

# Enhancing Large Language Model Reasoning via Latent Space Modeling for Pre-thinking and Pre-answering

Anonymous ACL submission

## Abstract

Chain-of-Thought (CoT) has proven effective for eliciting reasoning abilities in large language models, but standard autoregressive decoders treat reasoning and answer tokens uniformly as next-token prediction targets, despite their fundamentally different roles. This homogeneous modeling introduces conflicting inductive biases and significant inference overhead. Latent CoT methods attempt to improve efficiency by removing explicit reasoning tokens, yet often degrade performance, indicating that explicit reasoning supervision remains critical for preserving pretrained reasoning capabilities. To address these challenges, we propose **Pre<sup>2</sup>**, a phase-aware reasoning framework that decouples thinking and answering while retaining explicit CoT supervision. Our approach introduces independent pre-thinking and pre-answering latent states and integrates them with explicit token representations through a hybrid latent-token decoding paradigm. This design increases computation along both width and depth dimensions with minimal parameter overhead, enabling stronger reasoning without sacrificing inference efficiency. Experiments on one in-domain dataset and four out-of-domain benchmarks demonstrate consistent improvements over strong baselines and robust generalization. Code will be released at [github.com](https://github.com)

## 1 Introduction

Chain-of-Thought (CoT) prompting, introduced by [Wei et al. \(2022\)](#), has proven effective in enhancing the reasoning abilities of large language models by inducing intermediate reasoning steps, yielding strong gains on mathematical and multi-step reasoning tasks ([Guo et al., 2025](#); [OpenAI, 2024](#); [Muennighoff et al., 2025](#)). However, these prompt-based methods rely on explicitly generating reasoning tokens at inference time ([Wang et al., 2022](#); [Zhou et al., 2022](#); [Kojima et al., 2022](#); [Zhang et al., 2022](#)). Due to the autoregressive nature of language models, early reasoning errors can propagate

through subsequent predictions, often resulting in brittle and unstable reasoning trajectories.

To mitigate this limitation, a line of work has explored test-time scaling strategies that aggregate or refine reasoning paths through increased computation at inference. Representative approaches include self-consistency ([Wang et al., 2022](#)), which samples multiple reasoning trajectories and selects the most consistent answer, as well as more structured methods such as Tree-of-Thoughts (ToT, [Yao et al. 2023](#)), which explicitly searches over reasoning branches. While effective, these methods significantly increase inference cost, as they require extensive sampling, branching, or search procedures, making them impractical for latency- or budget-constrained settings ([Li et al., 2024](#); [Wang et al., 2025](#)).

More recently, latent CoT approaches have emerged as an alternative direction, aiming to bypass explicit reasoning token generation altogether ([Wei et al., 2025](#); [He et al., 2025](#); [Liu and Zhang, 2025](#)). Instead of producing natural language reasoning steps, these methods perform reasoning directly in a continuous latent space, thereby reducing inference-time token generation. Representative works such as Coconut ([Hao et al., 2024](#)), CoDi ([Shen et al., 2025](#)), and SoftCoT ([Xu et al., 2025](#)) demonstrate that latent representations can capture certain aspects of reasoning while improving efficiency. Nevertheless, by skipping or discarding explicit reasoning tokens, these approaches often struggle to fully preserve the reasoning capabilities of pretrained models, highlighting a fundamental trade-off between efficiency and reasoning fidelity ([Xia et al., 2025](#); [Liu et al., 2025](#)).

Traditional methods uniformly model reasoning and answer tokens, ignoring their fundamentally different patterns ([Nye et al., 2021](#); [Lampinen et al., 2022](#)). While latent reasoning approaches eliminate CoT for efficiency, we argue that CoT provides essential supervision signals aligned with pretrain-

ing. Fully eliminating CoT may weaken reasoning ability (Turpin et al., 2023). To address these limitations, we explicitly decouple the modeling of reasoning and answering by introducing independent pre-thinking and pre-answering phases, and our approach preserves explicit reasoning supervision while introducing structured latent phases, balancing inference efficiency and reasoning expressivity. Specifically, instead of forcing the model to autoregressively commit to explicit reasoning tokens, we supervise continuous latent representations that summarize the reasoning-ready and answer-ready states. This separation enables the model to first internalize a structured reasoning abstraction before transitioning into answer generation, thereby strengthening its reasoning capability while reducing reliance on verbose or brittle CoT generation.

We introduce **Pre<sup>2</sup>**, a novel framework that effectively addresses all of the above-mentioned problems simultaneously. **We model Chain-of-Thought reasoning as a structured latent trajectory composed of a reasoning-ready state and an answer-ready state. Instead of autoregressively generating reasoning tokens, we supervise continuous latent representations aligned with pooled reasoning and answer embeddings, enabling non-autoregressive and controllable reasoning abstraction.** Our contributions are summarized as follows: (1) We propose a phase-aware reasoning framework that independently models the pre-thinking and pre-answering processes, effectively decoupling the conflicting patterns between thinking and answering stages. (2) We introduce a hybrid paradigm that combines latent representations with explicit token-level supervision, requiring only a minimal increase in model parameters while preserving strong performance and high inference efficiency. (3) We conduct extensive experiments on one in-domain dataset and four out-of-domain benchmarks, where our method consistently outperforms strong baselines, demonstrating robust generalization across domains.

## 2 Related Works

**Explicit Chain-of-Thought Reasoning.** Explicit Chain-of-Thought (CoT) prompting has become a cornerstone for eliciting reasoning capabilities in large language models (LLMs). Pioneering works demonstrate that explicitly generating intermediate reasoning steps significantly boosts performance on arithmetic and symbolic tasks (Wei et al., 2022;

Kojima et al., 2022). Despite the successes, explicit CoT suffers from inherent inefficiency due to the autoregressive generation of long reasoning sequences. Furthermore, treating reasoning tokens and answer tokens as homogeneous prediction targets introduces structural conflicts; reasoning trajectories often require exploratory and verbose patterns, whereas answers demand conciseness (Nye et al., 2021; Lampinen et al., 2022).

**Latent Chain-of-Thought Reasoning.** To mitigate the inference overhead of explicit tokens, Latent CoT approaches aim to internalize reasoning into continuous hidden states or implicit processes, enabling models to “think before speaking” (Zelikman et al., 2024). Methods like implicit reasoning via stepwise knowledge internalization (Deng et al., 2024) and teacher-student distillation (e.g., CODI, Shen et al. 2025) attempt to transfer reasoning semantics from explicit traces to model embeddings. Other approaches, such as Coconut (Hao et al., 2024) and CCoT (Cheng and Van Durme, 2024), explore reasoning in continuous space by compressing steps into dense vectors or “contemplation tokens”. However, purely latent methods face significant challenges: they often suffer from performance degradation on complex tasks due to the lack of fine-grained supervision signals (Hao et al., 2024). Unlike explicit CoT, where every step is verified, methods like Coconut typically rely on answer-level supervision, while distillation methods like CODI focus on trajectory-level alignment, potentially losing critical step-wise logic.

**Computational Scaling across Width and Depth Dimensions.** Existing computational scaling methods face limitations: width-based approaches (e.g., Pause Tokens, Goyal et al. 2023) often lack explicit semantic guidance, while depth-based methods (e.g., MoR, Bae et al. 2025 and Ouro Zhu et al. 2025) keep reasoning implicit within recursive loops. We bridge these dimensions by scaling width via learnable latent tokens and depth through independent modules.

Our approach leverages a hybrid framework combining explicit and latent Chain-of-Thought reasoning. Unlike pause tokens that operate exclusively with static token-level embeddings, our framework employs learnable sentence-level embeddings to capture sequence representations of complete reasoning steps and answer tokens.

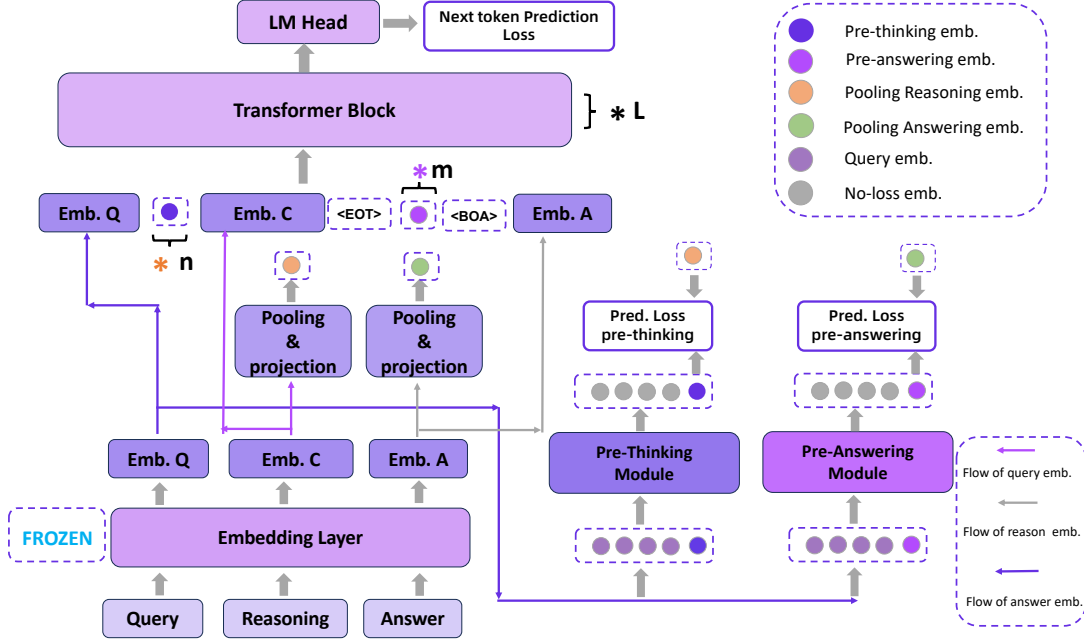


Figure 1: Overview of the proposed framework.

### 3 Methodology

#### 3.1 Preliminary of Pre<sup>2</sup>

Let  $\mathcal{V}$  denote a vocabulary of size  $V = |\mathcal{V}|$ . Given a single training sample, we define a CoT sequence as

$$\mathbf{X} = (\mathbf{Q} := x^{(q)}, \mathbf{R} := x^{(r)}, \mathbf{A} := x^{(a)}),$$

where  $x^{(q)} = (q_1, \dots, q_{L_q}) \in \mathcal{V}^{L_q}$ ,  $x^{(r)} = (r_1, \dots, r_{L_r}) \in \mathcal{V}^{L_r}$ , and  $x^{(a)} = (a_1, \dots, a_{L_a}) \in \mathcal{V}^{L_a}$  denote the query, reasoning steps, and answer, respectively.

Our framework introduces a novel approach where three built-in embeddings serve as initialization for the latent states of pre-thinking and pre-answering. We first introduce three continuous latent vectors that have been initialized within our model:

$$\mathbf{z}_{\text{think}}^{(0)} \in \mathbb{R}^d, \mathbf{z}_{\text{ans}}^{(0)} \in \mathbb{R}^d, \mathbf{z}_{\text{boa}}^{(0)} \in \mathbb{R}^d,$$

where  $\mathbf{z}_{\text{think}}^{(0)}$ ,  $\mathbf{z}_{\text{ans}}^{(0)}$  and  $\mathbf{z}_{\text{boa}}^{(0)}$  are learnable global parameters. Specifically,  $\mathbf{z}_{\text{think}}^{(0)}$  is trained to learn a latent representation of the complete reasoning process;  $\mathbf{z}_{\text{ans}}^{(0)}$  captures an abstract representation of the answer; and  $\mathbf{z}_{\text{boa}}^{(0)}$  serves as an answer indicator that signals the model to transition into answer generation mode.

Our model uses a frozen token embedding matrix  $\mathcal{E}$  to get different non-contextualized embeddings

for query, reasoning and answer tokens. This approach offers an additional advantage: the frozen embedding layer provides stable training targets while preserving model capacity, thus avoiding performance degradation. Let  $\mathcal{E} : \mathcal{V} \rightarrow \mathbb{R}^d$  be a shared token embedding function. The embedded sequences are given by:

$$\begin{aligned} \mathbf{Q} &= (\mathbf{q}_1, \dots, \mathbf{q}_{L_q}) \in \mathbb{R}^{L_q \times d}, \\ \mathbf{R} &= (\mathbf{r}_1, \dots, \mathbf{r}_{L_r}) \in \mathbb{R}^{L_r \times d}, \\ \mathbf{A} &= (\mathbf{a}_1, \dots, \mathbf{a}_{L_a}) \in \mathbb{R}^{L_a \times d}. \end{aligned}$$

We then define a pooling operator  $\mathcal{P}(\cdot)$  to obtain sequence-level representations. And the pooled reasoning and answers representations will be regarded as gold latent representations because we assume that these pooled representations include all necessary information to guide the model to generate complete reasoning tokens and answering tokens.

$$\begin{aligned} \bar{\mathbf{r}} &= \mathcal{P}(\mathbf{R}) \in \mathbb{R}^d, \\ \bar{\mathbf{a}} &= \mathcal{P}(\mathbf{A}) \in \mathbb{R}^d. \end{aligned}$$

#### 3.2 Pre-thinking and Pre-answering Modeling

The modeling of pre-thinking and pre-answering is centered on a key principle: we aim to obtain embeddings that encapsulate holistic information. Specifically: The **pre-thinking embeddings** are expected to capture complete information about the reasoning steps, encoding the entire reasoning

trajectory in a compact representation. The **pre-answering embeddings** should encapsulate the final answer information. With these global representations, the model can provide top-down guidance and predictions throughout the subsequent autoregressive generation of reasoning tokens and answer tokens, rather than relying solely on local, step-by-step predictions. We now address the key question: how are these compressed embeddings learned?

We make a key empirical observation: for a small proportion of complex reasoning tasks, models can produce correct answers directly from the query even when the generated reasoning steps are unreliable or incorrect. This finding suggests that *the query itself encapsulates sufficient information for answer prediction, independent of the reasoning quality*. Motivated by this insight, we design our approach to leverage query tokens as conditions, using the built-in embeddings as latent state initializations and employing two lightweight Transformer decoder blocks to causally model token-level dependencies. Moreover, recent studies have demonstrated that increasing per-token computation depth yields unique benefits that cannot be achieved by simply scaling the width of sequence modeling, particularly for certain reasoning tasks (Bae et al., 2025; Zhu et al., 2025). To enhance computational and inference efficiency, we selectively allocate additional depth-wise computation exclusively to reasoning and answer tokens.

Let  $f_{\text{think}}$  denote the pre-thinking module. The input sequence is constructed as

$$\mathbf{H}_{\text{think}}^{(0)} = [\mathbf{Q}; \mathbf{z}_{\text{think}}^{(0)}],$$

where  $[\cdot; \cdot]$  is concatenation operator.

The hidden states are obtained by

$$\mathbf{H}_{\text{think}} = f_{\text{think}}(\mathbf{H}_{\text{think}}^{(0)}).$$

We take the final hidden state as the predicted pre-thinking latent representation:

$$\hat{\mathbf{z}}_{\text{think}} = \mathbf{H}_{\text{think}}[-1].$$

Similarly, we define a pre-answering module  $f_{\text{ans}}$  with input

$$\mathbf{H}_{\text{ans}}^{(0)} = [\mathbf{Q}; \mathbf{z}_{\text{ans}}^{(0)}].$$

The predicted pre-answering latent is given by

$$\hat{\mathbf{z}}_{\text{ans}} = f_{\text{ans}}(\mathbf{H}_{\text{ans}}^{(0)})[-1].$$

Specifically, the two modules consist of two Transformer blocks and utilize a causal mask for attention computation. This design enables both  $\mathbf{z}_{\text{think}}^{(0)}$  and  $\mathbf{z}_{\text{ans}}^{(0)}$  to attend to the complete set of query token embeddings, allowing modeling pre-thinking and pre-answering separately. To enable  $\mathbf{z}_{\text{think}}^{(0)}$  and  $\mathbf{z}_{\text{ans}}^{(0)}$  to learn representations that can predict the overall architecture encompassing the thinking and answering stages, we let  $\bar{\mathbf{r}}$  and  $\bar{\mathbf{a}}$  be the supervised signal for pre-thinking modeling and pre-answering modeling, and the loss function will be either smooth L1 or L2 loss. To prevent gradient backpropagation and ensure unidirectional learning flow, we apply a stop-gradient operation ( $\mathcal{SG}(\cdot)$ ) on  $\bar{\mathbf{r}}$  and  $\bar{\mathbf{a}}$ . Thus, the pre-thinking alignment loss (e.g., L2 norm) is defined as

$$\mathcal{L}_{\text{think}} = \|\hat{\mathbf{z}}_{\text{think}} - \mathcal{SG}(\bar{\mathbf{r}})\|_2^2. \quad (1)$$

The pre-answering alignment loss is defined as

$$\mathcal{L}_{\text{ans}} = \|\hat{\mathbf{z}}_{\text{ans}} - \mathcal{SG}(\bar{\mathbf{a}})\|_2^2. \quad (2)$$

### 3.3 Training of Pre<sup>2</sup>

In the following section, we describe the training procedure. Let  $\mathbf{e}_{\text{eot}} \in \mathbb{R}^d$  denote the end-of-thought embedding. The decoder input sequence is constructed as

$$\mathbf{H}_{\text{dec}}^{(0)} = [\mathbf{Q}; \hat{\mathbf{z}}_{\text{think}}; \mathbf{R}; \mathbf{e}_{\text{eot}}; \hat{\mathbf{z}}_{\text{ans}}; \mathbf{z}_{\text{boa}}; \mathbf{A}] \in \mathbb{R}^{T_{\text{dec}} \times d},$$

where  $T_{\text{dec}} = L_q + 1 + L_r + 1 + 1 + 1 + L_a$ .

And the  $\mathbf{e}_{\text{eot}}$  is a special token that is already in the vocabulary of the pre-trained backbone model. And this  $\mathbf{e}_{\text{eot}}$  is essential for our inference pipeline because it can provide a proper signal that tells us the reasoning stage is completed.

The decoder  $f_{\text{dec}}$  produces hidden states

$$\mathbf{H}_{\text{dec}} = f_{\text{dec}}(\mathbf{H}_{\text{dec}}^{(0)}). \quad (3)$$

The corresponding supervision labels are defined as

$$\mathbf{Y} = [\underbrace{-100}_{\mathbf{Q}}, \mathbf{R}, \mathbf{e}_{\text{eot}}, \underbrace{-100}_{\hat{\mathbf{z}}_{\text{ans}}}, \underbrace{-100}_{\mathbf{z}_{\text{boa}}}, \mathbf{A}], \quad (4)$$

where  $-100$  indicates masked positions that do not contribute to the loss.

Let  $\mathcal{T}_{\text{sup}} \subset \{1, \dots, T_{\text{dec}}\}$  denote the set of supervised token positions, given by

$$\mathcal{T}_{\text{sup}} = \mathcal{T}_{\text{cot}} \cup \mathcal{T}_{\text{ans}}, \quad (5)$$

where

$$\mathcal{T}_{\text{cot}} = \{t \mid x_t \in \mathbf{R} \cup \{\mathbf{e}_{\text{cot}}\}\}, \mathcal{T}_{\text{ans}} = \{t \mid x_t \in \mathbf{A}\}.$$

The masked language modeling loss is defined as

$$\begin{aligned} \mathcal{L}_{\text{LM}} &= - \sum_{t \in \mathcal{T}_{\text{sup}}} \log p_{\theta}(y_t \mid \mathbf{H}_{\text{dec}}^{(<t)}) \\ &= - \sum_{t \in \mathcal{T}_{\text{cot}}} \log p_{\theta}(r_t \mid \mathbf{H}_{\text{dec}}^{(<t)}) \\ &\quad - \sum_{t \in \mathcal{T}_{\text{ans}}} \log p_{\theta}(a_t \mid \mathbf{H}_{\text{dec}}^{(<t)}). \end{aligned} \quad (3)$$

The overall training objective is given by

$$\mathcal{L} = (1 - \lambda_{\text{think}} - \lambda_{\text{ans}}) \mathcal{L}_{\text{LM}} + \lambda_{\text{think}} \mathcal{L}_{\text{think}} + \lambda_{\text{ans}} \mathcal{L}_{\text{ans}}. \quad (4)$$

### 3.4 Inference of Pre<sup>2</sup>

The inference procedure differs from vanilla autoregressive decoding in several ways. Notably, we must explicitly maintain: (1) the current query length, and (2) an End-of-Thought (EoT) signal to determine when CoT generation concludes and answer generation should begin. We provide a detailed pseudocode description in Algorithm 1, where  $\mathbf{q}$  is query embeddings,  $\hat{\mathbf{z}}_t$  is predicted pre-thinking embedding,  $\hat{\mathbf{z}}_a$  is predicted pre-answering embedding,  $\mathbf{z}_{\text{boa}}$  is the answering indicator embedding that initialized in the model. Specifically, during the prefilling phase, we execute the pre-thinking  $f_{\text{think}}(\cdot)$  and pre-answering modules  $f_{\text{answer}}(\cdot)$  to generate  $\hat{\mathbf{z}}_t$  and  $\hat{\mathbf{z}}_a$  before any token generation begins. In the autoregressive generation stages that follow, we dynamically reconstruct the input sequence  $\mathbf{X}$  through concatenation operations, ensuring that the model’s behavior remains consistent across training and inference.

## 4 Experimental Design

### 4.1 Datasets

We train the Pre<sup>2</sup> framework on **GSM8K-Aug** (Deng et al., 2024), a large-scale corpus of approximately 385,000 samples, which augments the original GSM8K (Cobbe et al., 2021) training set with questions, step-by-step reasoning chains, and final solutions to provide granular supervision for aligning our pre-thinking and pre-answering modules. For evaluation, we assess reasoning robustness and generalization across five benchmarks categorized into in-domain and out-of-domain settings. The in-domain evaluation uses the standard **GSM8K**

---

### Algorithm 1 Inference Process of Pre<sup>2</sup>

---

**Require:** Token embeddings  $\mathbf{x}$  from frozen embedding layer;

**Require:** Stage-flags:

is\_prefilling, is\_cot\_end, is\_answering,

**Ensure:** Updated embeddings  $\mathbf{x}$

1: **if** is\_prefilling **then**

2:    $\hat{\mathbf{z}}_t = f_{\text{think}}([x; \mathbf{z}_{\text{think}}^{(0)}])[-1]$

3:    $\hat{\mathbf{z}}_a = f_{\text{answer}}([x; \mathbf{z}_{\text{answer}}^{(0)}])[-1]$

4:    $\mathbf{x} \leftarrow [x; \hat{\mathbf{z}}_t]$

5: **else**

6:   **if** is\_cot\_end **then**

7:     **if**  $\neg$ is\_answering **then**

8:        $\mathbf{x} \leftarrow [\mathbf{q}; \hat{\mathbf{z}}_t; \text{CoT}; \langle \text{eot} \rangle; \hat{\mathbf{z}}_a; \mathbf{z}_{\text{boa}}]$

9:       is\_answering  $\leftarrow$  True

10:     **else**

11:        $\mathbf{x} \leftarrow [\mathbf{q}; \hat{\mathbf{z}}_t; \text{CoT}; \langle \text{eot} \rangle; \hat{\mathbf{z}}_a; \mathbf{z}_{\text{boa}}; \text{Ans}]$

12:     **end if**

13:     **else**

14:        $\mathbf{x} \leftarrow [\mathbf{q}; \hat{\mathbf{z}}_t; \text{CoT}]$

15:     **end if**

16: **end if**

17: **return**  $\mathbf{x}$

---

(Cobbe et al., 2021) test split, while out-of-domain evaluation employs four datasets: **GSM-Hard** (increased arithmetic difficulty through larger integers) (Gao et al., 2022), **MultiArith** (multi-step sequential reasoning) (Roy and Roth, 2015), **SVAMP** (structural and textual variations) (Patel et al., 2021), and **ASDiv** (diverse problem types and linguistic patterns) (Xu et al., 2025). Full dataset details are provided in the Appendix B.

### 4.2 Baselines

We evaluate Pre<sup>2</sup> against five representative baselines covering different reasoning paradigms: (1) **Zero-shot Pretrained**, which directly evaluates a pretrained model without fine-tuning; (2) **Few-shot In-Context Learning**, where three demonstration examples are provided as prompts without parameter updates; (3) **SFT with Explicit CoT**, which fine-tunes the model to generate the full query–reasoning–answer sequence; (4) **SFT with Frozen Embeddings**, identical to (3) but with the embedding layer fixed to assess the role of input representations; and (5) **Pause Token Baseline** (Goyal et al., 2023), which inserts fixed special tokens around the reasoning sequence following prior “thinking token” approaches, while keeping

Model	Method	Datasets					Overall	
		In-domain		Out-of-domain			Total	Avg
		GSM8K	GSM_H	Multi	ASDiv	SVAMP		
Qwen-0.5B	Zero-shot	43.06	12.15	82.78	77.55	46.33	261.87	52.37
	Few-shot	41.60	10.99	81.83	78.13	45.67	258.22	51.64
	Baseline	51.50	12.74	92.89	82.15	53.22	292.50	58.50
	Freeze-Baseline	51.96	12.36	92.61	82.30	<b>53.89</b>	293.12	58.62
	Pause Token	51.98	12.45	93.38	82.75	52.67	293.23	58.65
	<b>Pre<sup>2</sup>-1</b>	<b>53.78</b>	<b>12.76</b>	<b>94.17</b>	<b>83.33</b>	53.56	<b>297.60</b>	<b>59.52</b>
Qwen-1.5B	Zero-shot	46.90	12.40	90.00	81.31	51.33	281.94	56.38
	Few-shot	46.62	13.16	88.50	81.70	49.67	279.65	55.93
	Baseline	54.28	13.39	94.39	83.46	58.56	304.08	60.82
	Freeze-Baseline	54.01	13.57	94.40	83.49	58.33	303.80	60.76
	Pause Token	54.56	13.57	94.50	83.27	57.78	303.68	60.73
	<b>Pre<sup>2</sup>-1</b>	<b>54.76</b>	<b>14.53</b>	<b>94.50</b>	<b>83.53</b>	<b>59.11</b>	<b>306.43</b>	<b>61.29</b>

Table 1: Main performance comparison on in-domain (GSM8K) and out-of-domain reasoning benchmarks.

these tokens non-learnable during training.

### 4.3 Evaluation Metrics

We evaluate model performance using three metrics. **(1) Accuracy (Acc.)** measures answer correctness, defined as the percentage of correctly predicted answers. **(2) Generated Sequence Length** quantifies reasoning compactness by computing the average number of tokens in the generated reasoning chains across all examples, where shorter sequences indicate higher efficiency. **(3) Inference Throughput** evaluates efficiency by measuring the generation speed in tokens per second, computed as the total number of generated tokens divided by the total inference time. Higher throughput indicates better inference efficiency.

### 4.4 Research Questions

Our experiments aim to answer the following research questions: (1) Compared to traditional decoder-based models, what magnitude of improvement can our architecture achieve on downstream reasoning tasks? (2) More specifically, do pre-thinking modeling and pre-answering modeling enable the model to predict reasoning tokens and answer tokens with higher accuracy? (3) Does scaling up the pre-thinking module and pre-answering module along the depth dimension yield consistent performance gains? (4) Similarly, does scaling up the pre-thinking scope along the width dimension enhance the model’s reasoning capabilities? (5) Compared to the baseline model, we introduce two prediction modules and three internally initialized embeddings; relative to the baseline size, we add only a modest number of parameters, and

our inference mixes autoregressive decoding with latent-space prediction. Will these extra parameters and the new inference paradigm materially impact inference efficiency?

We show the experimental results for the five questions in five subsections of Sec. 5 respectively. Refer to the appendix A for implementation details.

## 5 Results and Discussions

### 5.1 Main Results

Our model consistently outperforms baselines across one in-domain and four out-of-domain benchmarks. Table 1 presents a comprehensive comparison between Qwen2.5-0.5B and Qwen2.5-1.5B models (Qwen et al., 2025) across various reasoning benchmarks. Several key observations emerge from these results: **(1) Consistent superiority of Pre<sup>2</sup>-1:** Our proposed Pre<sup>2</sup>-1 method consistently outperforms all baseline approaches across both model scales. For the 0.5B model, Pre<sup>2</sup>-1 achieves an average accuracy of 59.52%, representing a 1.02 percentage point improvement over the standard Baseline (58.50%). Similarly, for the 1.5B model, Pre<sup>2</sup>-1 attains 61.29% average accuracy, surpassing the Freeze-Baseline by 0.53 percentage points. **(2) Strong generalization to out-of-domain tasks:** Notably, Pre<sup>2</sup>-1 demonstrates robust performance not only on the in-domain GSM8K benchmark but also across all four out-of-domain datasets. For instance, on the 0.5B model, Pre<sup>2</sup>-1 improves Multi-Arith from 98.89% to 94.17%, and ASDIV from 82.15% to 83.33%, while maintaining competitive performance on SVAMP (53.56%). **(3) Comparison with Pause Token baseline:** While the Pause To-

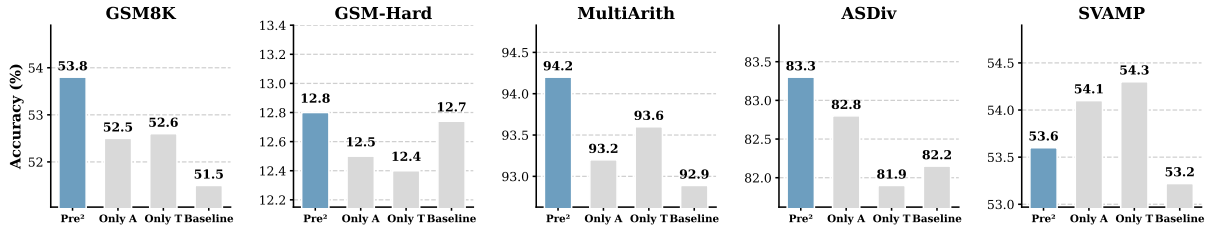


Figure 2: Ablation study on component effectiveness. Avg.acc: Pre<sup>2</sup>(59.52), Only A(59.01), Only T(58.97), Baseline(58.50)

ken method shows marginal improvements over the standard Baseline, it is consistently outperformed by our **Pre<sup>2</sup>-1** approach. This suggests that our pre-thinking and pre-answering mechanisms provide more effective modeling of the reasoning process than simple pause token insertion. **Although both the pause token method and our approach perform equivalent scaling along the sequence width dimension, we hypothesize that this advantage arises because pause tokens employ static embeddings, whereas our pre-thinking and pre-answering embeddings are learnable and dynamic. Our novel training approach further enables these embeddings to encapsulate rich information about future token sequences, providing stronger guidance for reasoning.** (4) **Scaling behavior:** The results indicate that our architectural improvements are complementary to model scaling, with benefits observed across both 0.5B and 1.5B parameter scales, demonstrating the generalizability of our approach.

## 5.2 Impact of Pre-thinking and Pre-answering on Performance

As shown in Figure 2, which presents an ablation study comparing our full model against the baseline and two single-component variants, revealing the individual and synergistic contributions of our architectural components. The full model (**Pre<sup>2</sup>**) achieves 59.52% average accuracy, outperforming the baseline by +1.02 points. Critically, this exceeds both **only\_t** (pre-thinking only) (58.97%, +0.47) and **only\_a** (pre-answering only) (59.01%, +0.51), demonstrating that the combined gain is substantially larger than either component alone, confirming their complementary nature. Besides, both ablated variants independently improve over the baseline, with **only\_a** marginally outperforming **only\_t**. The full model shows particularly strong gains on GSM8K (+2.28 points) and Multi-Arith (+1.28 points). Interestingly, single-

Model	GSM8K	GSM_H	Multi	ASDiv	SVAMP
<b>Pre<sup>2</sup>-1</b>	<b>53.78</b>	<b>12.76</b>	94.17	83.33	53.56
<b>Pre<sup>2</sup>-3</b>	52.31	12.11	<b>95.17</b>	83.30	<b>56.33</b>
<b>Pre<sup>2</sup>-5</b>	53.65	12.53	94.17	83.49	54.56
<b>Pre<sup>2</sup>-7</b>	52.97	12.26	94.89	<b>84.01</b>	55.11

Table 2: Ablation study on the Depth Dimension. Avg.acc: Pre<sup>2</sup>-1(59.52), Pre<sup>2</sup>-3(59.84), Pre<sup>2</sup>-5(59.68), Pre<sup>2</sup>-7(59.85)

component models slightly outperform the full model on SVAMP (54.33%/54.11% vs. 53.56%), suggesting task-specific sensitivities.

## 5.3 Exploring the Depth Scaling of Pre-thinking and Pre-answering Modules

Table 2 investigates the impact of scaling the pre-thinking and pre-answering modules along the depth dimension by varying the number of Transformer layers (1, 3, 5, 7). Increasing module depth consistently improves overall performance, with **Pre<sup>2</sup>-7** achieving the highest average accuracy of 59.85%, followed closely by **Pre<sup>2</sup>-3** (59.84%). This represents a +0.33 point improvement over the shallowest **Pre<sup>2</sup>-1** variant (59.52%), demonstrating that allocating additional depth-wise computation to the reasoning process yields meaningful gains. While deeper variants improve average performance, the benefits are not uniform across all datasets. **Pre<sup>2</sup>-1** excels on GSM8K (53.78%) and GSM\_hard (12.76%), whereas **Pre<sup>2</sup>-3** achieves the best results on Multi-Arith (95.17%) and SVAMP (56.33%). **Pre<sup>2</sup>-7** demonstrates the strongest performance on ASDiv (84.01%), suggesting that different reasoning tasks may benefit from varying computational depths.

## 5.4 On the Width Scaling of Pre-thinking Scope

Figure 3 examines the impact of scaling the pre-thinking scope (here we fix the depth at 1 transformer layer) (width dimension) from 1 to 5 to-

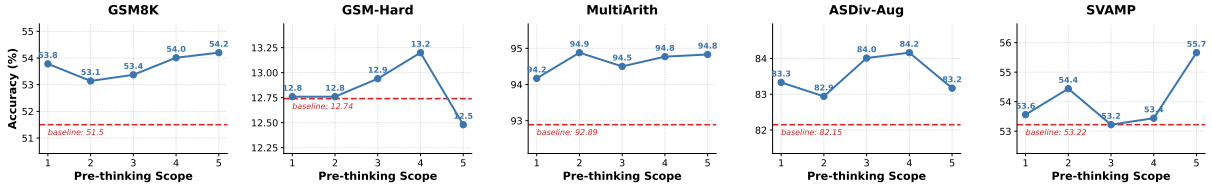


Figure 3: Ablation study on the Width Dimension (Pre-thinking Scope) Avg.acc: Scope-1(59.52), Scope-2(59.63), Scope-3(59.61), Scope-4(59.91), Scope-5(60.06), Baseline(58.50).

Model	GSM8K		GSMHard		MultiArith		SVAMP		ASDiv-Aug	
	Tok	T/s	Tok	T/s	Tok	T/s	Tok	T/s	Tok	T/s
<b>Pre<sup>2</sup>-0.5b</b>	47.6	76.7 (1.00)	78.9	76.6 (0.36)	23.6	77.2 (1.31)	25.9	76.5 (1.07)	15.6	75.9 (1.45)
Pause Token	44.5	76.8 (0.96)	77.2	76.5 (0.78)	24.6	77.1 (0.55)	24.7	76.6 (1.00)	15.6	75.1 (1.52)
<b>Pre<sup>2</sup>-0.5b-1-1</b>	46.1	76.9 (0.49)	78.6	76.5 (0.85)	24.6	76.3 (1.61)	24.4	76.3 (1.57)	15.6	75.4 (2.36)
<b>Pre<sup>2</sup>-0.5b-3-1</b>	45.1	75.0 (1.55)	71.3	75.8 (0.81)	24.6	75.3 (1.99)	23.3	75.7 (0.90)	15.6	77.0 (0.92)
<b>Pre<sup>2</sup>-0.5b-5-1</b>	44.4	74.5 (0.40)	74.9	74.4 (0.43)	24.6	74.6 (0.99)	28.5	74.3 (1.93)	15.6	75.3 (0.87)
<b>Pre<sup>2</sup>-0.5b-7-1</b>	46.2	74.1 (1.59)	72.3	73.6 (1.37)	24.6	74.5 (1.29)	28.5	73.8 (2.19)	15.6	73.1 (0.88)
<b>Pre<sup>2</sup>-0.5b-1-2</b>	45.0	75.8 (1.31)	77.8	76.1 (1.24)	24.6	75.7 (1.16)	26.3	76.0 (1.04)	15.6	74.2 (1.68)
<b>Pre<sup>2</sup>-0.5b-1-3</b>	45.3	75.3 (1.46)	75.5	76.5 (0.88)	24.6	75.6 (1.38)	25.8	75.6 (1.23)	16.4	75.1 (0.79)
<b>Pre<sup>2</sup>-0.5b-1-4</b>	45.3	75.1 (1.63)	71.0	75.2 (0.76)	24.6	75.3 (0.41)	26.4	75.4 (1.38)	15.6	75.3 (2.04)
<b>Pre<sup>2</sup>-0.5b-1-5</b>	49.0	74.6 (1.04)	73.8	74.5 (0.71)	24.6	75.1 (1.31)	25.7	75.6 (0.95)	15.6	73.8 (1.69)

Table 3: Inference efficiency analysis comparing **Pre<sup>2</sup>** variants against the Pause Token baseline.

522 kens, which means that we will expand the width  
523 of pre-thinking scope during training and inference  
524 time. Note that we fix the pre-answering  
525 scope at 1, as all benchmark answers are suffi-  
526 ciently short to be compressed into a single token  
527 representation. Expanding the pre-thinking scope  
528 demonstrates progressive performance improve-  
529 ments, with **Pre<sup>2</sup>-1-5** achieving the highest aver-  
530 age accuracy of 60.06%, representing a 0.54 point  
531 gain over **Pre<sup>2</sup>-1** and 1.56 points over the baseline  
532 (58.50%). This trend suggests that allocating wider  
533 latent space for reasoning compression enhances  
534 the model’s capacity to capture complex thought  
535 processes. Different benchmarks respond variably  
536 to scope expansion. GSM8K and GSMHard show  
537 steady improvements with wider scope (GSM8K:  
538 53.78%→54.20%, GSMHard: 12.76%→12.48%  
539 with fluctuations). SVAMP exhibits notable gains  
540 at scope-5 (55.66%), suggesting that certain rea-  
541 soning tasks benefit more from expanded latent  
542 capacity. Unlike depth scaling, width expansion  
543 doesn’t follow a strictly monotonic pattern. Scope-  
544 2 through scope-4 show incremental improvements  
545 (59.63%→59.61%→59.92%), with scope-3 show-  
546 ing a slight dip. This suggests potential trade-offs  
547 between compression capacity and optimization  
548 difficulty. These findings indicate that width-wise  
549 scaling of the pre-thinking module provides com-  
550plementary benefits to depth scaling, with wider

scopes enabling richer compression of reasoning  
551 chains. 552

### 5.5 Model Efficiency 553

554 Table 3 presents a comprehensive evaluation of  
555 inference throughput (tokens/second) across all  
556 model variants and benchmarks, demonstrating that  
557 our architectural do not significantly degrade infer-  
558 ence performance compared to the baseline.

## 6 Conclusion 559

560 We presented a phase-aware reasoning framework  
561 that decouples thinking and answering in large  
562 language models by combining latent representa-  
563 tions with explicit chain-of-thought supervision.  
564 By independently modeling pre-thinking and pre-  
565 answering phases, our approach mitigates the limi-  
566 tations of autoregressive reasoning while preserv-  
567 ing pretrained reasoning capabilities. Experiments  
568 across in-domain and out-of-domain benchmarks  
569 demonstrate consistent improvements in both per-  
570 formance and inference efficiency. This work  
571 highlights the importance of structured reasoning  
572 phases as a promising direction for building more  
573 efficient and robust reasoning models.

### Limitations 574

575 Despite its effectiveness, our approach has several  
576 limitations. First, the proposed phase-aware design

577	introduces additional architectural components and training objectives, which slightly increase implementation complexity compared to standard autoregressive decoders. Second, while our method improves inference efficiency relative to explicit chain-of-thought decoding, it still relies on explicit reasoning supervision during training, limiting its applicability to domains where high-quality reasoning annotations are available. Finally, our evaluation focuses on reasoning-centric benchmarks; the impact of phase-aware modeling on general-purpose language generation and non-reasoning tasks remains to be explored.		
578			
579			
580			
581			
582			
583			
584			
585			
586			
587			
588			
589			
590	<b>Ethical Considerations</b>		
591	Our use of publicly available datasets and developed models in this research raises no ethical concerns, and we confirm compliance with the ACL Code of Ethics.		
592			
593			
594			
595	<b>References</b>		
596	Sangmin Bae, Yujin Kim, Reza Bayat, Sungnyun Kim, Jiyouon Ha, Tal Schuster, Adam Fisch, Hrayr Harutyunyan, Ziwei Ji, Aaron Courville, and 1 others. 2025. Mixture-of-recursions: Learning dynamic recursive depths for adaptive token-level computation. <i>arXiv preprint arXiv:2507.10524</i> .		
597			
598			
599			
600			
601			
602	Jeffrey Cheng and Benjamin Van Durme. 2024. Compressed chain of thought: Efficient reasoning through dense representations. <i>arXiv preprint arXiv:2412.13171</i> .		
603			
604			
605			
606	Karl Cobbe, Vineet Kosaraju, Mohammad Bavarian, Mark Chen, Heewoo Jun, Lukasz Kaiser, Matthias Plappert, Jerry Tworek, Jacob Hilton, Reiichiro Nakano, and 1 others. 2021. Training verifiers to solve math word problems. <i>arXiv preprint arXiv:2110.14168</i> .		
607			
608			
609			
610			
611			
612	Yuntian Deng, Yejin Choi, and Stuart Shieber. 2024. From explicit cot to implicit cot: Learning to internalize cot step by step. <i>arXiv preprint arXiv:2405.14838</i> .		
613			
614			
615			
616	Luyu Gao, Aman Madaan, Shuyan Zhou, Uri Alon, Pengfei Liu, Yiming Yang, Jamie Callan, and Graham Neubig. 2022. Pal: Program-aided language models. <i>arXiv preprint arXiv:2211.10435</i> .		
617			
618			
619			
620	Sachin Goyal, Ziwei Ji, Ankit Singh Rawat, Aditya Krishna Menon, Sanjiv Kumar, and Vaishnavh Nagarajan. 2023. Think before you speak: Training language models with pause tokens. <i>arXiv preprint arXiv:2310.02226</i> .		
621			
622			
623			
624			
625	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. <i>arXiv preprint arXiv:2501.12948</i> .	628	629
626			630
627			
	Shibo Hao, Sainbayar Sukhbaatar, DiJia Su, Xian Li, Zhiting Hu, Jason Weston, and Yuandong Tian. 2024. Training large language models to reason in a continuous latent space. <i>arXiv preprint arXiv:2412.06769</i> .	631	632
		633	634
	Yinhan He, Wendy Zheng, Yaochen Zhu, Zaiyi Zheng, Lin Su, Sriram Vasudevan, Qi Guo, Liangjie Hong, and Jundong Li. 2025. Semcot: Accelerating chain-of-thought reasoning through semantically-aligned implicit tokens. <i>arXiv preprint arXiv:2510.24940</i> .	635	636
		637	638
		639	
	Takeshi Kojima, Shixiang Shane Gu, Machel Reid, Yutaka Matsuo, and Yusuke Iwasawa. 2022. Large language models are zero-shot reasoners. <i>Advances in neural information processing systems</i> , 35:22199–22213.	640	641
		642	643
		644	
	Andrew Lampinen, Ishita Dasgupta, Stephanie Chan, Kory Mathewson, Mh Tessler, Antonia Creswell, James McClelland, Jane Wang, and Felix Hill. 2022. Can language models learn from explanations in context? In <i>Findings of the Association for Computational Linguistics: EMNLP 2022</i> , pages 537–563.	645	646
		647	648
		649	650
	Yiwei Li, Peiwen Yuan, Shaoxiong Feng, Boyuan Pan, Xinglin Wang, Bin Sun, Heda Wang, and Kan Li. 2024. Escape sky-high cost: Early-stopping self-consistency for multi-step reasoning. <i>arXiv preprint arXiv:2401.10480</i> .	651	652
		653	654
		655	
	Jingyu Liu and Ce Zhang. 2025. Hamburger: Accelerating llm inference via token smashing. <i>arXiv preprint arXiv:2505.20438</i> .	656	657
		658	
	Yue Liu, Jiaying Wu, Yufei He, Ruihan Gong, Jun Xia, Liang Li, Hongcheng Gao, Hongyu Chen, Baolong Bi, Jiaheng Zhang, and 1 others. 2025. Efficient inference for large reasoning models: A survey. <i>arXiv preprint arXiv:2503.23077</i> .	659	660
		661	662
		663	
	Niklas Muennighoff, Zitong Yang, Weijia Shi, Xiang Lisa Li, Li Fei-Fei, Hannaneh Hajishirzi, Luke Zettlemoyer, Percy Liang, Emmanuel Candès, and Tatsunori B Hashimoto. 2025. s1: Simple test-time scaling. In <i>Proceedings of the 2025 Conference on Empirical Methods in Natural Language Processing</i> , pages 20286–20332.	664	665
		666	667
		668	669
		670	
	Maxwell Nye, Anders Johan Andreassen, Guy Gur-Ari, Henryk Michalewski, Jacob Austin, David Bieber, David Dohan, Aitor Lewkowycz, Maarten Bosma, David Luan, and 1 others. 2021. Show your work: Scratchpads for intermediate computation with language models.	671	672
		673	674
		675	676
	OpenAI. 2024. <i>Gpt-4o</i> . Large language model.	677	
	Arkil Patel, Satwik Bhattamishra, and Navin Goyal. 2021. Are NLP models really able to solve simple math word problems? In <i>Proceedings of the 2021 Conference of the North American Chapter of the Association for Computational Linguistics: Human</i>	678	679
		680	681
		682	

683	<i>Language Technologies</i> , pages 2080–2094, Online.	Shunyu Yao, Dian Yu, Jeffrey Zhao, Izhak Shafran,	737
684	Association for Computational Linguistics.	Tom Griffiths, Yuan Cao, and Karthik Narasimhan.	738
685	Qwen, :, An Yang, Baosong Yang, Beichen Zhang,	2023. Tree of thoughts: Deliberate problem solving	739
686	Binyuan Hui, Bo Zheng, Bowen Yu, Chengyuan	with large language models. <i>Advances in neural</i>	740
687	Li, Dayiheng Liu, Fei Huang, Haoran Wei, Huan	<i>information processing systems</i> , 36:11809–11822.	741
688	Lin, Jian Yang, Jianhong Tu, Jianwei Zhang, Jianxin	Eric Zelikman, Georges Harik, Yijia Shao, Varuna	742
689	Yang, Jiayi Yang, Jingren Zhou, and 25 oth-	Jayasiri, Nick Haber, and Noah D Goodman. 2024.	743
690	ers. 2025. <a href="#">Qwen2.5 technical report</a> . <i>Preprint</i> ,	Quiet-star: Language models can teach them-	744
691	arXiv:2412.15115.	selves to think before speaking. <i>arXiv preprint</i>	745
692	Subhro Roy and Dan Roth. 2015. <a href="#">Solving general arith-</a>	arXiv:2403.09629.	746
693	<a href="#">metic word problems</a> . In <i>Proceedings of the 2015</i>	Zhuosheng Zhang, Aston Zhang, Mu Li, and Alex	747
694	<i>Conference on Empirical Methods in Natural Lan-</i>	Smola. 2022. Automatic chain of thought prompting	748
695	<i>guage Processing</i> , pages 1743–1752, Lisbon, Portu-	in large language models. In <i>The eleventh interna-</i>	749
696	gal. Association for Computational Linguistics.	<i>tional conference on learning representations</i> .	750
697	Zhenyi Shen, Hanqi Yan, Linhai Zhang, Zhanghao Hu,	Denny Zhou, Nathanael Schärli, Le Hou, Jason Wei,	751
698	Yali Du, and Yulan He. 2025. Codi: Compress-	Nathan Scales, Xuezhi Wang, Dale Schuurmans,	752
699	ing chain-of-thought into continuous space via self-	Claire Cui, Olivier Bousquet, Quoc Le, and 1 oth-	753
700	distillation. <i>arXiv preprint arXiv:2502.21074</i> .	ers. 2022. Least-to-most prompting enables complex	754
701	Miles Turpin, Julian Michael, Ethan Perez, and Samuel	reasoning in large language models. <i>arXiv preprint</i>	755
702	Bowman. 2023. Language models don’t always say	arXiv:2205.10625.	756
703	what they think: Unfaithful explanations in chain-of-	Rui-Jie Zhu, Zixuan Wang, Kai Hua, Tianyu Zhang,	757
704	thought prompting. <i>Advances in Neural Information</i>	Ziniu Li, Haoran Que, Boyi Wei, Zixin Wen, Fan Yin,	758
705	<i>Processing Systems</i> , 36:74952–74965.	He Xing, Lu Li, Jiajun Shi, Kaijing Ma, Shanda Li,	759
706	Xuezhi Wang, Jason Wei, Dale Schuurmans, Quoc Le,	Taylor Kergan, Andrew Smith, Xingwei Qu, Mude	760
707	Ed Chi, Sharan Narang, Aakanksha Chowdhery, and	Hui, Bohong Wu, and 14 others. 2025. <a href="#">Scaling latent</a>	761
708	Denny Zhou. 2022. Self-consistency improves chain	<a href="#">reasoning via looped language models</a> . <i>Preprint</i> ,	762
709	of thought reasoning in language models. <i>arXiv</i>	arXiv:2510.25741.	763
710	<i>preprint arXiv:2203.11171</i> .	<b>A Experiment Details</b>	764
711	Yiming Wang, Pei Zhang, Siyuan Huang, Baosong	We provide the full hyperparameter settings, train-	765
712	Yang, Zhuosheng Zhang, Fei Huang, and Rui Wang.	ing procedures, and additional analysis used in our	766
713	2025. Sampling-efficient test-time scaling: Self-	experiments.	767
714	estimating the best-of-n sampling in early decoding.	<b>A.1 Implementation Details</b>	768
715	<i>arXiv preprint arXiv:2503.01422</i> .	Our framework, <b>Pre</b> <sup>2</sup> , was implemented using Py-	769
716	Jason Wei, Xuezhi Wang, Dale Schuurmans, Maarten	Torch and the Litgpt framework. We selected the	770
717	Bosma, Fei Xia, Ed Chi, Quoc V Le, Denny Zhou,	<b>Qwen2.5</b> series (Qwen et al., 2025) (0.5B and	771
718	and 1 others. 2022. Chain-of-thought prompting elic-	1.5B) as our backbone models due to their strong	772
719	its reasoning in large language models. <i>Advances</i>	mathematical reasoning capabilities.	773
720	<i>in neural information processing systems</i> , 35:24824–	<b>Module Initialization</b> As described in Section 3,	774
721	24837.	we freeze the original embedding layer of the back-	775
722	Xilin Wei, Xiaoran Liu, Yuhang Zang, Xiaoyi Dong,	bone model to maintain a stable semantic space.	776
723	Yuhang Cao, Jiaqi Wang, Xipeng Qiu, and Dahua	The pre-thinking ( $f_{\text{think}}$ ) and pre-answering ( $f_{\text{ans}}$ )	777
724	Lin. 2025. Sim-cot: Supervised implicit chain-of-	modules are initialized as lightweight Transformer	778
725	thought. <i>arXiv preprint arXiv:2509.20317</i> .	decoder blocks. Specifically, their hidden size and	779
726	Heming Xia, Chak Tou Leong, Wenjie Wang, Yongqi	number of attention heads align with the backbone	780
727	Li, and Wenjie Li. 2025. Tokenskip: Controllable	model to ensure compatibility with the frozen em-	781
728	chain-of-thought compression in llms. <i>arXiv preprint</i>	beddings. For the latent vectors $\mathbf{z}_{\text{think}}^{(0)}$ and $\mathbf{z}_{\text{ans}}^{(0)}$ ,	782
729	<i>arXiv:2502.12067</i> .	we use a Gaussian initialization with mean 0 and	783
730	Yige Xu, Xu Guo, Zhiwei Zeng, and Chunyan Miao.	standard deviation 0.02.	784
731	2025. <a href="#">SoftCoT: Soft chain-of-thought for efficient</a>	<b>Training Setup</b> All models were trained on a	785
732	<a href="#">reasoning with LLMs</a> . In <i>Proceedings of the 63rd</i>	cluster of 8 NVIDIA A100 (80GB) GPUs using	786
733	<i>Annual Meeting of the Association for Computational</i>	the Distributed Data Parallel (DDP) strategy. The	787
734	<i>Linguistics (Volume 1: Long Papers)</i> , pages 23336–		
735	23351, Vienna, Austria. Association for Computa-		
736	tional Linguistics.		

training process utilized the AdamW optimizer. To stabilize the learning of the newly introduced latent modules alongside the backbone, we applied a linear warmup strategy followed by a cosine decay scheduler. For the baselines (Freeze-Baseline and Pause Token), we strictly followed the same training configurations (batch size, learning rate, and scheduler) to ensure a fair comparison. For all evaluations, we use greedy decoding to ensure deterministic generation.

## A.2 Hyperparameters

Table 4 lists the detailed hyperparameters used for the main experiments. We determined the loss weights  $\lambda_{\text{think}}$  and  $\lambda_{\text{ans}}$  through a preliminary grid search on a hold-out validation set from GSM8K-Aug. We found that setting the weights to 1.0 provided a balanced gradient flow between the alignment objectives ( $\mathcal{L}_{\text{think}}$ ,  $\mathcal{L}_{\text{ans}}$ ) and the autoregressive prediction objective ( $\mathcal{L}_{\text{LM}}$ ).

## B Dataset Details

In this section, we provide detailed specifications of the datasets used in our experiments. We first present the statistical overview and then detail the specific characteristics and examples for each benchmark.

### B.1 Dataset Overview and Metadata

Table 5 summarizes the general metadata for each dataset. All selected datasets focus on the *Math* domain and require *Arithmetic* reasoning capabilities. The answer types are consistently numerical values.

### B.2 Dataset Statistics

Table 6 presents the statistical details regarding the length of inputs and reasoning paths. We report the average length of questions (tokens), the average length of the Chain-of-Thought (CoT) rationale, and the average number of reasoning steps.

### B.3 Benchmark Detail

In this section, we describe the specific characteristics of each evaluation benchmark and the training data.

**GSM8K-Aug (Deng et al., 2024)** GSM8K-Aug is an augmented dataset derived from GSM8K (Cobbe et al., 2021). It expands the original 8.5k training problems to roughly 385k examples through paraphrasing, numerical resampling, and

synthetic generation using GPT-4. The dataset exhibits a distribution of reasoning steps where the majority of problems require two to four steps, accompanied by a long tail of instances requiring six or more steps. This balance of common and complex instances makes the dataset suitable for training models to generalize across various reasoning difficulty levels.

**GSM-Hard (Gao et al., 2022)** GSM-Hard is a challenging variant of the GSM8K dataset constructed to test the ability of models to cope with harder numerical values. It retains the same problem statements as the original GSM8K but replaces many of the numbers with larger and less common numerical values, thereby increasing the difficulty of superficial arithmetic computation and reasoning. The dataset consists of approximately 1,319 examples.

**MultiArith (Roy and Roth, 2015)** MultiArith is a dataset of multi-step arithmetic word problems designed to challenge systems in correctly sequencing multiple operations. It contains 600 problems collected from educational sources, each solvable by a single equation involving two or more of the four basic operations (addition, subtraction, multiplication, division). The problems require reasoning over multiple sentences to extract and combine numeric quantities, understand implied operations, and compute the result. It is widely used as a benchmark for evaluating arithmetic reasoning generalization, particularly for handling compositional and sequential numerical reasoning.

**SVAMP (Patel et al., 2021)** SVAMP (Simple Variations on Arithmetic Math Word Problems) is a benchmark dataset designed to test the robustness of math word problem solvers to superficial changes. It contains 1,000 elementary-level arithmetic word problems (grade 4 and below), each involving a single unknown and solvable by an arithmetic expression with no more than two operators. The problems are generated through controlled variations in wording, structure, and number values applied to existing datasets (such as MAWPS and ASDiv-A) to reduce artifacts and superficial cues.

**ASDiv-AUG (Xu et al., 2025)** ASDiv (Academia Sinica Diverse MWP Dataset) is a diverse corpus designed to evaluate the capability of solvers across a wide range of language patterns and problem types. The original dataset contains approximately 2,305 math word problems, each annotated with its

Hyperparameter	<i>0.5b</i>										<i>1.5b</i>		
	1-1	1-2	1-3	1-4	1-5	3-1	5-1	7-1	pt	sft	1-1	pt	sft
Pre-thinking Scope	1	2	3	4	5	5	5	5	1	-	1	1	-
Thinking Layers	1	1	1	1	1	3	5	7	-	-	1	-	-
Answering Layers	1	1	1	1	1	3	5	7	-	-	1	-	-
Epochs					50						11	-	-
Micro Batch Size					32						16	-	-
Max Gradient Norm								1.0					
Optimizer								AdamW					
Learning Rate								2.0e-4					
Min LR								6.0e-5					
LR Warmup Steps								2000					
Precision								bf16-mixed					
Weight Decay								0.01					
Global Batch Size								256					
$\beta_1$								0.9					
$\beta_2$								0.999					
$\lambda^{\text{think}}$								0.15					
$\lambda_{\text{ans}}$								0.15					

Table 4: Hyperparameters for all model variants. (For the format 1-1, the first 1 means we use 1-layer for Pre-thinking module and Pre-answering module and the latter 1 means Pre-thinking scope.)

Dataset	Properties			Size		License
	Domain	Reasoning	Ans. Type	Train	Test	
GSM8K-AUG	Math	Arithmetic	Number	385,620	1,319	MIT
GSM-Hard	Math	Arithmetic	Number	1,319	-	MIT
MultiArith	Math	Arithmetic	Number	420	180	CC BY 4.0
SVAMP	Math	Arithmetic	Number	700	300	MIT
ASDiv-AUG	Math	Arithmetic	Number	4,183	1,038	CC BY 4.0

Table 5: Overview of the datasets. We categorize the datasets by their domain properties and sample sizes.

Dataset	Avg. Q Len.	Avg. CoT Len.	Avg. Steps
GSM8K-AUG	61.3	36.2	3.26
GSM-Hard	65.8	154.2	-
MultiArith	41.3	-	-
SVAMP	41.5	15.0	1.28
ASDiv-AUG	36.6	12.3	1.23

Table 6: Detailed statistics of the **test split** used in our experiments. Lengths are calculated based on token count, and steps represent the reasoning depth of the gold solutions.

884 problem type and grade level to indicate difficulty.  
885 ASDiv-AUG serves as an augmented or hardened  
886 version of this benchmark, typically constructed  
887 by filtering for higher-complexity problems or gener-  
888 ating challenging variations to demand stronger  
889 mathematical reasoning capabilities from the mod-  
890 els.

Dataset	Example Question	Rationale / Steps	Answer
GSM8K-Aug	Out of 600 employees in a company, 30% got promoted while 10% received bonus. How many employees did not get either a promotion or a bonus?	[ "«600*30/100=180»", "«600*10/100=60»", "«180+60=240»", "«600-240=360»" ]	360
GSM-Hard	Janet's ducks lay 16 eggs per day. She eats three for breakfast every morning and bakes muffins for her friends every day with 4933828. She sells the remainder at the farmers' market daily for \$2 per fresh duck egg. How much in dollars does she make every day at the farmers' market?	def solution(): ""Janet's ducks lay 16 eggs per day. She eats three for breakfast every morning and bakes muffins for her friends every day with four. She sells the remainder at the farmers' market daily for \$2 per fresh duck egg. How much in dollars does she make every day at the farmers' market?"" eggs_per_day = 16 eggs_eaten = 3 eggs_baked = 4933828 eggs_sold = eggs_per_day - eggs_eaten - eggs_baked price_per_egg = 2 money_made = eggs_sold * price_per_egg result = money_made return result	-9,867,630
MultiArith	There are 64 students trying out for the school's trivia teams. If 36 of them didn't get picked for the team and the rest were put into 4 groups, how many students would be in each group?	-	7
SVAMP	There are 87 oranges and 290 bananas in Philip's collection. If the bananas are organized into 2 groups and oranges are organized into 93 groups How big is each group of bananas?	( 290.0 / 2.0 )	145
ASDiv-AUG	gino has 64 popsicle sticks . i have 100 popsicle sticks . what is the sum of our popsicle sticks ?	«64+100=164»	164

Table 7: Examples from the evaluation and training datasets. Each entry shows a representative question, the reasoning steps (Rationale) required to solve it, and the final numerical answer.