

# Rethinking Thinking Tokens: Understanding Why They Underperform in Practice

Sreeram Vennam<sup>1</sup>, David Valente<sup>2</sup>, David Herel<sup>3</sup>, Ponnurangam Kumaraguru<sup>1</sup>,

<sup>1</sup>IIT Hyderabad

<sup>2</sup>Instituto Superior Técnico

<sup>3</sup>FEE, Czech Technical University in Prague,

## Abstract

Thinking Tokens (TT) (Herel and Mikolov, 2023) have been proposed as an unsupervised method to facilitate reasoning in language models. However, despite their conceptual appeal, our findings show that TTs marginally improve performance and consistently underperform compared to Chain-of-Thought (CoT) reasoning across multiple benchmarks. We hypothesize that this underperformance stems from the reliance on a single embedding for TTs, which results in inconsistent learning signals and introduces noisy gradients. This paper provides a comprehensive empirical analysis to validate this hypothesis and discusses the implications for future research on unsupervised reasoning in LLMs.

## 1 Introduction

Large Language Models (LLMs) have demonstrated unprecedented performance across a wide array of natural language processing tasks, from translation to creative text generation. However, reasoning remains one of the key challenges in LLM research. Recent innovations, such as Chain-of-Thought (CoT) prompting (Wei et al., 2022a), have shown considerable promise by breaking down complex tasks into sequential reasoning steps. This method has led to significant performance improvements on reasoning benchmarks like GSM8K and CommonsenseQA. The step-by-step nature of CoT enables explicit, interpretable reasoning but typically requires manual or supervised intervention through well-structured prompts.

Thinking Tokens (TTs) (Herel and Mikolov, 2023), in contrast, provide an unsupervised mechanism for reasoning. TTs work by introducing an intermediate “thinking” token, providing the model with more time to reason internally before generating output. This delay is theorized to allow deeper computation over hidden states, enhancing reasoning capacity by operating within the latent space

of the model, which is potentially more expressive than the token space where CoT functions.

Despite their theoretical advantages, Thinking Tokens have underperformed in comparison to CoT across multiple reasoning benchmarks. This paper seeks to address two key questions: (1) How do Thinking Tokens compare against Chain-of-Thought prompting? (2) Why do Thinking Tokens underperform, and what could be a potential solution?

We hypothesize that TTs’ underperformance stems from their reliance on a single-token embedding, which may introduce noise and inconsistency during gradient updates. To explore this, we conduct controlled experiments on a range of reasoning tasks. Our contributions include: 1) An empirical comparison between TT and CoT reasoning across arithmetic, and symbolic tasks. 2) Identification of the root cause of TT’s underperformance. 3) Empirical validation of this hypothesis through gradient analysis.

## 2 Related Work

Reasoning in Large Language Models (LLMs) has been an active area of research in recent years. Chain-of-Thought (CoT) prompting (Wei et al., 2022b) has emerged as a seminal technique, demonstrating that breaking down complex tasks into step-by-step reasoning improves performance across various benchmarks. CoT has proven particularly effective on tasks requiring multi-step reasoning, such as arithmetic problems, commonsense reasoning, and symbolic reasoning tasks (Zhang et al., 2023).

To mitigate the need for manually structured prompts, Thinking Tokens (TTs) (Herel and Mikolov, 2023) were introduced as an unsupervised approach to reasoning. TTs attempt to extend the reasoning capabilities of LLMs by inserting a “thinking” token, allowing the model more time to

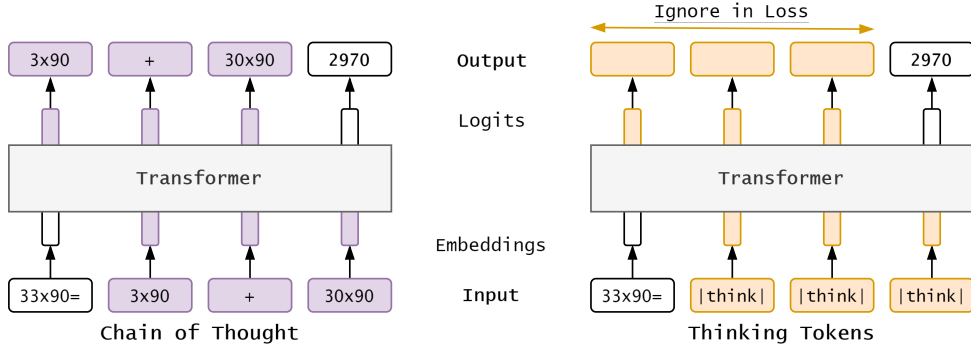


Figure 1: Chain of thought compared to thinking tokens. These approaches show striking similarity despite their differences.

compute over latent states before output generation. While this method promises deeper internal computation, empirical results show that TTs fall short of the performance achieved by CoT, likely due to training instabilities introduced by a single-token embedding. In parallel, other methods like Pause Tokens (Goyal et al., 2024) have been explored to simulate cognitive pauses, placing tokens at random intervals to mimic the effect of reasoning time. However, none of these alternatives have matched the structured reasoning power that CoT delivers. Moreover, studies have delved into understanding the nature of reasoning in LLMs by analyzing their circuit complexity, and internal representations, offering insights into how models handle tasks requiring logical and sequential thinking (Zhang et al., 2023; Malach et al., 2024). These analyses are crucial for understanding why approaches like CoT excel, while unsupervised methods like Thinking Tokens struggle to achieve similar success.

### 3 Hypothesis and Theoretical Analysis

We hypothesize that the core issue with Thinking Tokens lies in the embedding mechanism. When a single embedding is used for TTs, during back-propagation the model receives inconsistent learning signals, leading to noisy gradient updates. This noise disrupts learning, particularly in tasks that require structured intermediate steps, such as arithmetic reasoning or multi-hop commonsense tasks.

To formalize this, let  $\mathbf{h}_t$  represent the hidden state at time step  $t$ , and let  $\mathbf{e}_{TT}$  be the embedding associated with the Thinking Token. The gradient  $\nabla L(\mathbf{e}_{TT})$  with respect to the loss function  $L$  becomes inconsistent across training examples, as  $\mathbf{e}_{TT}$  serves multiple, contextually distinct roles across the token sequence. This contrasts with

CoT, where each reasoning step has an explicit, interpretable role in the output.

#### 3.1 Embedding Space, Gradient Dynamics, and Chain Representation

The core principles behind both Thinking Tokens (TTs) and Chain-of-Thought (CoT) lie in how they facilitate intermediate reasoning steps within a model. While CoT provides explicit intermediate representations, TTs introduce a single or multiple shared token embeddings to emulate intermediate steps. In CoT, reasoning steps are encoded as distinct tokens  $e_{CoT}^1, e_{CoT}^2, \dots, e_{CoT}^m$ , where  $m$  denotes the number of reasoning steps. The transformer model generates a sequence of hidden states:

$$h_t = f(e_{CoT}^t, h_{t-1}), \quad (1)$$

where  $h_{t-1}$  is the hidden state from the previous step, and  $f$  is a non-linear transformation (e.g., a multi-layer perceptron). The explicit decomposition of reasoning into distinct tokens yields structured and stable gradient updates:

$$\nabla L(e_{CoT}^t) = \frac{\partial L}{\partial h_t}, \quad (2)$$

where  $L$  is the loss function. Each  $e_{CoT}^t$  represents a unique reasoning subtask, allowing the model to isolate learning signals for each step and reduce noise during training.

Thinking Tokens (TTs) adopt a different mechanism by using one or more shared tokens  $e_{TT}^1, e_{TT}^2, \dots, e_{TT}^k$  to facilitate internal reasoning:

$$h_t = f(e_{TT}^i, h_{t-1}), \quad (3)$$

where  $i$  indexes the set of shared tokens. In this case, the gradient updates are noisier because the

|          | Digit Multiplication |      |     | Natural Language Datasets |            |
|----------|----------------------|------|-----|---------------------------|------------|
|          | 2d                   | 3d   | 4d  | GSM8k                     | OpenBookQA |
| Baseline | 0.0                  | 0.0  | 0.0 | 6.30                      | 37.2       |
| TT       | 7.32                 | 0.01 | 0.0 | 4.51                      | 37.2       |
| CoT      | 91.9                 | 66.3 | > 0 | 18.7                      | 42.0       |
| TT + CoT | 92.3                 | 67.8 | > 0 | 17.5                      | 39.6       |

Table 1: Performance of our four configurations across various tasks. We report accuracy over the integer value in Digit Multiplication. We report exact match accuracy on GSM8k and OpenBookQA. > 0 indicates the experiment produced non-zero results but was killed early due to resource constraints.

same token embeddings  $e_{TT}^i$  are reused across multiple reasoning steps:

$$\Delta e_{TT} = \sum_{i=1}^k \nabla L(e_{TT}^i). \quad (4)$$

This reuse of embeddings across different contexts introduces ambiguity into the learning signals, making it harder for the model to cleanly separate the contribution of each reasoning step, leading to a noisier gradient signal compared to CoT.

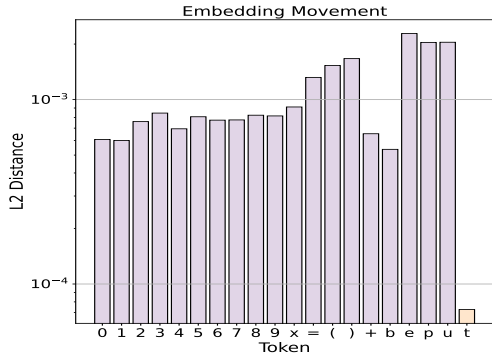


Figure 2: One TT embedding hardly moves from the initialized value.

## 4 Experiments

We first compare TT against CoT in section 4.1 on both synthetic data and popular natural language benchmarks. We perform our analysis on tasks in which intermediate computation can help neural networks learn, such as GSM8k. We then verify our hypothesis empirically in section 4.2 by monitoring both the embedding weights and the gradient of the vocabulary layer of the model.

### 4.1 TT vs CoT

We compare performance across the below 4 configurations. Our results can be found in Table 1.

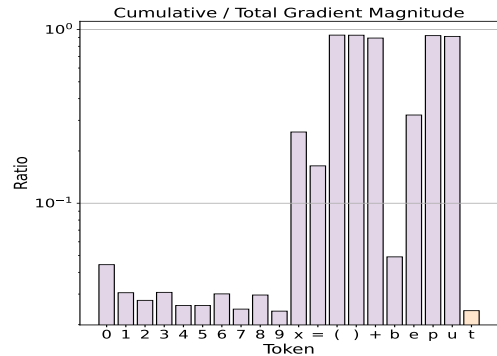


Figure 3: One TT embedding receives insufficient cumulative gradient.

**Baseline** The base configuration without intermediate steps, expecting the multiplication result immediately after the expression.

**CoT** Pretraining on CoT using the intermediate calculations. This involves supervising the model with the intermediate calculations during training.

**TT** Pretraining using TTs hoping that the model learns to use these tokens to reason through the problem.

**TT + CoT** Combining CoT with TT during pretraining.

#### 4.1.1 Digit Multiplication

We adopt a synthetic dataset for integer multiplication from (Malach et al., 2024) which includes intermediate calculations that takes the role of CoT. We conduct our experiments using four main configurations across 2-digit, 3-digit, and 4-digit multiplication. The think token is depicted as "t".

**Experimental Setup** We train a GPT-2 based transformer for 100 epochs, utilizing proportionally more data samples as the digit count increased.

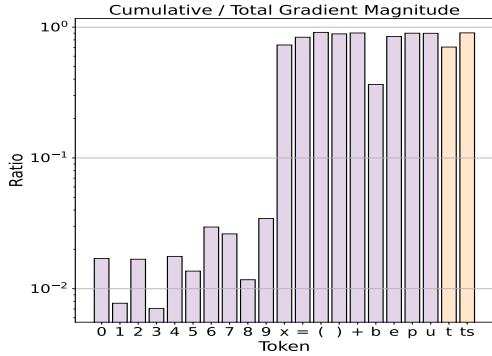


Figure 4: Two TT embeddings receive clear large cumulative gradients.

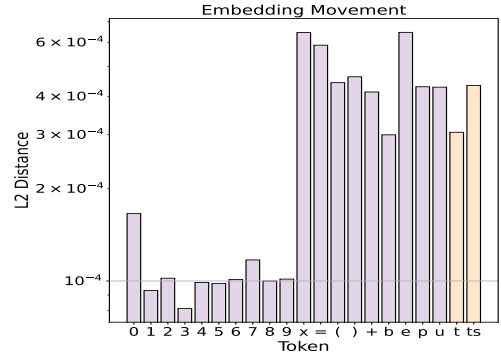


Figure 5: Two TT embeddings show clear deviation from initialization.

Implementation details can be found in the Appendix.

#### 4.1.2 Natural Language Tasks

CoT can aid NLP tasks as well. In this section, we compare TT against CoT for standard NLP tasks such as mathematical reasoning (GSM 8k) and question answering (OpenBookQA).

**Experimental Setup** We finetune Llama 3.2 (1B) (Dubey et al., 2024) on two datasets, GSM8k, and OpenBookQA. CoT data for these datasets are available through ThoughtSource (Ott et al., 2023). Implementation details can be found in the Appendix.

#### 4.2 Gradient Analysis

To validate our hypothesis (Section 3) empirically, we record the embeddings and gradients of the embedding layer regularly during training and analyze them. Furthermore, we introduce two thinking tokens with distinct embeddings, “*t*” and “*ts*” and monitor their embeddings as well, if our hypothesis is true, we should see clearer gradients using two distinct tokens.

**Noisy Gradients** We calculate how far the embedding travels. If the embedding hasn’t moved much, during training, that could be a strong indicator of noisy gradients. We also calculate the cumulative gradient. Noisy gradients result in a low cumulative gradient since the mean of noise is zero.

**Inconsistent Learning Signals** We calculate the gradient direction variance – high directional variance where the token embedding erratically changes to accommodate every context can be a sign of inconsistent learning signals.

## 5 Analysis

Table 1 shows us that TT performs worse than CoT across multiple benchmarks. Sometimes it can even hurt model performance, a finding echoed by (Goyal et al., 2024). In summary, TT offers only marginal gains while being overshadowed by CoT.

### 5.1 Gradient Analysis

**One Embedding** We see in Figure 2 that the embedding hardly travels from its initialization showing the tokens limited expressivity in practice. Figure 3 describes how the cumulative gradient is the smallest of the tokens explain why the embedding is stagnant.

**Two Embeddings** The addition of two unique embeddings results in significantly clearer gradients. Figure 4 shows that these tokens now receive solid gradients. Figure 5 reveals that the embeddings are no longer stagnant and move from initialization, verifying our hypothesis.

## 6 Discussion

Our results suggest that while Thinking Tokens offer a novel unsupervised mechanism for reasoning, they suffer from noisy gradients due to their single embedding mechanism. When paired with CoT, they perform as if they weren’t present. Although TT looks appealing due to unsupervised reasoning within the latent space, future approaches should aim to better make use of this space with richer vectors rather than a single embedding.

## References

- Abhimanyu Dubey, Abhinav Jauhri, Abhinav Pandey, Abhishek Kadian, Ahmad Al-Dahle, Aiesha Letman, Akhil Mathur, Alan Schelten, Amy Yang, Angela Fan, et al. 2024. The llama 3 herd of models. *arXiv preprint arXiv:2407.21783*.
- Sachin Goyal, Ziwei Ji, Ankit Singh Rawat, Aditya Krishna Menon, Sanjiv Kumar, and Vaishnavh Nagarajan. 2024. [Think before you speak: Training language models with pause tokens](#). In *International Conference on Learning Representations (ICLR)*.
- David Herel and Tomas Mikolov. 2023. [Thinking tokens for language modeling](#). In *Proceedings of the 8th Conference on Artificial Intelligence and Theorem Proving (AITP 2023)*.
- Eran Malach et al. 2024. [Auto-regressive next-token predictors are universal learners](#). In *Proceedings of the 41st International Conference on Machine Learning (ICML)*.
- Simon Ott, Konstantin Hebenstreit, Valentin Liévin, Christoffer Egeberg Hother, Milad Moradi, Maximilian Mayrhauser, Robert Praas, Ole Winther, and Matthias Samwald. 2023. [Thoughtsource: A central hub for large language model reasoning data](#). *arXiv preprint*.
- Jason Wei, Xuezhi Wang, Dale Schuurmans, Maarten Bosma, brian ichter, Fei Xia, Ed Chi, Quoc V Le, and Denny Zhou. 2022a. [Chain-of-thought prompting elicits reasoning in large language models](#). In *Advances in Neural Information Processing Systems*, volume 35, pages 24824–24837. Curran Associates, Inc.
- Jason Wei et al. 2022b. Chain of thought prompting elicits reasoning in large language models. *arXiv preprint arXiv:2201.11903*.
- Xue Zhang, Yuchen Sun, et al. 2023. Towards revealing the mystery behind chain of thought: A theoretical perspective. *arXiv preprint arXiv:2305.15408*.

## A Appendix

### A.1 Experimental Setup

#### A.1.1 Digit Multiplication Experiments

For our digit multiplication experiments, we utilized a synthetic dataset for integer multiplication as described in Malach et al. (2024). The following configurations were employed during the training process:

- **Model:** We implemented a GPT-2 based transformer model.
- **Epochs:** The model was trained for a total of 50 epochs.
- **Batch Size:** A batch size of 128 was used throughout the training.
- **Learning Rate:** The learning rate was determined using a learning rate finder as implemented in PyTorch Lightning (cite as lr-finder).
- **Hardware:** The experiments were conducted on 2 L40s (GPUs).

The dataset included multiple samples across different digit counts, specifically targeting 2-digit, 3-digit, and 4-digit multiplication tasks.

#### A.1.2 Natural Language Processing Tasks

For the NLP tasks, we focused on fine-tuning the Llama 3.2 (1B) model on established benchmarks:

- **Datasets:** The models were evaluated on GSM8K and OpenBookQA datasets, which are known for their reasoning complexity.
- **Epochs:** The training was performed for 5 epochs.
- **Batch Size:** An effective batch size of 512 was employed to ensure effective training dynamics.
- **Hardware:** The experiments were carried out using 2 L40s (GPUs).

The experiments aimed to assess the performance of Thinking Tokens (TTs) compared to Chain-of-Thought (CoT) prompting in various reasoning tasks, leveraging the structured reasoning benefits offered by CoT.