

MiT Attention: Efficient Fast-Weight Scaling via a Mixture of Top- k Activations

Qishuai Wen¹ Zhiyuan Huang¹ Xianghan Meng¹ Wei He¹ Chun-Guang Li¹

Abstract

The attention operator in Transformers can be viewed as a two-layer *fast-weight* MLP, whose weights are dynamically instantiated from input tokens and whose width equals sequence length N . As the context extends, the expressive capacity of such an N -width MLP increases, but scaling its fast weights becomes prohibitively expensive for extremely long sequences. Recently, this *fast-weight scaling* perspective has motivated the Mixture-of-Experts (MoE) attention, which partitions the sequence into *fast-weight experts* and sparsely routes the tokens to them. In this paper, we elevate this perspective to a unifying framework for a wide range of efficient attention methods by interpreting them as scaling fast weights through routing and/or compression. Then we propose a compress-and-route strategy, which compresses the N -width MLP into a narrower one using a small set of *landmark queries* and constructs *deformable experts* by gathering top- k activated key-value pairs for each landmark query. We call this strategy a *Mixture of Top- k Activations* (MiTA), and refer to the resulting efficient mechanism as *MiT attention*. Preliminary experiments on vision tasks demonstrate the promise of our MiTA attention and motivate further investigation on its optimization and broader applications in more challenging settings. Code is available at [this URL](#).

1. Introduction

Attention is the core operation of Transformers, which underpin today’s success and wide application of deep learning. Intuitively, attention learns to store the context as key-value associations and retrieve this short-term memory via queries (Bietti et al., 2023). However, such an all-to-all

¹School of Artificial Intelligence, Beijing University of Posts and Telecommunications, Beijing 100876, P.R. China. Correspondence to: Chun-Guang Li <lichunguang@bupt.edu.cn>.

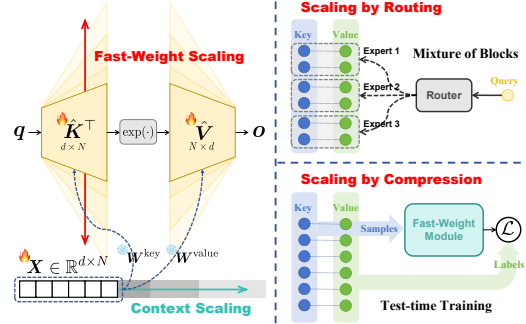


Figure 1. Fast-weight scaling and its two scaling strategies. As the context extends, the width of the two-layer fast-weight MLP induced by full attention increases accordingly. We categorize efficient fast-weight scaling approaches into two strategies: a) scaling by routing and b) scaling by compression, and illustrate each of them with a representative method: Mixture of Block Attention (MoBA) (Lu et al., 2025) and Test-time Training (TTT) (Sun et al., 2025). A pioneering method that combines both strategies is Native Sparse Attention (NSA) (Yuan et al., 2025).

lookup paradigm incurs quadratic computation and memory complexity in the sequence length, thereby hindering its scaling to long sequences. To this end, a plethora of efficient attention methods have been explored (Zhang et al., 2025).

In recent years, several lines of work have converged on a unifying viewpoint: *the key-value pairs in full attention can be viewed as the fast weights of an MLP* (Schlag et al., 2021; Han et al., 2025); *and the parameters of a two-layer MLP can be viewed as slow-weight key-value pairs* (Geva et al., 2021). Therefore, pursuing an efficient attention can be framed as a *fast-weight scaling* problem, and draw inspiration from slow-weight (i.e., parameter) scaling approaches, such as weight tying (Wu et al., 2025a), model pruning (Cheng et al., 2024), and conditional computation (Riquelme et al., 2021).

Recent work along this direction has demonstrated that sparse attention inspired by Mixture-of-Experts (MoE) can scale up the effective sequence length, deliver wall-clock speedups, and accommodate large-scale pretraining (Lu et al., 2025; Yuan et al., 2025). In contrast to traditional MoE, whose experts are structured, input-independent model parameters (i.e., slow weights), their experts are input-dependent fast weights (Schmidhuber, 1992), more precisely, subsets of key-value pairs. Routing each query to

the fast-weight experts sparsely, MoE attention reduces the complexity from quadratic to linear in sequence length N .

Roughly speaking, the central challenge in scaling the fast weights via MoE is to construct fast-weight experts from unstructured key-value pairs. For example, given the spatial and temporal locality priors in many modalities, a straightforward, hardware-friendly approach is to partition the sequence into contiguous, non-overlapping, fixed-size blocks, and regard these blocks as fast-weight experts (i.e., a mixture of blocks or chunks) (Lu et al., 2025; Yuan et al., 2025; Wu et al., 2025b; Cai et al., 2025). After that, routing vectors are obtained by aggregating each block into a single vector, e.g., via average pooling or parameterized modules.

Since splitting the N -width fast-weight MLP into a mixture of such blocks is coarse and suboptimal, subsequent work (Jia et al., 2025) has sought to improve upon the splitting scheme. Notably, top- k attention (You et al., 2025; Liu et al., 2025) can be interpreted as spawning N k -width sub-MLPs (i.e., experts) from the N -width MLP, via gathering top- k key-value pairs attended (i.e., activated) by each query. Therefore, the fixed-shape expert in prior MoE attention methods are replaced by top- k attention with deformable experts. However, instead of the width of the fast-weight MLP in full attention, it is the number of such deformable fast-weight experts that always equals N in their top- k attention.

On the other hand, unlike the above methods that scale fast weights by accessing a subset of them, linear attention (Katharopoulos et al., 2020) and Test-Time Training (TTT) (Sun et al., 2025) compress the N -width MLP into one (or several (Zhang et al., 2026)) light-weight module(s), analogous to model compression (Ba & Caruana, 2014) and knowledge distillation (Hinton et al., 2015). Therefore, from a fast-weight scaling perspective, we divide the scaling approaches of existing efficient attention methods discussed in this paper into two general categories: a) scaling by routing, and b) scaling by compression (see Fig. 1 for an illustration). Scaling by compression alone sacrifices a precise access to the original key-value pairs; whereas scaling by routing alone lacks a global summary of the full context. Although these two approaches are not mutually exclusive, most existing methods typically adopt only one of them.

In this paper, we elevate the fast-weight scaling perspective into a five-dimensional taxonomy that accommodates a broad spectrum of prior efficient-attention methods (see Tab. 1). While such a taxonomy is by no means complete, it paves the way towards a unified framework. Then, we propose to combine the two scaling strategies together and construct a tunable number of deformable fast-weight experts via **Mixture of Top- k Activations (MiTA)**, which uses a small, tunable number of *landmark queries* to look up key-value pairs. Specifically, MiTA compresses the val-

ues as *landmark values*, yielding a compact key-value set (i.e., a narrower MLP) as a shared expert, in which landmark queries serve as compressed keys and landmark values serve as compressed values; and reorganizes the original key-value pairs into deformable experts via gathering top- k key-value pairs activated by each landmark query. As shown in Fig. 2, MiTA attention concatenates these compressed key-value pairs with a routed, deformable subset of the original key-value pairs for each query.

Contributions. Contributions of the paper are summarized as follows.

1. We introduce a five-dimensional taxonomy to sort the existing efficient attention methods from a fast-weight scaling perspective.
2. We propose an efficient attention method—MiTA, which constructs deformable experts and combines the scaling by routing and scaling by compression strategies together.

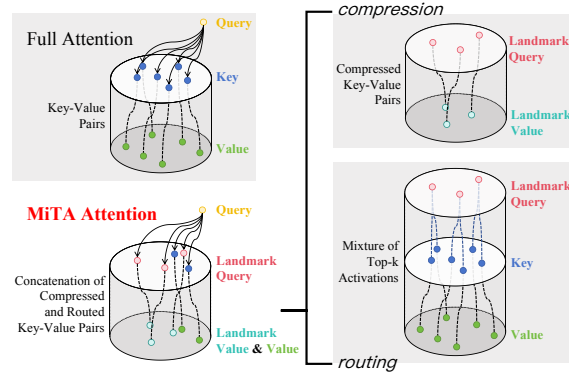


Figure 2. Illustration for our MiTA attention. In full attention, each query attends to all key-value pairs. In our MiTA attention, it attends to the concatenation of a small number of the compressed key-value pairs and a routed subset of the full key-value pairs.

2. Related Work

2.1. Fast-Weight Scaling

In contrast to slow weights (i.e., model parameters) encoding persistent knowledge, fast weights are input-conditioned and can be viewed as short-term memory (Behrouz et al., 2024), playing an important role in meta-learning (Kirsch et al., 2022) and in-context learning (Chan et al., 2022). While most scaling efforts have focused on expanding slow weights (e.g., increasing model width and depth) (Kaplan et al., 2020; Tan & Le, 2019), extending Transformers’ sequence length turns out to implicitly scale fast weights.

For instance, Test-Time Training (TTT) (Sun et al., 2025) effectively compresses the N -width two-layer fast-weight MLP in full attention into a smaller fast-weight module, likewise for linear attention (Katharopoulos et al., 2020) into a

Table 1. A five-dimensional taxonomy for efficient attention methods from the fast-weight scaling perspective. We present expert type and expert count together in the table. The analysis of this taxonomy and how MiTA fits into it can be found in Sec. 3.3.

METHOD	SCALING STRATEGY	EXPERT TYPE AND COUNT	EXPERT CONSTRUCTION	ROUTING TOPOLOGY
LINEAR ATTENTION (2020) MHSA (2026)	COMPRESSION	ONE SHARED LINEAR LAYER m LINEAR LAYERS	KERNELIZATION LOCAL PRIOR & KERNELIZATION	ALL-TO-ONE m TIMES $\frac{N}{m}$ -TO-ONE
LINFORMER (2020) PVT (2021)	COMPRESSION	ONE SHARED MLP	LEARNABLE PROJECTION	ALL-TO-ONE
AGENT ATTENTION (2024)	COMPRESSION	ONE SHARED MLP	LANDMARK TESTING	ALL-TO-ONE
TTT (2025) ViT ³ (2025)	COMPRESSION	ONE SHARED MODULE	TEST-TIME TRAINING	ALL-TO-ONE
DEFORMABLE DETR (2017) DAT (2022)	ROUTING	N MLPs ONE SHARED MLP	OFFSET PREDICTING	N TIMES ONE-TO-ONE N -TO-ONE
MoBA (2025) NSA (2025)	ROUTING COMPRESSION & ROUTING	$\frac{N}{n}$ MLPs	LOCALITY PRIOR	N -TO- $\frac{N}{n}$
SPARK ATTENTION (2025) DSA (2025)	ROUTING	N MLPs	LIGHT-WEIGHT SKIM	N TIMES ONE-TO-ONE
MiTA (OURS)	COMPRESSION & ROUTING	m MLPs	LANDMARK TESTING	N -TO- m

linear layer (Han et al., 2025) and Linformer (Wang et al., 2020) into a narrower MLP. Rather than scaling by compressing the full fast-weight MLP, MoE attention factorizes it into experts and access them sparsely via routing.

In our MiTA, the two strategies mentioned above—scaling by compressing and scaling by sparse routing—are combined. It follows the Agent attention (Han et al., 2024) to compress full key-value pairs into a smaller set (i.e., a narrower MLP), which offers a coarse yet global summary, and then leverages the top- k activations of the agent tokens to identify underlying experts, which enables precise retrieval.

2.2. MoE-Inspired Sparse Attention

Previous sparse attention methods typically focus on the design of fixed sparse patterns, such as the local window and the vertical-slash pattern (Jiang et al., 2024), or learnable ones, such as hash buckets (Kitaev et al., 2020) and k -means clusters (Roy et al., 2021). To maintain a global connectivity, shared memory (i.e., tokens that attend to and are attended by all tokens) is also introduced (Beltagy et al., 2020). However, these patterns are either task-dependent (Lai et al., 2025) or too burdensome to achieve wall-clock speedup (Dao et al., 2022).

Recently, sparse attention inspired by Mixture-of-Experts (MoE) has been applied to pre-training large language models, yielding extended effective context with reduced overhead (Lu et al., 2025; Yuan et al., 2025). Its key advantage is the routing mechanism, which enables query-aware selection within the key-value cache. Consequently, despite coarse selection granularity (e.g., evenly split, non-overlapping blocks), the sparse pattern is highly flexible and can vary across queries.

In this paper, we argue that, beyond routing, the experts themselves can also be query-aware, i.e., deformable.

Specifically, for an arbitrarily long sequence, we can build a fixed number of experts by gathering semantically related key-value pairs. A similarly motivated line of research is to replace the fixed, non-overlapping rectangular patchification in ViTs (Dosovitskiy et al., 2021) with a deformable, content-adaptive tokenization scheme (Chen et al., 2021).

2.3. Deformable Attention

More broadly, conditional computation, the idea underlying MoE, has been explored through the deformable convolution (Dai et al., 2017) in convolutional neural networks. And deformable attention was introduced for object detection by Zhu et al. (2021) and later generalized to Vision Transformers by Xia et al. (2022). Given a query, while being blind to keys, deformable attention predicts offsets relative to a small set of default positions of the keys, thereby inducing a query-aware, deformable sparse attention pattern (i.e., a fast-weight expert in our context).

Recent top- k attention methods, e.g., spark attention (You et al., 2025) and deepseek sparse attention (Liu et al., 2025), which have been applied to the training of Gemma 3n and DeepSeek-V3.2, respectively, can also be viewed as improved deformable attention. Instead of predicting the spatial positions solely from the query as in prior deformable attention, these methods locate the top- k key-value pairs by allowing each query to take a lightweight skim of the full keys, e.g., using partial features (You et al., 2025) or low precision (FP8) (Liu et al., 2025). Our MiTA retains the key advantage of the top- k attention by constructing deformable experts (i.e., sparse patterns) that depend on both queries and keys. However, unlike the deformable attention methods discussed above, we replace the per-query offset predicting (Zhu et al., 2021; Xia et al., 2022) or light-weight skim (You et al., 2025; Liu et al., 2025) with more efficient per-query routing.

3. Methods

3.1. Efficient Attention as Efficient Fast-Weight Scaling

This subsection will first present the mathematical formulation of the fast-weight MLP in full attention, and then reinterpret efficient attention as a fast-weight scaling problem. From this perspective, a taxonomy of previous efficient attention methods in terms of how they scale fast weights will be given.

Fast-weight MLP. Formally, full attention (i.e., scaled dot-product attention (Vaswani et al., 2017); SDPA) can be written as:

$$\text{SDPA}(\mathbf{q}, \mathbf{K}, \mathbf{V}) = \mathbf{V} \text{softmax} \left(\mathbf{K}^\top \mathbf{q} / \sqrt{d} \right), \quad (1)$$

where $\mathbf{q} \in \mathbb{R}^d$ is a query, and the columns of $\mathbf{K}, \mathbf{V} \in \mathbb{R}^{d \times N}$ are the keys and values, respectively. By contrast, a D -width two-layer MLP, i.e., the feed-forward network (FFN) in Transformers, is:

$$\text{MLP}_\sigma(\mathbf{x} \mid \mathbf{W}_1; \mathbf{W}_2) = \mathbf{W}_2 \sigma \left(\mathbf{W}_1^\top \mathbf{x} + \mathbf{b}_1 \right) + \mathbf{b}_2, \quad (2)$$

where $\mathbf{x} \in \mathbb{R}^d$ is the input vector, $\mathbf{W}_1, \mathbf{W}_2 \in \mathbb{R}^{d \times D}$ are weights, $\mathbf{b}_1 \in \mathbb{R}^D$ and $\mathbf{b}_2 \in \mathbb{R}^d$ are biases¹, and $\sigma(\cdot)$ is an element-wise activation function (e.g., ReLU). Note that full attention in Eq. (1) is equivalent to the following MLP:

$$\text{MLP}_{\text{exp}}(\mathbf{q} \mid \hat{\mathbf{K}}; \hat{\mathbf{V}}) = \hat{\mathbf{V}} \exp \left(\hat{\mathbf{K}}^\top \mathbf{q} \right), \quad (3)$$

$$\hat{\mathbf{K}} = \mathbf{K} / \sqrt{d}, \quad \hat{\mathbf{V}} = \mathbf{V} / \left(\exp(\hat{\mathbf{K}}^\top \mathbf{q})^\top \mathbf{1}_N \right), \quad (4)$$

where the weights are given by the scaled keys and values, the biases are all zero, and the activation function is the exponential function $\exp(\cdot)$. Therefore, we have demonstrated that full attention in Eq. (1) is equivalent to the N -width, fast-weight, two-layer MLP in Eq. (3).

Unbounded fast-weight scaling. Because the hidden dimension of the fast-weight MLP in Eq. (3) always equals to the sequence length N , the per-query overhead grows linearly with N rather than fixed as in the slow-weight MLP in Eq. (2). Thus, under the all-to-all lookup paradigm, processing N queries incurs an overall $\mathcal{O}(N^2)$ overhead. This unbounded, rapid growth prevents scaling the fast weights to arbitrarily long sequences. While this challenge is commonly described as the quadratic complexity of full attention in the efficient-attention literature, we refer to it as the *unbounded fast-weight scaling* issue.

Fast-weight scaling taxonomy. Although we have discussed relevant efficient-attention methods in Sec. 2, they have not been systematically organized into a unified fast-weight scaling perspective. We therefore propose a five-dimensional taxonomy: a) *scaling strategy*, i.e., scaling the

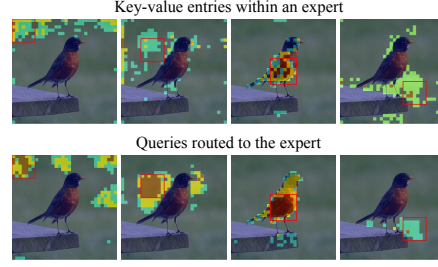


Figure 3. Visualization of experts’ gathered key-value pairs, and routed queries. The red box marks the local window from which the landmark query is obtained via average pooling. The attention heatmap (averaged over heads) indicates key-value pairs within each expert (top row) and the queries routed to it (bottom row). Notably, neither the expert’s key-value pairs nor the routed queries are confined to the local window.

N -width fast-weight MLP either by compressing it into a light-weight module or by routing each query to a subset of it; b) *expert count*, i.e., how many fast-weight experts are constructed (in particular, scaling by compression typically yields a single shared expert); c) *expert type*, i.e., what module each expert takes (e.g., an MLP or a linear layer); d) *expert construction*, i.e., how experts are formed from the key-value pairs; and e) *routing topology*, i.e., the query-expert assignment pattern. The mapping of existing methods to this taxonomy is summarized in Tab. 1. We defer the in-depth discussion until Sec. 3.3.

3.2. Mixture of Top- k Activations (MiTA) and MiTA Attention

Both scaling by compression and scaling by routing have inherent limitations. Compressing the full fast-weight MLP inevitably discards information, although the information loss can be mitigated by adopting milder compression schemes (e.g., piecewise compression (Zhang et al., 2026)) or by compressing into a more expressive module (e.g., via test-time training (Sun et al., 2025)). In contrast, scaling by routing is faithful and lossless, but it often lacks a global view of the context.

To combine the strengths of both strategies, i.e., to retain a compact global summary while enabling precise, token-level retrieval, we introduce a small set of *landmark queries* $\tilde{\mathbf{Q}} \in \mathbb{R}^{d \times m}$ with $m \ll N$. These landmark queries probe the full key-value cache and then, via cross-attention, compress it into a global fast-weight module, while simultaneously forming deformable fast-weight experts by gathering the top- k activated key-value pairs for each landmark query.

This design is motivated by two observations: a) register tokens can attend to distinct semantic regions (i.e., fast-weight experts) of an image (Darcet et al., 2024); and b) class embeddings (Strudel et al., 2021) (or object queries (Carion et al., 2020) and mask embeddings (Cheng et al., 2021)) provide a compact global summary for dense prediction

¹For notion simplicity, biases are omitted from the argument.

tasks (Wen & Li, 2024). The landmark queries can be obtained in various ways, e.g., by assigning a set of learnable slow weights or by downsampling the sequence using a convolutional module. In this work, we simply apply average pooling over uniformly spaced, equal-sized windows, as shown in Fig. 3.

Specifically, the i -th landmark query \tilde{q}_i defines the i -th expert \mathcal{E}_i as follows:

$$\mathcal{E}_i(q) = \text{SDPA}(q, K^{(i)}, V^{(i)}), i \in \{1, \dots, m\}, \quad (5)$$

$$K^{(i)} = K_{:, \mathcal{I}_i}, V^{(i)} = V_{:, \mathcal{I}_i} \in \mathbb{R}^{d \times k}, \quad (6)$$

$$\mathcal{I}_i = \text{Top}_k(K^\top \tilde{q}_i) \in \{1, \dots, N\}^k, \quad (7)$$

where q is a query routed to expert \mathcal{E}_i , and $K^{(i)}, V^{(i)}$ are the top- k key-value pairs activated by \tilde{q}_i . This approach effectively reorganizes the full key-value pairs into a **Mixture of Top- k Activations (MiTA)** (i.e., deformable fast-weight experts). We directly use \tilde{q}_i as the routing vector for expert \mathcal{E}_i . Then the routing logits from the N queries $Q \in \mathbb{R}^{d \times N}$ to the m experts is $Q^\top \tilde{Q} \in \mathbb{R}^{N \times m}$. We denote the index of the rank- j expert that a query q is routed to under these logits as $e_j(q) \in \{1, \dots, m\}$.

Moreover, regarding the compressed global module, which can be viewed as a shared expert (Dai et al., 2024), we extract a set of *landmark values* $\tilde{V} \in \mathbb{R}^{d \times m}$ via cross-attention using the landmark queries. The i -th landmark value corresponding to the i -th landmark query is given by:

$$\tilde{v}_i = \text{SDPA}(\tilde{q}_i, K, V). \quad (8)$$

Then, the shared expert $\tilde{\mathcal{E}}$ is defined as:

$$\tilde{\mathcal{E}}(q) = \text{SDPA}(q, \tilde{Q}, \tilde{V}), \quad (9)$$

where landmark queries act as the keys in standard attention. Note that the computations in Eqs. (7) and (8) can be largely shared, where landmark queries are both querying the keys. Likewise, Eq. (9) can reuse the routing logits. Therefore, the computations of the two scaling strategies are tightly coupled in our MiTA.

Rather than a straightforward MoE implementation, which isolates the computation within each expert and aggregates their outputs via a weighted sum, we concatenate the experts as a single standard attention, as shown in Fig. 2. Therefore, our **MiTA attention** can be written as:

$$\text{MiTA}(q) = V^* \text{softmax}(K^{*\top} q / \sqrt{d}), \quad (10)$$

$$K^* = [\tilde{Q}, K^{(e_1(q))}, \dots, K^{(e_s(q))}], \quad (11)$$

$$V^* = [\tilde{V}, V^{(e_1(q))}, \dots, V^{(e_s(q))}], \quad (12)$$

where $K^*, V^* \in \mathbb{R}^{d \times (m+ks)}$ are the key-value pairs that a query q attends to in our MiTA, and s is the number of

Algorithm 1 Mixture of Top- k Activations (MiTA) attention

Require: Query, key, value matrices $Q, K, V \in \mathbb{R}^{d \times N}$; the number of landmark queries m ; the width of each fast-weight expert k ; each query is routed to the shared expert and to $s = 1$ additional expert.

- 1: // Obtain landmark queries
- 2: $\tilde{Q} = \text{AdaptiveAvgPool}(Q, \text{output_size}=m)$ // [d, m]
- 3: // Lookup key-value pairs via landmark queries
- 4: $S^{\text{kv}} = K^\top \tilde{Q} / \sqrt{d}$ // [N, m]
- 5: // Gather the top- k activated key-value pairs
- 6: $\mathcal{I}^{\text{kv}} = \text{Flatten}(\text{TopK}(S^{\text{kv}\top}, k, \text{dim}=1))$ // [m*k]
- 7: $K^{\text{expt}}, V^{\text{expt}} = K[:, \mathcal{I}^{\text{kv}}], V[:, \mathcal{I}^{\text{kv}}]$ // [d, m*k]
- 8: // Obtain landmark values (construct the shared expert)
- 9: $\tilde{V} = V \text{softmax}(S^{\text{kv}})$ // [d, m]
- 10: // Always route queries to the shared expert
- 11: $O^{\text{share}} = \text{FlashAttention}(Q, \tilde{Q}, \tilde{V})$
- 12: // Sparsely route queries to other experts
- 13: $\mathcal{I}^{\text{expt}} = \text{ArgSort}(\text{ArgMax}(\tilde{Q}^\top Q, \text{dim}=0))$ // [N]
- 14: $O^{\text{expt}} = \text{FlashAttention}(Q[:, \mathcal{I}^{\text{expt}}], K^{\text{expt}}, V^{\text{expt}}, \text{cu_seqlens_q} = \text{CumSum}(\text{BinCount}(\mathcal{I}^{\text{expt}})), \text{cu_seqlens_k} = [0, k, 2k, \dots, mk])$
- 15: // Combine results via online softmax
- 16: $O = \text{Combine}(O^{\text{share}}, O^{\text{expt}})$ // [d, N]
- 17: **return** O

routed experts (not counting the shared expert) per query. Note that, in practice, the attention can still be computed expert-wise and then combined, thanks to the online softmax (Milakov & Gimelshein, 2018), as in the FlashAttention (Dao, 2024) implementation of Lu et al. (2025).

Our MiTA attention in Eq. (10) can be summarized as follows: a) obtain landmark queries \tilde{Q} from queries Q via average pooling; b) use landmark queries \tilde{Q} to query keys K , and hence b.1) construct deformable fast-weight experts (i.e., $K^{(i)}, V^{(i)}$) via the top- k activations and b.2) extract landmark values \tilde{V} via cross-attention; c) route queries Q sparsely to the experts, while keeping the shared expert (i.e., \tilde{Q}, \tilde{V}) always active.

Implementations. In particular, we implement our MiTA attention in Eq. (10) with $s = 1$, i.e., each query is dispatched to the shared expert and exactly one additional fast-weight expert. Algorithm 1 presents the pseudocode of our MiTA, which matches queries to experts by sorting queries according to their expert assignments $e_1(q)$. The more general but more complex case with $s > 1$ can be implemented via MoBA (Lu et al., 2025).

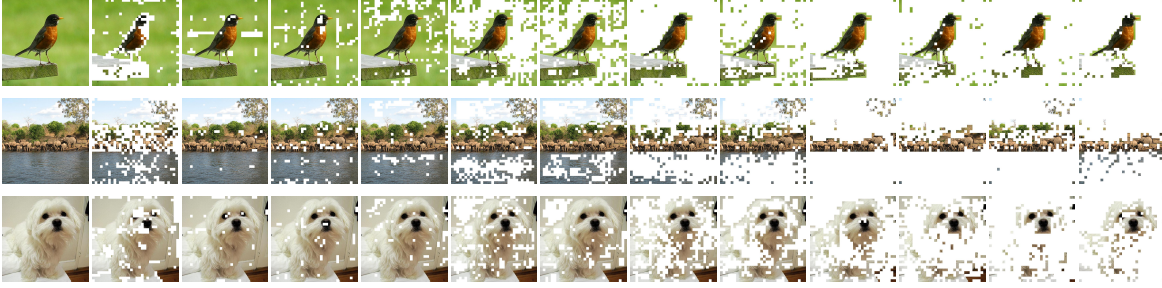


Figure 4. The token pruning effect of our MiTA. Each row visualizes, for each layer, the positions of key-value pairs (aggregated over heads) selected as experts; the leftmost image shows the original input. In later layers, most tokens are effectively “pruned” (i.e., not selected as experts), and attention concentrates on class-relevant regions. The examples are sampled from the ImageNet-1K training set.

Complexity analysis. The computational complexity of Eq. (10) alone is $\mathcal{O}(N(m + ks))$, whereas full attention incurs $\mathcal{O}(N^2)$; in practice, $N \gg m + ks$. Moreover, the cross-attention operations in Eqs. (8) and (10) can be accelerated by the FlashAttention family (Dao et al., 2022; Dao, 2024; Shah et al., 2024). The primary bottleneck comes from the gather operation in Eq. (6): it entails irregular (random) memory access and is therefore likely to make the overall pipeline memory-bound. Nonetheless, Deepseek Sparse Attention suggests that this limitation can be mitigated via its optimized implementation (Liu et al., 2025, Fig. 3). Moreover, compared with prior top- k attention (You et al., 2025; Liu et al., 2025), which instantiate a private fast-weight expert per query (N experts in total) without routing, our MiTA uses a *fixed* (or tunable) number (i.e., m) of fast-weight experts, which is more hardware-friendly, at the cost of an N -to- m routing step, which is explored in MoBA (Lu et al., 2025). In summary, our MiTA is composed of operations that have already been validated by engineering-centric prior work, allowing us to focus on the fast-weight scaling perspective, its taxonomy, the combination of the two scaling strategies, and the construction of a fixed number of deformable fast-weight experts.

3.3. Fit MiTA into the Fast-Weight Scaling Taxonomy

While MiTA is primarily motivated by combining the two complementary scaling strategies, we find that it also enriches the remaining dimensions of the fast-weight scaling taxonomy proposed in Sec. 3.1. In this subsection, we briefly discuss each of them.

Expert construction. This determines the expert type and count, and hence the routing topology. An important consideration is what the construction conditions on: the query only (e.g., DAT (Xia et al., 2022)), the key-value pairs only (e.g., MHLA (Zhang et al., 2026)), or both. When using key-value pairs, the construction may be content-dependent and thus deformable (e.g., DSA (Liu et al., 2025)) or merely position-driven (e.g., MoBA (Lu et al., 2025)). While most efficient attention methods condition on one of the above aspects, our MiTA conditions on all of them by extract-

ing landmark queries from queries and then probe the full key-value pool. We illustrate a resulting benefit of this design in Fig. 4.

Expert type and count. Expert types can be categorized in three classes: a) linear layers (e.g., linear attention), b) MLPs (e.g., sparse attention and PVT (Wang et al., 2021)), c) arbitrary modules (e.g., test-time training). The scaling-by-routing branch of MiTA is restricted to MLP experts, as in sparse attention. However, the scaling-by-compression branch can be implemented through test-time training, thereby generalizing to other modules. As for the expert count, MiTA controls it via the hyperparameter m . One benefit of this design has been discussed in complexity analysis. Another benefit is that it enables expert composition through the routing mechanism.

Routing topology. In MiTA, there is a fixed number of base sparse patterns (i.e., fast-weight experts), which is a substantial improvement over efficient attention with only a single expert, yet is still far fewer than prior top- k attention (You et al., 2025; Liu et al., 2025), which has N experts. Crucially, to compensate for this limitation, MiTA can resort to the routing mechanism to compose these m base experts: when routing each query to s experts, the number of effective sparse patterns (i.e., the combinatorial number of experts) is $\binom{m}{s}$. In particular, when $s = 1$, routing may degenerate into clustering, restricting information flow within each cluster. We thus quantify the positional overlap between an expert’s gathered key-value pairs and the queries routed to it. As shown in Fig. 5, the overlap remains consistently modest across layers, suggesting that MiTA with $s = 1$ performs routing rather than hard clustering.

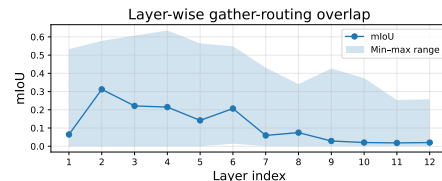


Figure 5. The layer-wise positional overlap (between the key-value pairs gathered by an expert and the queries routed to it) is quantified by mIoU, averaged over experts and heads.

4. Experiments

To verify the efficacy of our proposed MiTA, we conduct image classification experiments on ImageNet-1K (Deng et al., 2009) and semantic segmentation experiments on ADE20K (Zhou et al., 2019). Moreover, we assess the long sequence modeling capability of MiTA on the Long Range Arena benchmark (Tay et al., 2021), and report its wall-clock inference speed against full attention on extremely long sequences. Furthermore, we study the robustness (more precisely, generalization) of MiTA under varying expert count (m) and width (k) when the model is fixed after training. Additionally, we explore a broader problem of interest: the *algorithmic generalization* of Transformers across attention mechanisms.

Table 2. Results on ImageNet-1K training from-scratch. Δ : using depth-wise convolution (DWC); Bias : using Agent bias; $^{25+25}$: MiTA attention with $m+k = 25+25$ (and similarly for other superscripts).

METHODS	ACC.	MACS	# PARAMS
DEiT-T (2021)	72.2	1.2G	5.7M
FLATTEN-DEiT-T Δ (2023)	74.1	1.1G	6.1M
AGENT-DEiT-T $^{49,\Delta,\text{Bias}}$ (2024)	74.9	1.2G	6.0M
AGENT-ViT-T 25	68.9	1.1G	5.7M
AGENT-ViT-T 49	71.7	1.1G	5.7M
ViT-T	73.1	1.2G	5.7M
DEiT-T Δ	75.8	1.2G	5.7M
MiTA-ViT-T $^{25+25}$	72.9	1.1G	5.7M
MiTA-ViT-T $^{25+25,\Delta}$	<u>74.8</u>	1.1G	5.7M
DEiT-S (2021)	79.8	4.6G	22.1M
AGENT-DEiT-S $^{49,\Delta,\text{Bias}}$ (2024)	80.5	4.4G	22.7M
AGENT-ViT-S 49	76.0	4.4G	22.1M
ViT-S	77.9	4.6G	22.1M
ViT-S Δ	79.5	4.6G	22.1M
MiTA-ViT-S $^{25+25}$	78.5	4.4G	22.1M
MiTA-ViT-S $^{25+25,\Delta}$	<u>79.3</u>	4.4G	22.1M

4.1. Classification Evaluated on ImageNet-1K

We integrate MiTA attention into the standard Vision Transformer (ViT) (Dosovitskiy et al., 2021) to obtain MiTA-ViT, and compare its top-1 accuracy on the validation set against ViT, Agent-ViT, and Linear-ViT, where Agent-ViT implements Agent attention (Han et al., 2024), which can be viewed as the case $k = 0$ of MiTA, while Linear-ViT adopts linear attention (Katharopoulos et al., 2020).

Implementation details. For fair comparison and to isolate the effect of the attention operator, unless otherwise specified, we remove auxiliary components used in prior designs, such as depth-wise convolutions (DWC) and the Agent bias (Han et al., 2024). By default, we set $s = 1$, $m = 25$, and $k = 25$ for our MiTA. Under this setting, each query attends to $m + k = 50$ key-value pairs. And the maximum number of distinct key-value pairs that are

potentially attentive is $m \times k = 625$, which exceeds the typical sequence length used in ImageNet-1K classification.

Results. In Tab. 2, the results in black are obtained by training models from scratch on the ImageNet-1K training set with an identical training recipe, whereas results in gray are taken directly from the original papers. Note that MiTA-ViTs perform on par with ViTs, both with and without DWC augmentation. In contrast, Agent-ViT with pure Agent attention falls behind by a large margin, as it scales fast weights by compressing them only. In Tab. 3, we fine-tune the released ViT models (Wightman, 2019) pretrained on ImageNet-21K, on ImageNet-1K, while replacing the attention mechanism. It shows that MiTA-ViT performs consistently better than Agent-ViT and is comparable to ViT (at least being more compatible with full attention’s pretrained weights), and Agent-ViT struggles to narrow this gap even by raising m from 49 to 64.

Table 3. ImageNet-1K finetune results. ViTs are pretrained on ImageNet-21K, and then finetuned with the attention mechanism replaced. We finetune all models for 50 epochs.

METHOD	TINY	SMALL	BASE	LARGE
ViT	76.9	81.2	84.4	85.9
LINEAR-ViT	73.2	78.6	79.8	80.7
AGENT-ViT 49	74.6	79.7	81.4	83.5
AGENT-ViT 64	74.7	79.9	81.5	83.8
MiTA-ViT $^{25+25}$	<u>75.6</u>	<u>80.9</u>	<u>82.8</u>	<u>85.5</u>

Table 4. Segmentation results on ADE20K. The image resolution and patch size are fixed to 512×512 and 16×16 for all models, yielding a sequence length of 1,024.

METHOD	mIoU (SS/MS)	# PARAMS
SEG-T-MASK (2021)	38.1 / 38.8	7M
SEG-T-MiTA $^{36+36}$	38.8 / 39.6	6M
SEG-S-MASK	45.3 / 46.9	27M
SEG-S-MiTA $^{36+36}$	45.5 / 46.8	24M
SEG-B-MASK	48.8 / 50.0	106M
SEG-B-MiTA $^{36+36}$	48.8 / 49.9	96M

4.2. Semantic Segmentation

For semantic segmentation, we apply our MiTA attention in the decoder while keeping the encoder as a ViT pretrained on ImageNet-21K, and fairly compare with Segmenter (Strudel et al., 2021), whose decoder is a mask transformer built with full attention. The results are reported in Tab. 4, showing that our MiTA attention can effectively handle long sequences and dense prediction, achieving performance comparable to, or even surpassing, full attention. We attribute this in part to the faster convergence of deformable attention on long sequences, as suggested by (Zhu et al., 2021).

Table 5. Long Range Arena (LRA) benchmark results. We report accuracy for each individual task and average accuracy across all tasks. We use $m = 16$, $k = 16$ for MiTA, and $m = 32$ for Agent attention. For fair comparison, all methods are implemented *without* DWC.

METHOD	LISTOPS (2K)	TEXT (4K)	RETRIEVAL (4K)	IMAGE (1K)	PATHFINDER (1K)	AVG
STANDARD	37.10	65.02	80.11	40.45	74.16	59.37
REFORMER (2020)	19.05	64.88	78.64	43.29	69.36	55.04
LINFORMER (2020)	37.25	55.91	79.37	37.84	67.60	55.59
PERFORMER (2021)	18.80	63.81	78.62	37.07	69.87	53.63
NYSTRÖMFORMER (2021)	18.70	65.64	80.37	39.70	73.46	55.57
AGENT ATTENTION (2024)	37.15	64.41	<u>79.75</u>	38.17	72.57	58.40
MiTA ATTENTION	<u>37.20</u>	63.24	79.37	<u>42.19</u>	74.32	<u>59.26</u>

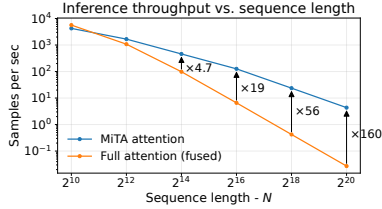


Figure 6. Inference throughput. The results are measured on an NVIDIA RTX 4090 (24GB) with a three-layer Transformer (model dimension 128) and a synthetic dataset of 10K samples.

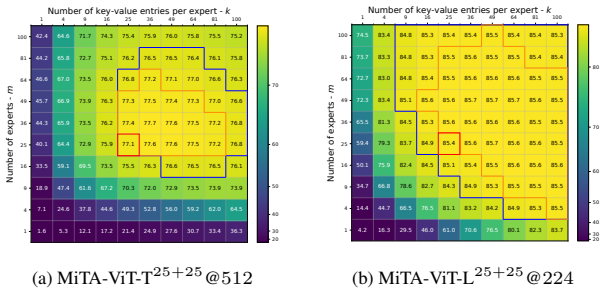


Figure 7. Generalization across m and k . Red boxes mark the training expert count and width, i.e., (m, k) , and the corresponding baseline top-1 accuracy on ImageNet-1K. Orange boxes mark the inference (m, k) that exceed the baseline accuracy, while blue boxes mark the inference (m, k) that achieve 99% of the baseline accuracy. Both models are finetuned from pretrained ViTs. Best viewed when zoomed in.

4.3. Long Sequence Modeling and Wall-Clock Speed

To assess the performance of MiTA attention on long-sequence tasks, we use the Long Range Arena (LRA) benchmark (Tay et al., 2021) and compare MiTA with other efficient attention mechanisms. The results are reported in Tab. 5. Similar to full attention, MiTA performs consistently well across the five tasks.

We also report MiTA’s inference throughput (samples per second) in Fig. 6 across sequence lengths up to 2^{20} ($\approx 1M$), and compare it against full attention implemented with PyTorch’s fused kernel. While the full context is compressed into an m -width fast-weight MLP, we set $m \times k \approx N/4$, which allows MiTA attention to precisely access the most important $1/4$ tokens. Note that no custom CUDA kernel optimization is applied to the top- k gathering operation here.

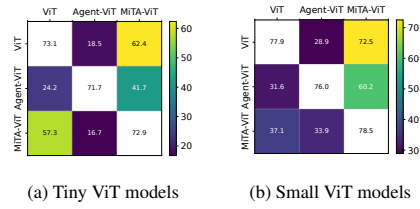


Figure 8. Checkpoint swapping. The y-axis indexes the target model (i.e., the attention mechanism used at inference time), while the x-axis indexes the source weights (i.e., the attention mechanism under which the weights were trained). We omit the diagonal entries from the heatmap since they are not of interest. All models are trained from scratch on ImageNet-1K. Best viewed when zoomed in.

4.4. Generalization Across Model Designs

Unlike length extrapolation (Zhao et al., 2024), we study MiTA’s generalization across the expert count m and expert width k . Fig. 7 reports MiTA’s inference performance when interpolating within, and extrapolating beyond, the training (m, k) . We find that MiTA is notably more robust to increasing m and k than to decreasing them. A promising direction is to train MiTA with smaller m and k for efficiency, and then scale up m or k to obtain substantial gains.

Additionally, we explore the generalization of Transformers across attention mechanisms. Given a checkpoint that learns slow weights under an attention mechanism, we swap in a different attention mechanism at inference time and evaluate the resulting performance. The results are reported in Fig. 8. We observe that checkpoints trained with MiTA transfer better to full attention.

5. Conclusion

We adopted a fast-weight scaling as a unifying perspective for efficient attention methods and introduced a five-dimensional taxonomy. Moreover, we proposed an efficient attention method, which is termed as the mixture of top- k activations (MiTA), to bridge the two scaling strategies—scaling by routing and scaling by compression—and construct a tunable number of deformable experts. We provided a detailed analysis of both the taxonomy and our MiTA, suggesting a promising path for developing efficient attention from the fast-weight scaling perspective.

References

- Ba, J. and Caruana, R. Do deep nets really need to be deep? In *NeurIPS*, 2014.
- Behrouz, A., Zhong, P., and Mirrokni, V. Titans: Learning to memorize at test time. *arXiv preprint arXiv:2501.00663*, 2024.
- Beltagy, I., Peters, M. E., and Cohan, A. Longformer: The long-document transformer. *arXiv preprint arXiv:2004.05150*, 2020.
- Bietti, A., Cabannes, V., Bouchacourt, D., Jégou, H., and Bottou, L. Birth of a transformer: A memory viewpoint. In *NeurIPS*, 2023.
- Cai, S., Yang, C., Zhang, L., Guo, Y., Xiao, J., Yang, Z., Xu, Y., Yang, Z., Yuille, A., Guibas, L., et al. Mixture of contexts for long video generation. *arXiv preprint arXiv:2508.21058*, 2025.
- Carion, N., Massa, F., Synnaeve, G., Usunier, N., Kirillov, A., and Zagoruyko, S. End-to-end object detection with transformers. In *ECCV*, 2020.
- Chan, S. C. Y., Santoro, A., Lampinen, A. K., Wang, J. X., Singh, A. K., Richemond, P. H., McClelland, J. L., and Hill, F. Data distributional properties drive emergent in-context learning in transformers. In *NeurIPS*, 2022.
- Chen, Z., Zhu, Y., Zhao, C., Hu, G., Zeng, W., Wang, J., and Tang, M. DPT: Deformable patch-based transformer for visual recognition. In *ACMMM*, 2021.
- Cheng, B., Schwing, A., and Kirillov, A. Per-pixel classification is not all you need for semantic segmentation. In *NeurIPS*, 2021.
- Cheng, H., Zhang, M., and Shi, J. Q. A survey on deep neural network pruning: Taxonomy, comparison, analysis, and recommendations. *TPAMI*, 46(12):10558–10578, 2024.
- Choromanski, K. M., Likhoshesterov, V., Dohan, D., Song, X., Gane, A., Sarlós, T., Hawkins, P., Davis, J. Q., Mohiuddin, A., Kaiser, L., Belanger, D. B., Colwell, L. J., and Weller, A. Rethinking attention with performers. In *ICLR*, 2021.
- Dai, D., Deng, C., Zhao, C., Xu, R. X., Gao, H., Chen, D., Li, J., Zeng, W., Yu, X., Wu, Y., Xie, Z., Li, Y. K., Huang, P., Luo, F., Ruan, C., Sui, Z., and Liang, W. DeepSeekMoE: Towards ultimate expert specialization in mixture-of-experts language models. In *ACL*, 2024.
- Dai, J., Qi, H., Xiong, Y., Li, Y., Zhang, G., Hu, H., and Wei, Y. Deformable convolutional networks. In *ICCV*, 2017.
- Dao, T. FlashAttention-2: Faster attention with better parallelism and work partitioning. In *ICLR*, 2024.
- Dao, T., Fu, D., Ermon, S., Rudra, A., and Ré, C. FlashAttention: Fast and memory-efficient exact attention with io-awareness. *NeurIPS*, 2022.
- Darcet, T., Oquab, M., Mairal, J., and Bojanowski, P. Vision transformers need registers. In *ICLR*, 2024.
- Deng, J., Dong, W., Socher, R., Li, L., Li, K., and Fei-Fei, L. ImageNet: A large-scale hierarchical image database. In *CVPR*, 2009.
- Dosovitskiy, A., Beyer, L., Kolesnikov, A., Weissenborn, D., Zhai, X., Unterthiner, T., Dehghani, M., Minderer, M., Heigold, G., Gelly, S., Uszkoreit, J., and Houtsby, N. An image is worth 16x16 words: Transformers for image recognition at scale. In *ICLR*, 2021.
- Geva, M., Schuster, R., Berant, J., and Levy, O. Transformer feed-forward layers are key-value memories. In *EMNLP*, 2021.
- Han, D., Pan, X., Han, Y., Song, S., and Huang, G. Flatten transformer: Vision transformer using focused linear attention. In *ICCV*, 2023.
- Han, D., Ye, T., Han, Y., Xia, Z., Pan, S., Wan, P., Song, S., and Huang, G. Agent attention: On the integration of softmax and linear attention. In *ECCV*, 2024.
- Han, D., Li, Y., Li, T., Cao, Z., Wang, Z., Song, J., Cheng, Y., Zheng, B., and Huang, G. ViT³: Unlocking test-time training in vision. *arXiv preprint arXiv:2512.01643*, 2025.
- Hinton, G., Vinyals, O., and Dean, J. Distilling the knowledge in a neural network. *arXiv preprint arXiv:1503.02531*, 2015.
- Jia, W., Lu, Y., Huang, M., Wang, H., Huang, B., Chen, N., Liu, M., Jiang, J., and Mao, Z. MoGA: Mixture-of-groups attention for end-to-end long video generation. *arXiv preprint arXiv:2510.18692*, 2025.
- Jiang, H., Li, Y., Zhang, C., Wu, Q., Luo, X., Ahn, S., Han, Z., Abdi, A. H., Li, D., Lin, C., Yang, Y., and Qiu, L. MInference 1.0: Accelerating pre-filling for long-context llms via dynamic sparse attention. In *NeurIPS*, 2024.
- Kaplan, J., McCandlish, S., Henighan, T., Brown, T. B., Chess, B., Child, R., Gray, S., Radford, A., Wu, J., and Amodei, D. Scaling laws for neural language models. *arXiv preprint arXiv:2001.08361*, 2020.
- Katharopoulos, A., Vyas, A., Pappas, N., and Fleuret, F. Transformers are rnns: Fast autoregressive transformers with linear attention. In *ICML*, 2020.
- Kirsch, L., Harrison, J., Sohl-Dickstein, J., and Metz, L. General-purpose in-context learning by meta-learning transformers. *arXiv preprint arXiv:2212.04458*, 2022.
- Kitaev, N., Kaiser, L., and Levskaya, A. Reformer: The efficient transformer. In *ICLR*, 2020.
- Lai, X., Lu, J., Luo, Y., Ma, Y., and Zhou, X. FlexPrefill: A context-aware sparse attention mechanism for efficient long-sequence inference. In *ICLR*, 2025.
- Liu, A., Mei, A., Lin, B., Xue, B., Wang, B., Xu, B., Wu, B., Zhang, B., Lin, C., Dong, C., et al. DeepSeek-v3. 2: Pushing the frontier of open large language models. *arXiv preprint arXiv:2512.02556*, 2025.
- Lu, E., Jiang, Z., Liu, J., Du, Y., Jiang, T., Hong, C., Liu, S., He, W., Yuan, E., Wang, Y., et al. MoBA: Mixture of block attention for long-context llms. *arXiv preprint arXiv:2502.13189*, 2025.
- Milakov, M. and Gimelshein, N. Online normalizer calculation for softmax. *arXiv preprint arXiv:1805.02867*, 2018.
- Riquelme, C., Puigcerver, J., Mustafa, B., Neumann, M., Jenatton, R., Pinto, A. S., Keysers, D., and Houtsby, N. Scaling vision with sparse mixture of experts. In *NeurIPS*, 2021.

- Roy, A., Saffar, M., Vaswani, A., and Grangier, D. Efficient content-based sparse attention with routing transformers. *Trans. Assoc. Comput. Linguistics*, 9:53–68, 2021.
- Schlag, I., Irie, K., and Schmidhuber, J. Linear transformers are secretly fast weight programmers. In *ICML*, 2021.
- Schmidhuber, J. Learning to control fast-weight memories: An alternative to dynamic recurrent networks. *Neural Comput.*, 4 (1):131–139, 1992.
- Shah, J., Bikshandi, G., Zhang, Y., Thakkar, V., Ramani, P., and Dao, T. FlashAttention-3: Fast and accurate attention with asynchrony and low-precision. In *NeurIPS*, 2024.
- Strudel, R., Garcia, R., Laptev, I., and Schmid, C. Segmenter: Transformer for semantic segmentation. In *ICCV*, 2021.
- Sun, Y., Li, X., Dalal, K., Xu, J., Vikram, A., Zhang, G., Dubois, Y., Chen, X., Wang, X., Koyejo, S., Hashimoto, T., and Guestrin, C. Learning to (learn at test time): Rnns with expressive hidden states. In *ICML*, 2025.
- Tan, M. and Le, Q. EfficientNet: Rethinking model scaling for convolutional neural networks. In *ICML*, 2019.
- Tay, Y., Dehghani, M., Abnar, S., Shen, Y., Bahri, D., Pham, P., Rao, J., Yang, L., Ruder, S., and Metzler, D. Long range arena : A benchmark for efficient transformers. In *ICLR*, 2021.
- Touvron, H., Cord, M., Douze, M., Massa, F., Sablayrolles, A., and Jégou, H. Training data-efficient image transformers & distillation through attention. In *ICML*, 2021.
- Vaswani, A., Shazeer, N., Parmar, N., Uszkoreit, J., Jones, L., Gomez, A. N., Kaiser, Ł., and Polosukhin, I. Attention is all you need. *NeurIPS*, 2017.
- Wang, S., Li, B. Z., Khabsa, M., Fang, H., and Ma, H. Linformer: Self-attention with linear complexity. *arXiv preprint arXiv:2006.04768*, 2020.
- Wang, W., Xie, E., Li, X., Fan, D., Song, K., Liang, D., Lu, T., Luo, P., and Shao, L. Pyramid vision transformer: A versatile backbone for dense prediction without convolutions. In *ICCV*, 2021.
- Wen, Q. and Li, C.-G. Rethinking decoders for transformer-based semantic segmentation: A compression perspective. In *NeurIPS*, 2024.
- Wightman, R. Pytorch image models. <https://github.com/rwightman/pytorch-image-models>, 2019.
- Wu, B., Chen, M., Luo, X., Yan, S., Yu, Q., Xia, F., Zhang, T., Zhan, H., Zhong, Z., Zhou, X., et al. Parallel loop transformer for efficient test-time computation scaling. *arXiv preprint arXiv:2510.24824*, 2025a.
- Wu, J., Hou, L., Yang, H., Tao, X., Tian, Y., Wan, P., Zhang, D., and Tong, Y. VMoBA: Mixture-of-block attention for video diffusion models. *arXiv preprint arXiv:2506.23858*, 2025b.
- Xia, Z., Pan, X., Song, S., Li, L. E., and Huang, G. Vision transformer with deformable attention. In *CVPR*, 2022.
- Xiong, Y., Zeng, Z., Chakraborty, R., Tan, M., Fung, G., Li, Y., and Singh, V. Nyströmformer: A nyström-based algorithm for approximating self-attention. In *AAAI*, 2021.
- You, C., Wu, K., Jia, Z., Chen, L., Bhojanapalli, S., Guo, J., Evci, U., Wassenberg, J., Netrapalli, P., Willcock, J. J., et al. Spark transformer: Reactivating sparsity in ffn and attention. In *NeurIPS*, 2025.
- Yuan, J., Gao, H., Dai, D., Luo, J., Zhao, L., Zhang, Z., Xie, Z., Wei, Y., Wang, L., Xiao, Z., et al. Native sparse attention: Hardware-aligned and natively trainable sparse attention. In *ACL*, 2025.
- Zhang, J., Su, R., Liu, C., Wei, J., Wang, Z., Zhang, P., Wang, H., Jiang, H., Huang, H., Xiang, C., Xi, H., Yang, S., Li, X., Hu, Y., Fu, T., Zhao, T., Zhang, Y., Jiang, Y., Chen, C., Jiang, K., Chen, H., Zhao, M., Xu, X., Zhu, J., and Chen, J. A survey of efficient attention methods: Hardware-efficient, sparse, compact, and linear attention. 2025.
- Zhang, K., Huang, Y., Deng, Y., Yu, J., Chen, J., Ling, H., Xie, E., and Zhou, D. MHLA: Restoring expressivity of linear attention via token-level multi-head. *arXiv preprint arXiv:2601.07832*, 2026.
- Zhao, L., Feng, X., Feng, X., Zhong, W., Xu, D., Yang, Q., Liu, H., Qin, B., and Liu, T. Length extrapolation of transformers: A survey from the perspective of positional encoding. In *EMNLP*, 2024.
- Zhou, B., Zhao, H., Puig, X., Xiao, T., Fidler, S., Barriuso, A., and Torralba, A. Semantic understanding of scenes through the ade20k dataset. *IJCV*, 127(3):302–321, 2019.
- Zhu, X., Su, W., Lu, L., Li, B., Wang, X., and Dai, J. Deformable DETR: deformable transformers for end-to-end object detection. In *ICLR*, 2021.