

HighNoon LLM: Advancing Sequence Processing with Hierarchical Spatial Neural Memory

Anonymous

July 2025

Abstract

Transformer-based models have significantly advanced natural language processing (NLP) by capturing long-range dependencies, yet their quadratic computational complexity limits scalability for extended sequences. This paper presents HighNoon LLM, a novel large language model integrating Hierarchical Spatial Neural Memory (HSMN) to overcome these challenges. HSMN employs a hierarchical memory tree, constructed via chunking and aggregation, reducing complexity to $O(n \cdot c)$ —where c is a fixed chunk size—while explicitly modeling nested dependencies. We detail HSMN’s architecture, including its ChunkEncoder, Aggregator, and ReasoningModule, and provide rigorous mathematical formulations. Comparative analyses with standard transformers, efficient variants (e.g., Longformers), and other hierarchical models demonstrate HSMN’s superior efficiency and expressiveness. Additional features, such as continual learning via Elastic Weight Consolidation and on-device processing, enhance its adaptability and practicality. Although empirical results are preliminary, HSMN’s design promises transformative applications in document translation, summarization, and beyond, aligning with ethical innovation principles.

1 Introduction

Transformer-based models [1] have redefined sequence processing in NLP, excelling in tasks like machine translation and text generation. However, their self-attention mechanism incurs a quadratic complexity of $O(n^2)$ with respect to sequence length n , rendering them impractical for long sequences. Efforts to mitigate this—such as Longformers [2], Reformers [3], and Linformers [4]—reduce complexity but often sacrifice explicit hierarchical modeling.

HighNoon LLM introduces Hierarchical Spatial Neural Memory (HSMN), a scalable architecture that processes long sequences efficiently while capturing hierarchical structures inherent in language (e.g., phrases within sentences). By dividing sequences into chunks, constructing a hierarchical memory tree, and enabling context-aware reasoning, HSMN achieves a complexity of $O(n \cdot c)$, where c is the chunk size (e.g., 128). This paper elucidates HSMN’s design, compares it with existing models, and explores its potential in applications requiring deep contextual understanding, all while emphasizing ethical considerations inspired by principles of wisdom and stewardship [8].

2 Background

2.1 Standard Transformers

Transformers rely on self-attention to model token dependencies:

$$\text{Attention}(Q, K, V) = \text{softmax} \left(\frac{QK^T}{\sqrt{d_k}} \right) V, \quad (1)$$

where Q , K , and V are query, key, and value matrices, and d_k is the key dimension. This operation scales as $O(n^2)$, limiting scalability.

2.2 Efficient Transformer Variants

Variants address this inefficiency: - **Longformers** [2]: Use sparse attention (e.g., sliding windows), achieving $O(n \cdot k)$ complexity, where k is the window size. - **Reformers** [3]: Apply locality-sensitive hashing for $O(n)$ complexity. - **Linformers** [4]: Use low-rank approximations, also achieving $O(n)$.

These models, however, do not explicitly model hierarchical structures.

2.3 Memory-Augmented and Hierarchical Models

Neural Turing Machines (NTM) [5] and Differentiable Neural Computers (DNC) [6] use external memory for reasoning tasks, while Hierarchical Attention Networks (HAN) [7] model document hierarchies. HSMN builds on these concepts, tailoring hierarchical memory for sequence generation.

3 HSMN Architecture

HSMN comprises three components: ChunkEncoder, Aggregator, and ReasoningModule.

3.1 ChunkEncoder

The ChunkEncoder processes the input sequence $X = [x_1, x_2, \dots, x_n]$: - Divide X into $m \approx \lceil n/c \rceil$ chunks $C = [C_1, C_2, \dots, C_m]$, each of size c . - Encode each chunk: $e_i = \text{Encoder}(C_i)$, where $e_i \in \mathbb{R}^d$, yielding $E = [e_1, e_2, \dots, e_m]$. - Complexity per chunk is $O(c^2)$, total $O(n \cdot c)$.

3.2 Aggregator

The Aggregator builds a binary memory tree: - Level 0 nodes: $m_{0,i} = e_i$. - Higher levels: $m_{l,j} = f(m_{l-1,2j-1}, m_{l-1,2j})$, where f is a learned function (e.g., MLP). - Tree height is $O(\log m)$, total nodes $O(m)$.

3.3 ReasoningModule

The ReasoningModule generates output Y autoregressively: - Query vector q_t attends to memory nodes: $c_t = \sum_k \alpha_k m_k$, where $\alpha_k = \frac{\exp(q_t^T m_k)}{\sum_{k'} \exp(q_t^T m_{k'})}$. - Output probability: $p(y_t | y_{<t}, M) = \text{softmax}(g(s_t, c_t))$. - Process: $X \rightarrow E \rightarrow M \rightarrow Y$.

4 Mathematical Innovations

HSMN reduces complexity to $O(n \cdot c)$: - Transformers: $O(n^2)$ due to pairwise token interactions. - HSMN: $O(c^2)$ per chunk, $O(n/c)$ chunks, plus $O(n)$ tree construction. For $n = 10,000$, $c = 128$, HSMN requires $\sim 1.28\text{M}$ operations versus 100M for transformers—a 78x reduction.

The hierarchical tree explicitly captures nested dependencies, enhancing expressiveness over implicit attention mechanisms.

5 Comparison with Existing Models

HSMN outperforms in hierarchical modeling and scalability, balancing efficiency and expressiveness.

Table 1: Comparison of HSMN with Other Architectures

Model	Complexity	Hierarchical Modeling	Scalability
Transformer	$O(n^2)$	Implicit	Limited
Longformer	$O(n \cdot k)$	Implicit	Moderate
Reformer	$O(n)$	Implicit	High
Linformer	$O(n)$	Implicit	High
NTM/DNC	Task-dependent	Limited	Task-dependent
RNN	$O(n)$	Limited	Moderate
HAN	$O(n)$	Explicit	High
HSMN	$O(n \cdot c)$	Explicit	High

6 Additional Features

- **Continual Learning:** Uses Elastic Weight Consolidation (EWC) to mitigate catastrophic forgetting, with loss term: $\sum_k \lambda F_k(w - w_k^*)^2$. - **Local Processing:** Optimized for on-device operation, enhancing privacy and reducing latency.

7 Experimental Results

Due to proprietary constraints, empirical data is limited. Preliminary tests suggest HSMN excels in: - Document-level translation - Long-form summarization - Code generation Future benchmarks will compare HSMN with GPT-4 and Longformers.

8 Ethical Considerations

HSMN aligns with ethical innovation, promoting wisdom and stewardship [8]. Its efficiency reduces energy consumption, and local processing enhances user privacy.

9 Conclusion

HighNoon LLM’s HSMN architecture advances sequence processing by reducing complexity and explicitly modeling hierarchies. Its versatility and ethical design position it as a promising solution for scalable NLP. Future work will focus on empirical validation and adaptive memory enhancements.

References

- [1] Vaswani, A., et al. (2017). Attention is All You Need. *NeurIPS*.
- [2] Beltagy, I., et al. (2020). Longformer: The Long-Document Transformer. *arXiv:2004.05150*.
- [3] Kitaev, N., et al. (2020). Reformer: The Efficient Transformer. *ICLR*.
- [4] Wang, S., et al. (2020). Linformer: Self-Attention with Linear Complexity. *NeurIPS*.
- [5] Graves, A., et al. (2014). Neural Turing Machines. *arXiv:1410.5401*.
- [6] Graves, A., et al. (2016). Hybrid Computing Using a Neural Network with Dynamic External Memory. *Nature*.

- [7] Yang, Z., et al. (2016). Hierarchical Attention Networks for Document Classification. *NAACL*.
- [8] Proverbs 2:6, King James Bible.