jingren Zhou, and Junyang Lin. 2024a. Processbench: Identifying process errors in mathematical reasoning. <i>Preprint</i> , arXiv:2412.06559.	
	Lianmin Zheng, Siyuan Zhuang, Zhuohan Li, Cody Hao Yu, Lequn Li, Haotian Chen, Joseph E. Gonzalez, Ion Stoica, and Jonathan Ragan-Kelley. 2024b. SGLang:	

- 814 Efficient and expressive structured generation for
815 large language models. In *Proceedings of the 18th*
816 *Conference of the European Chapter of the Association*
817 *for Computational Linguistics (EACL 2024)*,
818 pages 1053–1071, St. Julian’s, Malta. Association for
819 Computational Linguistics.
- 820 Jianyuan Zhong, Zeju Li, Zhijian Xu, Xiangyu Wen,
821 and Qiang Xu. 2025. *Dyve: Thinking fast and*
822 *slow for dynamic process verification*. *Preprint*,
823 arXiv:2502.11157.
- 824 Aojun Zhou, Ke Wang, Zimu Lu, and 1 others. 2024a.
825 Solving challenging math word problems using gpt-4
826 code interpreter with code-based self-verification. In
827 *International Conference on Learning Representations*
828 *(ICLR)*.
- 829 Jin Peng Zhou, Charles Staats, Wenda Li, and 1 others.
830 2024b. Don’t trust: Verify-grounding llm quantitative
831 reasoning with autoformalization. *arXiv preprint*
832 arXiv:2403.18120.

A Appendix

The Use of Large Language Models

We use Large Language Models (LLMs), including ChatGPT and Gemini, solely for the purpose of editing and polishing the writing in this paper.

Broader Impact

The development of FlexiVe and the Solve-Detect-Verify pipeline represents a significant step toward making advanced AI reasoning systems more practical, reliable, and efficient. By designing a verifier that dynamically allocates computational resources—switching between rapid “fast thinking” and meticulous “slow thinking”—our framework directly confronts the critical trade-off between accuracy and efficiency that currently limits the deployment of large models. This approach promotes a more sustainable and scalable paradigm for AI reasoning, reducing the reliance on computationally expensive, brute-force methods like Best-of-N sampling with process-based verifiers. Our work has the potential to enhance trust and safety in AI systems. By not only identifying but also pinpointing the exact location of errors and providing targeted feedback for correction, our pipeline improves the interpretability and debuggability of the reasoning process. This iterative refinement is crucial for high-stakes domains where reliability is paramount, such as automated scientific discovery, medical diagnostics, and educational tools. By making state-of-the-art reasoning more computationally accessible, our work also helps democratize advanced AI, enabling powerful capabilities to run in more resource-constrained environments. This research paves the way for future investigations into more sophisticated self-correcting systems and adaptive computation, pushing the frontier of efficient and trustworthy artificial intelligence.

A.1 Implementation Details and Experimental Setup

This section provides comprehensive details regarding the training of FlexiVe, the implementation of the Solve-Detect-Verify pipeline, evaluation benchmarks, and specific implementation clarifications.

A.1.1 FlexiVe Training

Training Protocol and Rationale We train FlexiVe using Group Relative Policy Optimization (GRPO) (Shao et al., 2024a) initialized from the DeepSeek-R1-Distill-Qwen-14B model (Shao et al., 2024a). We utilize the BIG-Bench Mistake dataset (Tyen et al., 2024), using 1,526 samples for training and 170 for testing, derived from a 90%/10% split. The objective is to predict the first error index (idx_{gt}) or -1 if correct, optimized using the composite reward (Section 3.2, main paper). Training initially focused on optimizing the “Fast Thinking” mode (NoThinking activated) to instill efficient, accurate error detection with minimal verbosity. This strategy established a strong, low-cost baseline and promoted data efficiency, providing a robust foundation that generalized well to the “Slow Thinking” mode. Statistics for the training data are provided in Table 4.

Table 4: Details of the model and dataset used for training.

Items	Values
Model	FlexiVe-14B
Benchmark	BIG-Bench Mistake
Train Set Size	1,526
Test Set Size	170

RL vs. SFT Generalization As discussed in the main paper (Section 4.2), our RL approach demonstrated superior generalization compared to Supervised Fine-Tuning (SFT). An SFT baseline trained on 10,000 complex reasoning paths showed poor generalization when evaluated on the diverse, often simpler traces in ProcessBench. In contrast, FlexiVe, RL-trained on only 1,526 samples, generalized effectively. This highlights RL’s advantage in fostering robust verifiers capable of handling diverse reasoning styles and complexities, even with significantly less data.

Hyperparameters and Optimization We employed LoRA (Hu et al., 2022) targeting attention projection layers and used the AdamW (Loshchilov and Hutter, 2019) optimizer with gradient checkpointing.

874
875
876

Training utilized the transformers (Wolf et al., 2020) and trl (von Werra et al., 2020-2024) libraries, tracked via Weights & Biases (Biewald, 2020). The key hyperparameters and optimization settings are summarized in Table 5.

Table 5: Training details.

Parameter	Value	Description
Base Model	DeepSeek-R1-Distill-Qwen-14B	Base model for initialization
Learning Rate	5×10^{-6}	Initial learning rate
Batch Size	1	Per-device batch size
Num Train Epochs	3	Number of training epochs
Gradient Accum. Steps	8	Gradient accumulation steps
PEFT / LoRA	True (r=16, α =32)	Adapter fine-tuning (LoRA)
LR Scheduler Type	Linear	Learning Rate Scheduler Type
Optimizer	AdamW	Optimization algorithm
Warmup Steps	100	Number of warmup steps
GRPO Group Size	14	Number of generations per prompt
KL Coefficient	0.04	KL penalty coefficient for GRPO

877
878
879
880
881
882

A.1.2 Evaluation Benchmarks

We assessed our framework on a suite of challenging mathematical reasoning benchmarks. For evaluating step-level verification performance, we used the four standard splits of the **ProcessBench** benchmark: GSM8K, MATH, Olympiad-Bench, and OmniMATH. For evaluating the end-to-end performance of the full Solve-Detect-Verify pipeline, we used problems from the **AIME 2024** and **AIME 2025** competitions. The number of problems in the test set for each benchmark is detailed in Table 6.

Table 6: Details of datasets used for model evaluation.

Benchmark	Test Set Size
<i>ProcessBench Splits</i>	
GSM8K	400
MATH	1,000
Olympiad-Bench	1,000
OmniMATH	1,000
<i>End-to-End Evaluation</i>	
AIME 2024	30
AIME 2025	30

883
884
885
886
887
888
889
890

A.1.3 Solve-Detect-Verify Pipeline

The Solve-Detect-Verify pipeline employs an adaptive, iterative strategy. We detail the specific logic for the streaming ‘Solve-Detect’ mechanism in Algorithm 1 and the full iterative refinement control flow in Algorithm 2.

Solve and Detect Logic Algorithm 1 outlines the streaming detection mechanism. It generates tokens step-by-step, pausing upon encountering a hesitation keyword ($\mathcal{K}_{hesitation}$) to check for solution completeness using the solver’s own log-probabilities. This granular control allows the system to exit the generation loop early if the model signals it has finished thinking.

Algorithm 1 Solve-Detect Stage Implementation

Require: Problem P , Solver \mathcal{M}_{solve} **Ensure:** Candidate Solution S_1

```
1: procedure SOLVEDetect( $P, \mathcal{M}_{solve}$ )
2:    $S_1 \leftarrow \emptyset$ 
3:    $stop\_flag \leftarrow false$ 
4:   for  $k = 1$  to  $L_{max}$  do ▷ Token-by-token generation
5:      $t_k \sim \mathcal{M}_{solve}(\cdot | P, S_1^{(k-1)})$ 
6:      $S_1^{(k)} \leftarrow S_1^{(k-1)} \oplus t_k$ 
7:     if  $t_k = EOS$  then
8:        $stop\_flag \leftarrow true$ 
9:     if  $S_1^{(k)}$  ends with  $kw \in \mathcal{K}_{hesitation}$  then
10:       $logp_{Yes} \leftarrow \log p_{\mathcal{M}_{solve}}(Yes | PromptComplete(S_1^{(k)}))$ 
11:       $logp_{No} \leftarrow \log p_{\mathcal{M}_{solve}}(No | PromptComplete(S_1^{(k)}))$ 
12:      if  $logp_{Yes} > logp_{No}$  then ▷ Check solution completeness
13:         $stop\_flag \leftarrow true$ 
14:      if  $stop\_flag$  then
15:        break
16:       $S_1 \leftarrow S_1^{(k)}$ 
17:   return  $S_1$ 
```

Full Iterative Pipeline Algorithm 2 describes the higher-level control flow, integrating the SolveDetect procedure with the FlexiVe verification module to perform iterative refinement up to T attempts. 891
892

Algorithm 2 Solve-Detect-Verify Pipeline Implementation Flow

Require: Problem P , Verification Parameters $\Theta_V = (k_{fast}, \tau_{agree}, k_{slow})$, Max Attempts T

```
1:  $S_{current} \leftarrow NIL$ 
2:  $F_{prev} \leftarrow NIL$ 
3: for  $t = 1$  to  $T$  do
4:   ▷ — Solve and Detect (See Algorithm 1) —
5:    $Prompt_t \leftarrow FormatPrompt(P, S_{current}, F_{prev})$ 
6:    $S_t \leftarrow SOLVEDetect(P, \mathcal{M}_{solve})$ 
7:    $S_{current} \leftarrow S_t$ 
8:   if  $t < T$  then ▷ Verify if not the last attempt
9:     ▷ — Verify (FlexiVe) —
10:     $(is\_valid_t, error\_step_t, F_t) \leftarrow AdaptiveVerify(P, S_t, \Theta_V)$ 
11:    if  $is\_valid_t = True$  then
12:      break ▷ Solution verified, terminate early
13:    else
14:       $F_{prev} \leftarrow F_t$  ▷ Prepare feedback for refinement
15:   return  $S_{current}$ 
```

Solve Module and Prompts We employ DeepSeek-R1-14B/32B as the solver LLM. The initial prompt (Figure 6) guides the model to generate a structured solution. If refinement is required ($t > 1$), a retry prompt (Figure 7) incorporates feedback from FlexiVe (F_{prev}). 893
894
895

Detect Module The GenerateSolutionWithDetection function implements a streaming detection framework to identify and curtail “overthinking.” 896
897

- **Hesitation Keywords:** Generation is monitored for hesitation cues (Figure 8). These keywords were derived empirically by observing common phrases signaling a pause or self-correction in LLM 898
899

LLM Initial Solver Prompt

```
The following is a math problem:  
[Math Problem]  
{question}  
Solve it step by step. For each step, you  
should use \n\n in the end.  
Please put your final answer (i.e., the  
index) in \boxed{{}}.
```

Figure 6: LLM Initial Solver Prompt.

LLM Retry Prompt with Feedback

```
The following is a math problem:  
[Math Problem] {question}  
  
You previously attempted to solve this:  
[Previous Solution]  
{previous_solution}  
  
The feedback is:  
[Verification Feedback]  
{verifier_feedback}  
  
Please correct your solution.  
Provide a complete, new solution.  
Put your final answer in \boxed{{}}.
```

Figure 7: LLM Retry Prompt with Feedback.

900 outputs.

- 901 • **Completeness Check:** Upon detecting a keyword, the proposer is suspended. A Detector LLM
902 (the same base model) evaluates the context using the prompt in Figure 9. We compare the log-
903 probabilities of “Yes” and “No” to determine completeness.
- 904 • **Efficiency (KV Cache Reuse):** The ‘Detect’ module achieves high efficiency by leveraging
905 vLLM (Kwon et al., 2023) with prefix caching. Since the detection prompt is a continuation
906 of the existing generation context, vLLM automatically reuses the KV cache from preceding steps,
907 leading to minimal overhead (more than 90% reuse).
- 908 • **“Continue-after-detected” Logic:** If completeness is detected, the generation might be briefly
909 continued to ensure the current thought segment is fully articulated before truncation, facilitating
910 better context for potential sequential revision.

Hesitation Keywords

```
Wait, double-check, Alternatively, Hmm,  
Let me check, Alright, make sure,  
Another way, Let me verify, to confirm,  
Looking back, But wait
```

Figure 8: Hesitation keywords monitored for detection.

LLM Detection Prompt

You are a solution completeness checker.
Given current solution to a math problem,
determine if it is a complete solution
(i.e., contains a final answer).
Respond with exactly one word: `Yes` if
complete, `No` otherwise.

Figure 9: LLM Detection Prompt.

Verify Module (AdaptiveVerify) This function implements the Flexible Allocation of Verification Budget (Section 3.2). It conducts k_{fast} “Fast Thinking” runs. If the agreement ratio meets τ_{agree} , the consensus is returned. Otherwise, it escalates to k_{slow} “Slow Thinking” runs. Across all experiments, k_{slow} is consistently set to $\lceil k_{fast}/8 \rceil$, balancing cost reduction with sufficient analysis to resolve ambiguities.

A.1.4 Evaluation Benchmarks and Baselines

FlexiVe Evaluation We assess step-level verification capabilities (F1 score) using ProcessBench (Zheng et al., 2024a) (GSM8K, MATH, OlympiadBench, OmniMATH). We compare against SOTA Process Reward Models (PRMs), including GenPRM (Zhao et al., 2025b) and Dyve (Zhong et al., 2025).

Pipeline Evaluation We evaluate the end-to-end effectiveness of the Solve-Detect-Verify pipeline on challenging mathematical datasets: AIME (2024, 2025) (Hugging Face, 2024; OpenCompass, 2025), AMC, CNMO (Liu et al., 2024) (China’s National Mathematical Olympiad), and OlympiadBench. We measure accuracy and efficiency (tokens, TFLOPS). We use DeepSeek-R1 14B/32B (Shao et al., 2024a) as the base worker LLMs, comparing against direct generation and Self-Consistency (Wang et al., 2023).

Compute Categories In Table 10, models are categorized by computational effort:

- **Moderate Compute:** Involves a reasonable number of samples without code execution (e.g., GenPRM Maj@8 w/o code, FlexiVe Flex@k). The adaptive nature of Flex@k keeps the average compute moderate.
- **High Compute:** Prioritizes maximal accuracy using extensive verification or intensive techniques (e.g., GenPRM Maj@8 w/ Code Exec, FlexiVe Think@64).

A.2 Detailed Experimental Results

A.2.1 FlexiVe Performance Scaling

Tables 7 (Think@k), 8 (NoThinking@k), and 9 (Flex@k) provide detailed F1 scores and total token consumption (in Millions, M) for FlexiVe across ProcessBench subsets.

Table 7: Performance of FlexiVe “With Thinking” (Think@k) on ProcessBench subsets. Tokens are total generated (Millions) across the respective test set.

k	GSM8K		MATH		OlympiadBench		OmniMATH	
	F1 (%)	Tokens (M)	F1 (%)	Tokens (M)	F1 (%)	Tokens (M)	F1 (%)	Tokens (M)
2	82.3	2.4	81.9	5.2	78.0	8.4	71.3	7.1
4	86.7	4.8	86.4	10.4	84.3	16.8	76.9	14.3
8	86.4	9.5	88.9	20.9	85.4	33.4	78.9	28.6
16	87.6	19.2	89.7	41.8	86.5	66.9	80.1	57.1
32	87.7	38.1	89.7	83.8	86.7	133.6	80.6	114.2
64	87.8	76.3	90.1	167.5	86.7	267.3	80.4	228.4
128	88.1	152.7	90.0	335.4	86.7	534.1	80.5	456.4

Table 8: Performance of FlexiVe “Without Thinking” (NoThinking@k) on ProcessBench subsets.

k	GSM8K		MATH		OlympiadBench		OmniMATH	
	F1 (%)	Tokens (M)	F1 (%)	Tokens (M)	F1 (%)	Tokens (M)	F1 (%)	Tokens (M)
2	61.5	0.4	57.2	1.5	49.0	1.9	50.5	1.6
4	66.8	0.7	61.3	3.0	53.8	3.7	52.5	3.3
8	66.7	1.5	62.8	6.1	55.2	7.5	53.6	6.6
16	66.8	3.0	64.3	12.1	55.9	15.0	54.2	13.3
32	66.5	5.9	64.4	24.2	55.9	29.9	54.7	26.5
64	66.8	11.8	64.2	48.5	56.1	59.8	54.0	52.9
128	66.7	23.7	65.0	96.8	56.3	119.8	54.1	105.9

Table 9: Performance of FlexiVe with Flexible Allocation (Flex@k) on ProcessBench subsets.

k	GSM8K		MATH		OlympiadBench		OmniMATH	
	F1 (%)	Tokens (M)	F1 (%)	Tokens (M)	F1 (%)	Tokens (M)	F1 (%)	Tokens (M)
2	72.97	0.2	72.92	1.0	67.43	1.3	61.41	1.3
4	78.43	0.3	77.67	1.5	72.41	2.1	67.34	2.1
8	75.75	0.5	78.86	3.1	70.06	4.3	66.57	4.2
16	76.88	0.9	78.20	6.1	73.07	8.2	68.94	7.9
32	82.84	2.1	83.30	13.9	79.23	19.5	73.40	18.8
64	82.00	4.3	83.63	28.7	79.26	39.6	74.67	38.6
128	83.02	8.9	84.96	59.1	79.98	80.8	75.23	78.5

Analysis of Trade-offs and Efficiency The data demonstrates the distinct trade-offs. Think@k establishes the accuracy upper bound at the highest cost. NoThinking@k is the most efficient but has the lowest accuracy ceiling. Flex@k effectively balances these extremes. On MATH@128, Flex@k (84.96% F1, 59.1M tokens) achieves an 82.4% token reduction compared to Think@128 (90.0% F1, 335.4M tokens).

At $k = 128$, Flex@k uses approximately 86.1% fewer tokens on average than Think@k. Notably, at higher k values, Flex@k can be both more accurate and more token-efficient than NoThinking@k (e.g., on GSM8K and MATH).

Visualizing F1 Scaling Figure 10 visualizes the F1 score scaling corresponding to the data above. Flex@k consistently outperforms NoThinking@k and generally matches or exceeds the DS14B baseline, confirming the effectiveness of the adaptive approach.

A.2.2 Comprehensive ProcessBench Results

Table 10 provides a comprehensive comparison of FlexiVe on ProcessBench. FlexiVe demonstrates strong performance and remarkable sample efficiency, achieving SOTA results despite being trained on only 1,526 samples, compared to 23K-404K samples for other models. In the Moderate Compute category, Flex@128 achieves the best average F1 (80.8%). In the High Compute category, Think@64 establishes a new SOTA for open-source models (86.3% Avg F1).

A.2.3 Statistical Significance and Stability

To validate robustness, we simulated the voting process 10 times from a pool of 512 cached completions to generate 95% confidence intervals for FlexiVe’s performance (Tables 15, 16, and 17, showing selected k values for brevity).

The analysis finds: (1) **Think@k** shows high stability (tight intervals $\leq 1\%$). (2) **NoThink@k** exhibits higher variance (wider intervals 2-5%). (3) **Flex@k** achieves a balanced trade-off (moderate intervals 1-4%), validating the reliability of the adaptive approach.

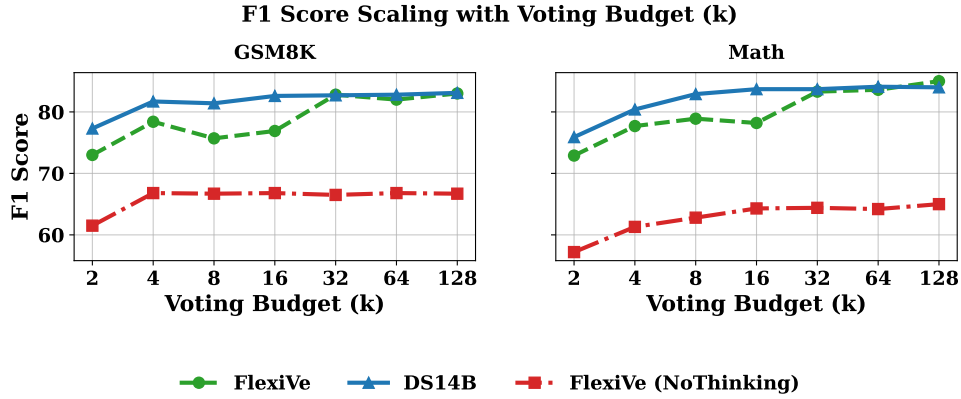


Figure 10: F1 score scaling with voting budget k on GSM8K (left) and MATH (right). FlexiVe (Flex@ k , green circles) improves with larger k , performing comparably or better than DS14B (blue triangles, baseline verifier), while both surpass the FlexiVe (NoThinking variant, red squares).

Table 10: ProcessBench results reported with F1 scores. Results for *FlexiVe* are highlighted. **bold** indicates the best in the sub category. All *FlexiVe* variants are trained on only 1526 samples.

Model	# Samples	GSM8K	MATH	Olympiad Bench	Omni-MATH	Avg.
<i>Proprietary Models</i>						
GPT-4o-0806	unk	79.2	63.6	51.4	53.5	61.9
o1-mini	unk	93.2	88.9	87.2	82.4	87.9
<i>Open Source Models (1.5B)</i>						
Skywork-PRM-1.5B	unk	59.0	48.0	19.3	19.2	36.4
GenPRM-1.5B (Pass@1) w/ Code Exec	23K	52.8	66.6	55.1	54.5	57.3
<i>Open Source Models (7-8B)</i>						
Math-Shepherd-PRM-7B	445K	47.9	29.5	24.8	23.8	31.5
RLHFlow-PRM-Mistral-8B	273K	50.4	33.4	13.8	15.8	28.4
EurusPRM-Stage2	30K	47.3	35.7	21.2	20.9	31.3
Qwen2.5-Math-PRM-7B	~344K	82.4	77.6	67.5	66.3	73.5
RetrievalPRM-7B	404K	74.6	71.1	60.2	57.3	65.8
Universal-PRM-7B	unk	85.8	77.7	67.6	66.4	74.3
Direct Generative PRM-7B	23K	63.9	65.8	54.5	55.9	60.0
GenPRM-7B w/ Code Exec (Pass@1)	23K	78.7	80.3	72.2	69.8	75.2
GenPRM-7B w/ Code Exec (Maj@8)	23K	81.0	85.7	78.4	76.8	80.5
<i>Open Source Models (14-32B) w/ Moderate Compute</i>						
Dyve-14B	117K	68.5	58.3	49.0	47.2	55.8
GenPRM-32B w/o Code Exec (Maj@8)	23K	78.8	85.1	78.7	74.9	79.3
<i>FlexiVe</i> (Flex@32)	1526	82.8	83.3	79.2	73.4	79.7
<i>FlexiVe</i> (Flex@128)	1526	83.0	85.0	80.0	75.2	80.8
<i>Open Source Models (14-32B) w/ High Compute</i>						
GenPRM-32B (Pass@1) w/ Code Exec	23K	83.1	81.7	72.8	72.8	77.6
GenPRM-32B (Maj@8) w/ Code Exec	23K	85.1	86.3	78.9	80.1	82.6
<i>FlexiVe</i> (Think@64)	1526	88.1	90.1	86.7	80.4	86.3

A.3 Extended Pareto Frontier and Efficiency Analysis

This section provides the detailed data and extended analysis corresponding to the Pareto frontier discussion in Section 4.3 of the main paper. We analyze the accuracy-efficiency trade-off of *FlexiVe* against

957

958

959

GenPRM (?), evaluating its dominance in both theoretical TFLOPS and empirical wall-clock time.

A.3.1 Underlying Data

Table 11 lists the exact data points for the Pareto frontier presented in Figure 11 (main paper). The comparison metrics are F1 scores and Relative TFLOPS on the MATH split of ProcessBench (?).

Table 11: Detailed Data for Pareto Frontier Analysis on MATH Dataset (Corresponding to Figure 11 in main paper).

Model	Config (@k)	F1 Score (%)	Relative TFLOPS
FlexiVe (Flex)	@2	72.9	1.9
	@4	75.8	3.8
	@8 (Best Trade-off)	78.9	7.5
	@16	78.2	13.2
	@32	83.3	27.2
FlexiVe (Think)	@1	81.9	6.1
	@2	82.3	10.2
	@4	86.4	12.3
	@8	88.9	22.8
GenPRM-7B (Maj)	@1	80.3	15.1
	@8	83.1	28.1
GenPRM-32B (Maj)	@1	80.0	13.4
	@8	83.5	29.4

A.3.2 Efficiency Analysis: TFLOPS vs. Wall-Clock Time

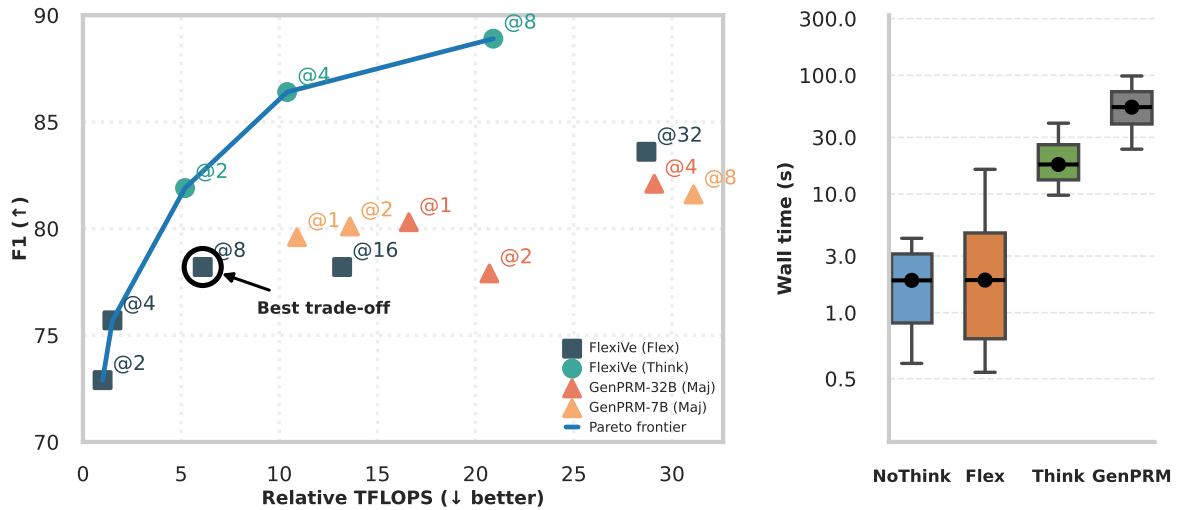


Figure 11: Pareto frontier analysis on the MATH dataset. (Left) F1 Score versus Relative TFLOPS. *FlexiVe* (Think@k) establishes the state-of-the-art frontier, achieving higher F1 scores at a lower computational cost relative to GenPRM-7B and GenPRM-32B. (Right) A comparison of wall-clock time. *FlexiVe* demonstrates substantially lower latency; its Flex mode is comparable to the NoThink baseline, while its Think mode is approximately 2.8x faster than GenPRM. The *FlexiVe* (Flex@8) configuration is identified as an optimal trade-off point.

The performance gap between *FlexiVe* and process-based baselines stems from key architectural differences: *FlexiVe* is a **holistic verifier** that processes traces in a single pass, while GenPRM is a **process-based verifier** that re-evaluates an expanding context at each step, leading to non-linear cost scaling.

TFLOPS Efficiency. As illustrated in Figure 11, *FlexiVe* is significantly more compute-efficient. The *FlexiVe* (Think@k) configurations define a new state-of-the-art Pareto frontier[cite: 19]. For instance, *FlexiVe* (Think@4) attains a higher F1 score ($\sim 87\%$) than the best GenPRM-32B model ($\sim 84\%$) while using less than half the computation (~ 12 vs. ~ 29 Relative TFLOPS).

Wall-Clock Latency Advantage. While TFLOPS provide a theoretical measure of compute, wall-clock time reveals the distinct advantage of the **Flex** mode. While *FlexiVe* (Think@2) offers competitive TFLOPS efficiency, it requires executing high-latency “Slow Thinking” traces sequentially. In contrast, *FlexiVe* (Flex@8) executes eight low-latency “Fast Thinking” runs in parallel, escalating to slow thinking only when necessary.

The left of Figure 11 presents the empirical wall-clock measurements. The Flex mode results in drastically lower latency: it achieves a median wall time of ~ 2 s, which matches the NoThink baseline[cite: 179]. Conversely, GenPRM exhibits a median latency of ~ 50 s, making it approximately $25\times$ slower than *FlexiVe* (Flex) and $2.8\times$ slower than *FlexiVe* (Think). Thus, while TFLOPS are comparable between certain configurations, Flex provides the superior accuracy-latency trade-off essential for real-world deployment.

Optimal Trade-off Selection. Based on this analysis, we identify the *FlexiVe* (Flex@8) configuration as the optimal trade-off. It achieves a substantial F1 score of 78.9% with a modest computational cost of 7.5 relative TFLOPS (Table 11)[cite: 1618, 1744]. This point represents the “knee” of the performance curve—securing most of the accuracy gains without the high expense of premium Think modes. Consequently, we adopt the Flex@8 setting for *FlexiVe* in all Solve-Detect-Verify pipeline experiments.

A.3.3 Solve-Detect-Verify Pipeline Performance Data (Figure 5)

This section provides the underlying data supporting the analysis presented in Section 4.4 and Figure 5 (main paper), focusing on the AIME 2024 benchmark.

Scaling Performance (BoN vs. SDV) Table 12 details the accuracy scaling as the number of samples (N) increases. The SDV pipeline consistently outperforms both simple majority voting and BoN ranking using external verifiers.

Table 12: Test-time Accuracy Scaling on AIME 2024 (Data supporting Figure 5, Top-Left).

Method	N=2	N=4	N=8	N=16
Solver Only (Maj Vote)	53.3	70.0	73.3	80.0
GenPRM-32B (BoN)	63.3	66.7	70.0	66.7
FlexiVe (BoN)	43.3	53.3	70.0	70.0
Solve-Detect-Verify	66.7	73.3	76.7	83.3

Iterative Gains Table 13 demonstrates the monotonic accuracy improvements achieved through the iterative refinement process of the SDV pipeline.

Table 13: Iterative Refinement Gains (Data supporting Figure 5, Top-Right).

Iterations	AIME 2024 Accuracy (%)	AIME 2025 Accuracy (%)
2	60.0	46.7
3	66.7	46.7
4	73.3	53.3

Token Efficiency Breakdown Table 14 details the average token usage and accuracy at each stage of the pipeline, illustrating the efficiency gains from the ‘Detect’ stage and the accuracy boost from the ‘Verify’ stage.

Table 14: Token Efficiency Breakdown on AIME 2024 (Data supporting Figure 5, Bottom).

Configuration	Average Tokens	Accuracy (%)
Solver LLM only	12,788	60.0
Solve + Detect	8,204	53.3
Solve-Detect-Verify	10,532	66.7

A.3.4 Scaling Properties

We explore scaling *Solve-Detect-Verify* along two dimensions: the verifier budget (Flex@N) and the solver budget (Number of Solutions).

Scaling Verifier Budget (Flex@N): We analyze scaling *FlexiVe*’s budget within a single pipeline run (Figure 12). The ‘w/o Flex’ setup significantly cuts token usage (e.g., 0.67 ratio on AIME2024) but reduces accuracy. Integrating ‘Flex@8’ substantially boosts accuracy over the baseline (e.g., 73.3% vs. 56.6% on AIME2024) while still using fewer tokens (0.96 ratio).

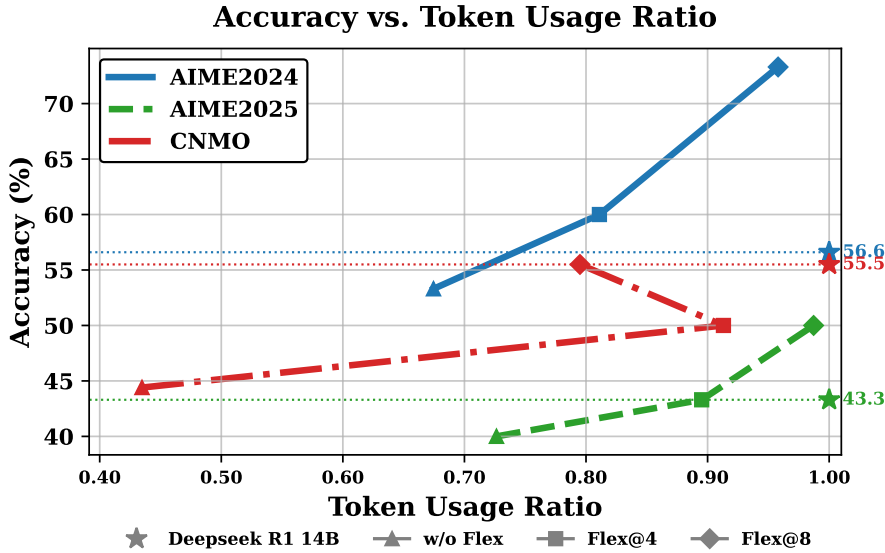


Figure 12: Impact of scaling *FlexiVe*’s verification budget (Flex@N) within a single *Solve-Detect-Verify* execution on Pass@1 Accuracy vs. Token Usage Ratio relative to DeepSeek R1 14B.

Scaling Solver Budget: To achieve higher peak accuracies, we scale compute by generating multiple solutions from the solver. On AIME2024 (Figure 5, top left panel), this strategy yields significant improvements, climbing from 67.5% (1 solution) to over 83% (16 solutions), requiring approximately 4x fewer solutions than the baseline to reach similar accuracy levels.

A.3.5 FlexiVe Performance Scaling Details

This section provides a more detailed breakdown of the performance scaling for the different configurations of our *FlexiVe* method. We present the 95% confidence intervals for accuracy on four benchmark datasets as the number of voting samples (N) increases.

The results are detailed for the ‘‘With Thinking’’ configuration (Think@k) in Table 15, the ‘‘Without Thinking’’ configuration (NoThink@k) in Table 16, and our primary *FlexiVe* method (Flex@k) in Table 17.

A consistent trend is evident across all tables: performance generally improves as the number of voting samples (N) increases from 2 to 128. For example, for the main Flex@k method on the math dataset, accuracy climbs from 72.9% to 85.0%. Concurrently, the confidence intervals tend to narrow with a larger N , indicating more stable and reliable results. These tables also quantitatively show that the Think@k method consistently achieves the highest performance, while NoThink@k establishes a performance baseline.

Table 15: 95% Confidence Intervals for FlexiVe “With Thinking” (Think@k).

Voting N	gsm8k	math	olympiadbench	omnimath
2	82.3 ± 0.89	81.9 ± 0.67	78.0 ± 0.38	71.3 ± 0.56
8	86.4 ± 0.50	88.9 ± 0.21	85.4 ± 0.40	78.9 ± 0.19
32	87.7 ± 0.44	89.7 ± 0.24	86.7 ± 0.33	80.6 ± 0.21
128	88.1 ± 0.32	90.0 ± 0.15	86.7 ± 0.15	80.5 ± 0.09

Table 16: 95% Confidence Intervals for FlexiVe “Without Thinking” (NoThink@k).

Voting N	gsm8k	math	olympiadbench	omnimath
2	61.5 ± 2.36	57.2 ± 4.28	49.0 ± 2.56	50.5 ± 3.55
8	66.7 ± 2.63	62.8 ± 4.91	55.2 ± 3.32	53.6 ± 4.05
32	66.5 ± 2.35	64.4 ± 4.96	55.9 ± 2.98	54.7 ± 3.93
128	66.7 ± 2.46	65.0 ± 5.09	56.3 ± 3.23	54.1 ± 4.08

Table 17: 95% Confidence Intervals for FlexiVe (Flex@k).

Voting N	gsm8k	math	olympiadbench	omnimath
2	73.0 ± 2.74	72.9 ± 4.08	67.4 ± 2.62	61.4 ± 3.34
8	75.8 ± 2.48	78.9 ± 2.85	70.1 ± 2.05	66.6 ± 2.57
32	82.8 ± 1.17	83.3 ± 2.38	79.2 ± 1.41	73.4 ± 2.39
128	83.0 ± 1.32	85.0 ± 1.48	80.0 ± 1.51	75.2 ± 2.47

A.4 Extended Discussions and Analysis

1024

A.4.1 Sensitivity Analysis of Consensus Threshold τ

1025

The consensus threshold τ governs the escalation from “Fast Thinking” to “Slow Thinking” in the Flex@k strategy. We performed a sensitivity analysis (Table 18) to validate our choice of $\tau = 0.8$.

1026

1027

Table 18: Sensitivity analysis for the consensus threshold τ in Flex@8, averaged across ProcessBench datasets. As τ varies, performance shifts between the ‘NoThink@8’ and ‘Think@8’ baselines.

Consensus Threshold (τ)	Slow Thinking Escalation (%)	Avg. F1 Score (%)	Avg. Total Tokens (M)
<i>NoThink@8 Baseline</i>	0%	59.6	5.4
0.5	5%	61.0	5.2
0.7	18%	69.5	4.1
0.8 (Chosen)	28%	72.9	3.0
0.95	80%	83.5	19.5
<i>Think@8 Baseline</i>	100%	84.9	23.1

At a low threshold ($\tau = 0.5$), escalation is minimal (5%), and performance approaches the ‘NoThink@8’ baseline. At a strict threshold ($\tau = 0.95$), the system escalates 80% of cases, approaching the ‘Think@8’ baseline but at a massive computational cost. Our chosen value, $\tau = 0.8$, represents the optimal balance, significantly raising the F1 score (72.9%) while maintaining high efficiency (3.0M tokens, nearly 8x lower than Think@8).

1028

1029

1030

1031

1032

A.4.2 Robustness of the Detection Mechanism

1033

As discussed in Section 4.5 (main paper), the robustness of the hesitation keyword detector depends on the model’s training paradigm (Table 19). On RL-distilled models (e.g., Qwen3-8B), the mechanism behaves predictably. However, on SFT-trained models (e.g., S1 14B), its behavior is erratic, sometimes increasing token usage and causing unpredictable accuracy shifts. This suggests RL instills a more reliable link between hesitation keywords and model uncertainty.

1034

1035

1036

1037

1038

Table 19: Sensitivity of Hesitation Keyword Detection Across Training Paradigms. (Table 2 in main paper).

Model (Training Paradigm)	Dataset	Baseline Acc. (%)	Solve+Detect Acc. (%)	Acc. $\Delta(pp)$	Token Δ
Qwen3-8B (RL-distilled)	AIME 2024	83.3	60.9	-22.4	-1,144
	AIME 2025	73.3	66.7	-6.6	-3,576
S1 14B (SFT-trained)	AIME 2024	30.0	26.7	-3.3	+2,206
	AIME 2025	13.3	33.3	+20.0	+2,374

1039 A.4.3 Effectiveness of the Refinement Loop

1040 We analyzed the refinement success rate on the AIME 2024 dataset. Out of 16 initial solutions that were
 1041 incorrect (S1), our pipeline successfully corrected 4 of them (S2), yielding a **25% success rate**. This
 1042 demonstrates the practical utility of the refinement mechanism, particularly noteworthy as FlexiVe was
 1043 not fine-tuned on the solver’s specific traces, indicating good generalization.

1044 A.4.4 Qualitative Analysis of the Feedback Mechanism

1045 We analyzed successful and failed feedback attempts to provide deeper insight into the correction process.

Successful S1 \rightarrow S2 Correction

Problem: Every morning Aya goes for a 9-kilometer-long walk... When she walks at a constant speed of s ... the walk takes her 4 hours, including t min...

S1 Error at Step: 2

FlexiVe Feedback (Excerpt): ...Understanding the problem: Aya walks 9 km at two different speeds... We need to find the total time when she walks at $(s + \frac{1}{2})$ km/h.

Setting up equations:

- First scenario: $4 = \frac{9}{s} + \frac{t}{60}$
- Second scenario: $2.4 = \frac{9}{s+2} + \frac{t}{60}$

Subtracting equations: ...

Result: S2 was correct

Ineffective Feedback (Failed Correction)

Problem: Let B be the set of rectangular boxes with surface area 54 and volume 23. Let r be the radius of the smallest sphere that can contain...

Error Location by FlexiVe: Step 13

FlexiVe Feedback (Excerpt): The solution starts by understanding that the radius... is half the space diagonal... $r^2 = \frac{a^2+b^2+c^2}{4}$. The goal is to maximize $a^2 + b^2 + c^2$ given the constraints...

Outcome: Correction failed

Analysis: For complex geometry problems, FlexiVe may fail to produce a corrective pathway, highlighting a limitation in advanced spatial and geometric problem-solving capabilities.