# Advanced Attention Variants

*Efficient Attention, Sparse Attention, Flash Attention & Beyond*

## 1. The Quadratic Problem

Standard self-attention has $O(n^2)$ time and memory complexity with respect to sequence length n. For a sequence of 1000 tokens, the attention matrix has 1,000,000 entries. For 10,000 tokens, it has 100,000,000 entries. This makes standard attention prohibitively expensive for long documents, high-resolution images, or genomic sequences. This challenge spurred a wave of research into efficient attention variants.

## 2. Sparse Attention

Sparse attention restricts each token to attend to only a subset of other tokens, reducing complexity to $O(n * \sqrt{n})$ or $O(n * \log(n))$. The key insight is that not all token pairs are equally important — most of the attention weight is concentrated in a small fraction of positions.

### 2.1 Longformer Attention

Longformer (Beltagy et al., 2020) combines local windowed attention with global attention. Each token attends to a local window of w/2 tokens on each side (linear complexity), while special tokens like [CLS] attend globally to all tokens. This allows Longformer to process documents with thousands of tokens efficiently.

### 2.2 BigBird Attention

BigBird (Zaheer et al., 2020) uses a combination of random attention (each query attends to r random keys), window attention (local context), and global tokens. This mixture is theoretically proven to be a universal approximator of sequence functions and reduces complexity to $O(n)$.

## 3. Linear Attention

Linear attention methods aim to approximate or reformulate softmax attention in $O(n)$ time. The key idea is to decompose the softmax kernel using feature maps phi(x) such that:

```
softmax(q^T k) ≈ phi(q)^T phi(k)
```

This allows the attention computation to be rewritten using matrix associativity, computing KV aggregations first and then applying Q — avoiding the n×n attention matrix entirely. Examples include Performer (using random Fourier features) and Linear Transformer.

## 4. Flash Attention

Flash Attention (Dao et al., 2022) takes a different approach — instead of approximating attention, it computes exact attention but uses a hardware-aware algorithm to minimize memory I/O. The key insight is that the bottleneck in standard attention is not compute but memory bandwidth between HBM (high-bandwidth memory) and SRAM (fast on-chip memory).

Flash Attention uses tiling to split Q, K, V into blocks and processes them in SRAM, never materializing the full n×n attention matrix in HBM. This achieves 2-4x speedup and $O(n)$ memory usage while computing mathematically identical results to standard attention. Flash Attention 2 and 3 further improved performance with better parallelism strategies.

# 5. Multi-Query and Grouped-Query Attention

## 5.1 Multi-Query Attention (MQA)

Multi-Query Attention (Shazeer, 2019) uses multiple query heads but shares a single key and value head across all query heads. This dramatically reduces the size of the KV cache during inference — critical for autoregressive generation where cached keys and values consume large amounts of memory. MQA is used in models like PaLM and Falcon.

## 5.2 Grouped-Query Attention (GQA)

Grouped-Query Attention (Ainslie et al., 2023) is a middle ground between MHA and MQA. Query heads are divided into G groups, and each group shares one key and value head. GQA achieves quality close to MHA with inference speed close to MQA. It is used in Llama 2, Llama 3, Mistral, and many modern open-source LLMs.

# 6. Relative Positional Encodings

## 6.1 Rotary Position Embedding (RoPE)

RoPE (Su et al., 2021) encodes position by rotating query and key vectors in 2D planes. The dot product between a query at position m and a key at position n depends only on their relative distance (m-n), giving the model a natural sense of relative position. RoPE generalizes well to longer sequences than seen during training and is used in GPT-NeoX, LLaMA, Falcon, and most modern LLMs.

## 6.2 ALiBi (Attention with Linear Biases)

ALiBi (Press et al., 2021) adds a linear bias to attention scores based on the distance between query and key positions: score($q_i$, $k_j$) -= m * |i - j|, where m is a head-specific slope. This penalizes attending to distant tokens, encouraging locality. ALiBi models extrapolate well to longer sequences without any retraining.

# 7. Comparison of Attention Variants

| Method | Complexity | Exact? | Use Case |
|---|---|---|---|
| Standard MHA | O(n^2) | Yes | General purpose, short sequences |
| Longformer | O(n * w) | Approx | Long documents |
| BigBird | O(n) | Approx | Very long sequences |
| Flash Attention | O(n^2) compute, O(n) mem | Yes | GPU-efficient training |
| Linear Attention | O(n) | Approx | Extreme length, efficiency |
| MQA / GQA | O(n^2) | Yes | Fast inference, small KV cache |

## 8. Summary & Takeaways

The field of attention mechanisms has evolved rapidly since 2017. Standard multi-head attention remains the gold standard for quality, while efficient variants like Flash Attention, GQA, and sparse attention methods make it practical to scale to longer sequences and larger models. Choosing the right attention variant depends on your sequence length, hardware constraints, inference requirements, and acceptable quality trade-offs. Modern LLMs like Llama 3 and Mistral combine Flash Attention + GQA + RoPE for an optimal balance of quality, speed, and memory efficiency.