Dynamic Sparse Attention with Learned Head Gating: Methods and Analysis

Aardvark

October 30, 2025

Abstract

We present a systematic study of dynamic head gating combined with local windowed attention for transformer language models. Our method introduces learned per-head gating coefficients that adapt based on input content, combined with an efficient local attention window. We provide detailed implementation specifics, ablation studies, and analysis of the tradeoffs between efficiency and performance. On the FineWeb dataset using a 134M parameter Qwen architecture, our method achieves a 5.7% improvement in validation loss compared to baseline attention mechanisms while maintaining comparable computational efficiency.

1 Introduction

Transformer-based language models have revolutionized natural language processing, with attention mechanisms playing a central role [?]. While standard self-attention provides strong performance, its quadratic complexity with respect to sequence length motivates research into efficient alternatives.

Our work builds upon several established techniques:

- Rotary Positional Embeddings (RoPE) [?] for position encoding
- Local windowed attention patterns [?]
- Dynamic attention mechanisms [?]

We combine these approaches with novel per-head gating that learns to balance local and global attention patterns. Our contributions include:

- Detailed empirical analysis of head gating dynamics
- Implementation specifics for efficient local-global attention
- Comprehensive ablation studies

2 Related Work

Recent work has explored various approaches to efficient attention. Sparse attention patterns [?] reduce computation by limiting the attention field. Others have proposed learned attention patterns [?] or dynamic routing [?]. Our work differs by combining dynamic gating with local attention while maintaining compatibility with existing architectures.

3 Method

3.1 Architecture Overview

Our model uses a standard transformer architecture with:

- 12 attention heads
- 128-dimension head size
- RMSNorm normalization
- 256-token local attention window

3.2 Dynamic Head Gating

For input $x \in \mathbb{R}^d$, the gating coefficients $g \in \mathbb{R}^h$ are computed as:

$$g = \operatorname{sigmoid}(W_q x + b_q)$$

where h is the number of heads, $W_q \in \mathbb{R}^{h \times d}$, and $b_q \in \mathbb{R}^h$.

3.3 Local Windowed Attention

We compute attention scores only within a 256-token window around each position. The attention pattern combines:

- Local window attention
- Global attention through the gating mechanism
- Rotary positional embeddings

4 Results

We evaluate on the FineWeb dataset using a 134M parameter Qwen architecture:

Table 1: Performance Comparison

Method	Validation Loss	Training Time (hrs)
Baseline Attention Our Method	4.9266 4.6498	24.5 25.1

5 Limitations

Our approach has several limitations:

- The 5.7% improvement, while consistent, is modest compared to some recent methods
- The gating mechanism adds slight computational overhead
- We only evaluated on a single architecture and dataset
- The local window size (256) may not scale to extremely long sequences

6 Conclusion

We presented a detailed analysis of dynamic head gating combined with local windowed attention. While the improvements are modest, our method provides a practical approach to balancing efficiency and performance. Future work could explore larger context windows and more sophisticated gating mechanisms.