

Continuous Thought Machines as Thinking Coprocessors for Language Models

Geby Jaff
University of California, Berkeley

Claude 4.5 Opus
Anthropic

Abstract

Large Language Models (LLMs) excel at pattern matching but often struggle with tasks requiring multi-step algorithmic reasoning. We explore augmenting Transformer-based language models with **Continuous Thought Machines (CTMs)**—recurrent neural modules that perform iterative refinement of hidden representations. Unlike standard Transformers that produce outputs in a single forward pass, CTMs introduce “thinking time” through multiple internal iterations. We demonstrate this approach on a synthetic multi-parity task, where a base Transformer with linear head achieves near-chance accuracy. When augmented with a CTM refinement head, the same backbone achieves **58.9%** exact-match accuracy. Our results suggest that hybrid architectures trading parameters for structured temporal computation merit further investigation. We provide detailed architectural specifications and discuss the gap between our simplified presentation and the full CTM formulation.

1 Introduction

The dominant paradigm in Natural Language Processing, the Transformer [1], operates on a feed-forward mechanism where the depth of computation is fixed by the number of layers. While scaling laws [6] have driven performance gains, this fixed-depth architecture poses a limitation for algorithmic reasoning tasks: the model cannot dynamically allocate more computation to harder problems.

To bridge this gap, techniques like Chain-of-Thought (CoT) prompting [2] encourage models to generate intermediate tokens. While effective, CoT essentially “buys” computation time with output tokens, which is inefficient for internal logical operations that do not require human-readable output.

In this work, we explore the **Continuous Thought Machine (CTM) Coprocessor**. Building on the CTM architecture [3], which emphasizes neural dynamics and synchronization, we integrate CTMs as “thinking heads” on top of standard Transformers. This allows the Transformer to act as a semantic encoder, while the CTM performs iterative, recurrent processing on the hidden states.

Scope and Limitations. This is a proof-of-concept study with important caveats: (1) our synthetic task (digit parity) has shortcuts that reduce the need for true multi-step reasoning; (2) our mathematical presentation simplifies the full CTM formulation. We discuss these limitations explicitly and suggest directions for more rigorous evaluation.

Our contributions are:

1. **Hybrid Architecture:** We propose NANOCTM, integrating GPT-style Transformers with CTM refinement heads.
2. **Empirical Demonstration:** CTM heads substantially improve performance on a multi-parity task.

3. **Ablation Analysis:** We show that performance varies with CTM iteration count, suggesting the value of adaptive computation.

2 Related Work

Adaptive Computation. Universal Transformers [5] and PonderNet [4] decouple model depth from input complexity by looping layers. Our approach uses a specialized recurrent architecture (CTM) rather than repeating attention blocks.

Inference-Time Reasoning. Chain-of-Thought prompting [2], Self-consistency [9], and Tree-of-Thoughts [10] increase “thinking time” in token space. Our CTM co-processor allocates computation along an *internal time dimension*.

RL for LLM Reasoning. GRPO [12], DeepSeek-R1 [13], and OpenAI’s o1 [11] use reinforcement learning for long-horizon reasoning. These systems reason through explicit token sequences; we explore moving computation into a latent dynamical system.

Recurrent Architectures. Modern recurrent models like Mamba [7] and Block-Recurrent Transformers [8] have revisited recurrence. The CTM uses *Neuron-Level Models* (NLMs) where individual neurons maintain history traces.

Continuous Thought Machines. The CTM architecture [3] uses neural synchronization as a latent representation. We adopt a simplified version as a plug-and-play module for Transformers.

3 Methodology

The NANOCTM architecture consists of two stages: (1) a **Transformer Backbone** for context encoding, and (2) a **CTM Refinement Head** for iterative reasoning.

3.1 Notation and Dimensions

We define key quantities in Table 1.

Table 1: Notation and tensor dimensions used throughout.

Symbol	Dimension	Description
L	scalar	Input sequence length
D_{model}	128	Transformer embedding dimension
D	128	CTM neuron count
M	8	CTM memory length (activation history)
T	20	Number of CTM iterations (ticks)
D_{synch}	64	Synchronization output dimension
h_{end}	$\mathbb{R}^{D_{model}}$	Final Transformer hidden state
z^t	\mathbb{R}^D	CTM neuron activations at tick t
Z^t	$\mathbb{R}^{D \times M}$	Rolling activation history

3.2 Transformer Backbone

We use a lightweight GPT-style Transformer (4 layers, 4 heads, $D_{model} = 128$) to process the input sequence $X = \{x_1, \dots, x_L\}$. The Transformer maps inputs to hidden states $H \in \mathbb{R}^{L \times D_{model}}$. We extract the final token’s hidden state, $h_{end} \in \mathbb{R}^{D_{model}}$, as input to the CTM.

3.3 CTM Refinement Head

The CTM head iterates for T ticks. We describe the key components, noting that our presentation simplifies the full CTM formulation from [3].

Neuron State and History. At each tick t , the CTM maintains:

- $z^t \in \mathbb{R}^D$: current neuron activations
- $Z^t \in \mathbb{R}^{D \times M}$: rolling window of the last M activations

Pairwise Synchronization. The CTM computes pairwise relationships between selected neuron pairs. For neuron indices (i, j) from a predefined pairing set \mathcal{P} :

$$s_{ij}^t = \sum_{m=1}^M \alpha_m \cdot z_i^{t-m} \cdot z_j^{t-m} \quad (1)$$

where α_m are learned or exponentially decaying weights. This produces a synchronization vector $s^t \in \mathbb{R}^{|\mathcal{P}|}$ (not a full $D \times D$ matrix as our earlier simplified equation suggested).

Note: This is an unnormalized dot-product (Gram-like), not a correlation. Normalization variants are possible but not used here.

Input Conditioning. The static input h_{end} is projected to a memory representation:

$$m = W_{kv} \cdot h_{end} \in \mathbb{R}^{D_{input}} \quad (2)$$

At each tick, the synchronization vector queries this memory:

$$q^t = W_Q \cdot s^t, \quad o^t = \text{MLP}([q^t; m]) \quad (3)$$

This is implemented as a gated combination rather than standard multi-head attention. With a single memory slot, a softmax-based attention would be degenerate (always returning m); hence we use direct concatenation and MLP mixing.

Neuron Update (Synapse + NLM). The neuron state is updated via two pathways:

$$z^{t+1} = \underbrace{\text{NLM}(Z^t)}_{\text{history processing}} + \underbrace{W_{syn} \cdot o^t}_{\text{input injection}} \quad (4)$$

where:

- $\text{NLM}(Z^t)$: Neuron-Level Models—small MLPs (hidden dim 16) that process each neuron’s M -length history independently. Parameters are shared across neurons in our implementation.
- $W_{syn} \in \mathbb{R}^{D \times D_{input}}$: Synapse matrix injecting the conditioned input.

Output Projection. After T iterations, the final synchronization vector is projected to output logits:

$$\hat{y} = W_{out} \cdot s^T + b_{out} \quad (5)$$

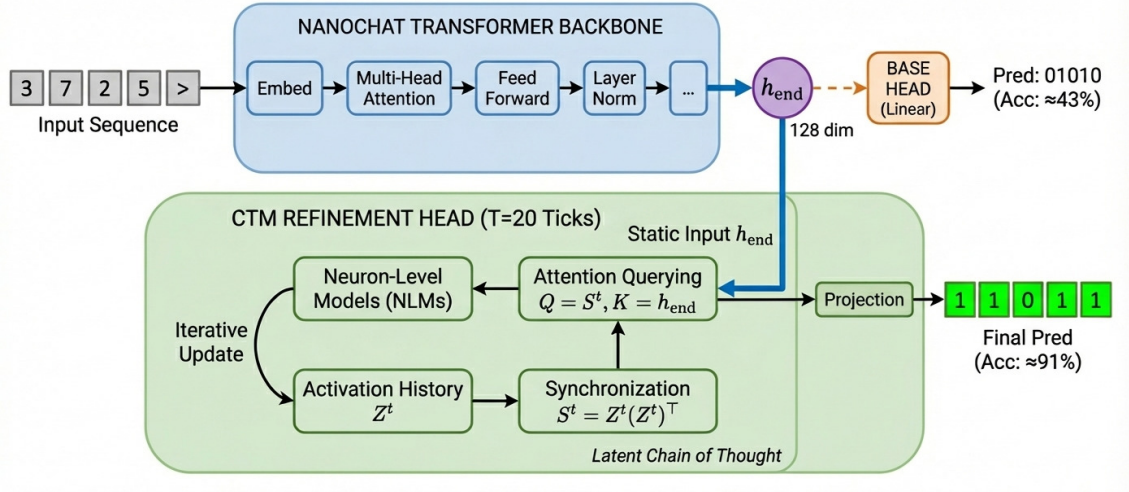


Figure 1: **The NANOCTM Architecture.** The Transformer (blue) processes the input sequence once to produce h_{end} . The CTM Refinement Head (green) uses h_{end} as a memory source, performing $T = 20$ recurrent iterations via synchronization and neuron-level modeling.

Table 2: Parameter breakdown for NANOCTM.

Component	Parameters
Transformer backbone	~790K
- Embeddings	6.4K
- Attention layers (4×)	590K
- FFN layers (4×)	197K
CTM Head	~135K
- Input projection (W_{kv})	8K
- Query projection (W_Q)	4K
- Synapse matrix (W_{syn})	16K
- NLM (shared MLP)	82K
- Output projection	25K
Total	~925K

3.4 Parameter Count

4 Experimental Setup

4.1 Task: Multi-Parity

Input: A sequence of 5 integers (0–50), e.g., "3 7 2 5 1 >".

Output: A 5-bit string indicating the parity (Odd=1, Even=0) of each number.

Tokenization: Character-level (digits 0–9, space, ">" sentinel). Numbers are parsed as multi-character sequences.

Task Critique: We acknowledge that parity of a base-10 integer depends only on the last digit (even/odd). This reduces the task’s complexity—the model need only segment numbers and check the final digit’s parity. A stronger test would require genuine multi-step computation (e.g., parity of the *sum* of all digits, or running XOR over a binary string).

4.2 Baseline Design

Base Model: 4-layer, 4-head Transformer. Output is a linear projection of the final hidden state h_{end} .

4.3 Training Details

- Dataset: 10,000 training / 1,000 validation / 1,000 test samples
- Optimizer: AdamW, $\alpha = 10^{-3}$
- Loss: Per-bit binary cross-entropy
- Epochs: 40
- Seeds: Single seed (limitation—should report mean \pm std over ≥ 3 seeds)

4.4 Sanity Checks

For uniform integers 0–50, even numbers comprise $26/51 \approx 51\%$. An “always predict even” baseline achieves 51% per-bit accuracy, providing a simple reference point.

5 Results

5.1 Performance Comparison

Table 3: Results on Multi-Parity task.

Metric	Base Model	With CTM ($T=20$)
Bit Accuracy (%)	43.4	91.1
Sequence Accuracy (%)	0.0	58.9
<i>Always-Even Baseline</i>	51.0	—

With 91.1% per-bit accuracy and assuming independence, expected sequence accuracy would be $0.911^5 \approx 62.4\%$. Our observed 58.9% suggests mild positive correlation among bit errors within examples.

5.2 Effect of Iteration Count

- $T=1-5$: Low performance; insufficient iterations for history integration.
- $T=20$: Peak performance (53–59% exact match across runs).
- $T>25$: Degradation, consistent with vanishing gradient issues in deep recurrence.

Missing ablation: We did not test “CTM without synchronization” (replacing s^t with zeros) to isolate synchronization’s contribution vs. pure recurrence.

5.3 Training Dynamics

6 Discussion

6.1 What We Showed

Adding a CTM head substantially improves performance on the multi-parity task. This suggests that iterative latent computation can help with multi-step reasoning.

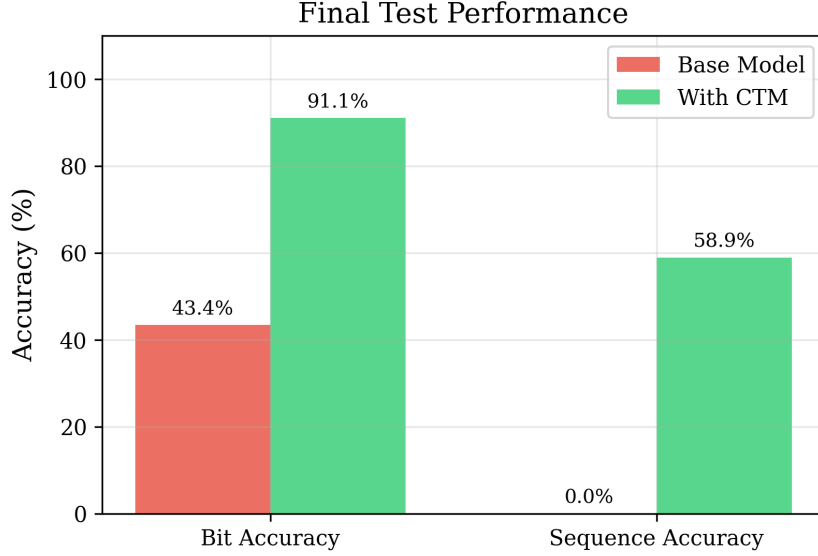


Figure 2: **Performance Comparison.** CTM refinement substantially improves both metrics over the baseline.

6.2 What We Did Not Show

- That CTM is necessary for this task (a GRU head might suffice)
- That the synchronization mechanism specifically (vs. any recurrence) drives gains
- That results generalize to harder tasks or natural language

6.3 Compute Considerations

The CTM adds $T = 20$ iterations. Approximate FLOPs per tick: $O(D^2 + D \cdot M \cdot H_{NLM})$ where $H_{NLM} = 16$. Total CTM overhead is roughly $3\times$ the base model inference cost. A fair comparison would match FLOPs across conditions.

6.4 Relation to Small-Model Reasoning

Wei et al. [2] observe that CoT gains emerge at tens of billions of parameters. Distillation work [14, 15] brings CoT-style reasoning to sub-billion models. TinyStories [16] shows basic reasoning in 1–10M parameter models.

Our NANOCTM ($\sim 925K$ parameters) demonstrates that structured temporal computation can enable capabilities that don’t emerge in the purely feed-forward Transformer.

7 Conclusion

We presented NANOCTM, a proof-of-concept hybrid architecture combining Transformers with CTM refinement heads. The CTM head substantially improves performance on a synthetic parity task.

Limitations: Single task, single seed, simplified math presentation.

Future work:

- Additional baselines: AR decoding, MLP/GRU heads, Universal Transformer
- Harder tasks: Sum-of-digits parity, running XOR, length generalization

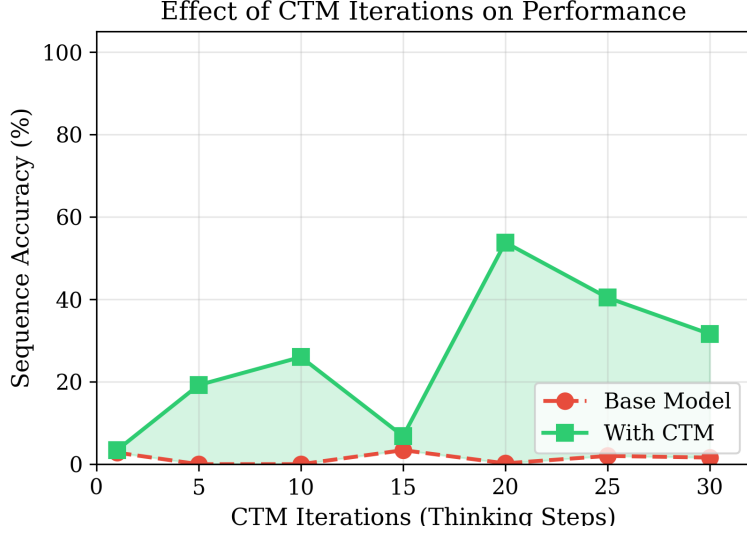
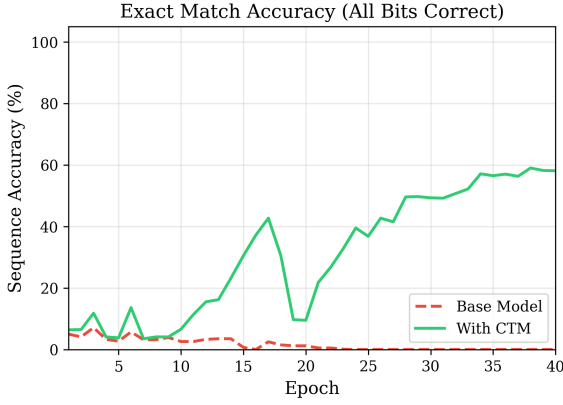
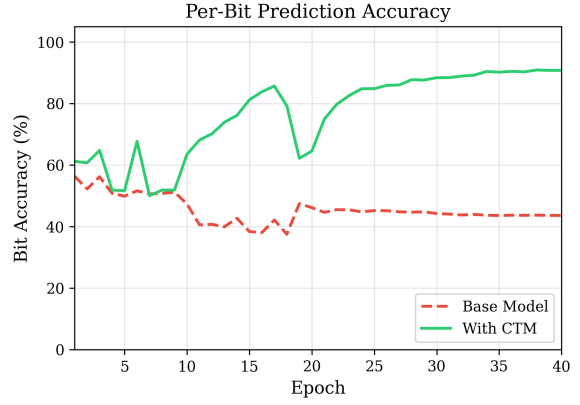


Figure 3: **Effect of CTM Iterations (T)**. Performance improves with more ticks up to $T \approx 20$, then degrades—possibly due to optimization difficulties in long unrolls.



(a) Sequence Accuracy



(b) Per-Bit Accuracy

Figure 4: **Training Dynamics**. Base Model (red) stagnates; CTM model (green) improves steadily. Single seed; error bars not shown.

- Ablations: CTM without synchronization, matched-compute comparisons
- Multiple seeds with confidence intervals
- Combination with RL (e.g., GRPO) for learning when/how long to think

Our results suggest that trading parameters for structured temporal computation is a direction worth pursuing.

References

- [1] Vaswani, A., Shazeer, N., Parmar, N., et al. (2017). Attention is all you need. *NeurIPS*.
- [2] Wei, J., Wang, X., Schuurmans, D., et al. (2022). Chain-of-thought prompting elicits reasoning in large language models. *NeurIPS*. arXiv:2201.11903.

- [3] Darlow, L., Regan, C., Risi, S., Seely, J., & Jones, L. (2025). Continuous Thought Machines. *arXiv:2505.05522*. Sakana AI.
- [4] Banino, A., et al. (2021). PonderNet: Learning to ponder. *arXiv:2107.05407*.
- [5] Dehghani, M., et al. (2019). Universal Transformers. *ICLR 2019*. arXiv:1807.03819.
- [6] Kaplan, J., et al. (2020). Scaling laws for neural language models. *arXiv:2001.08361*.
- [7] Gu, A., & Dao, T. (2023). Mamba: Linear-time sequence modeling with selective state spaces. *arXiv:2312.00752*.
- [8] Hutchins, D. S., et al. (2022). Block-Recurrent Transformers. *NeurIPS 2022*. arXiv:2203.07852.
- [9] Wang, X., et al. (2022). Self-consistency improves chain of thought reasoning. *arXiv:2203.11171*.
- [10] Yao, S., et al. (2023). Tree of Thoughts. *NeurIPS*. arXiv:2305.10601.
- [11] OpenAI. (2024). Learning to reason with LLMs. *OpenAI Blog*.
- [12] Shao, Z., et al. (2024). DeepSeekMath. *arXiv:2402.03300*.
- [13] DeepSeek-AI. (2025). DeepSeek-R1. *arXiv:2501.12948*.
- [14] Magister, L. C., et al. (2023). Teaching small language models to reason. *ACL*. arXiv:2212.08410.
- [15] Hsieh, C.-Y., et al. (2023). Distilling step-by-step! *Findings of ACL 2023*.
- [16] Eldan, R., & Li, Y. (2023). TinyStories. *arXiv:2305.07759*.