DeBiFormer: Vision Transformer with Deformable Agent Bi-level Routing Attention

NguyenHuu BaoLong^{*1}, Chenyu Zhang^{*1}, Yuzhi Shi¹, Tsubasa Hirakawa¹, Takayoshi Yamashita¹, Tohgoroh Matsui¹, and Hironobu Fujiyoshi¹

Chubu University, Kasugai, Japan maclong01@gmail.com {zhang1121, shi, hirakawa}@mprg.cs.chubu.ac.jp {takayoshi, matsui, fujiyoshi}@isc.chubu.ac.jp

Abstract. Vision Transformers with various attention modules have demonstrated superior performance on vision tasks. While using sparsityadaptive attention, such as in DAT, has yielded strong results in image classification, the key-value pairs selected by deformable points lack semantic relevance when fine-tuning for semantic segmentation tasks. The query-aware sparsity attention in BiFormer seeks to focus each query on top-k routed regions. However, during attention calculation, the selected key-value pairs are influenced by too many irrelevant queries, reducing attention on the more important ones. To address these issues, we propose the Deformable Bi-level Routing Attention (DBRA) module, which optimizes the selection of key-value pairs using agent queries and enhances the interpretability of queries in attention maps. Based on this, we introduce the Deformable Bi-level Routing Attention Transformer (DeBiFormer), a novel general-purpose vision transformer built with the DBRA module. DeBiFormer has been validated on various computer vision tasks, including image classification, object detection, and semantic segmentation, providing strong evidence of its effectiveness. Code is available at https://github.com/maclong01/DeBiFormer

Keywords: Vision Transformer, Self-Attention Mechanism, Image Recognition

1 Introduction

The Vision Transformer has recently demonstrated significant promise in the realm of computer vision [15,29,44]. It can capture long-range dependency in data [29,41], and is almost leading to a convolution-free model more flexible for fitting tons of data [44]. In addition, it enjoys high parallelism, which benefits training and inference for large models [11,41]. The computer vision community has observed a surge in the adoption and development of Vision Transformers [1,14,15,29,44,45].

^{*} Equal contribution.



Fig. 1. Vanilla attention and its sparse variants are illustrated in the diagram: (a) vanilla attention functions globally, leading to increased computational complexity and substantial memory footprint. (b)-(c) Multiple strategies strive to reduce complexity by incorporating sparse attention with various handcrafted patterns, such as local window [50] and dilated window [45,40,24]. (d) Deformable attention [47] facilitates image-adaptive sparsity by deforming the regular grid. (e) Bi-level routing attention [56] begins by searching for top-k (k = 3 in this case) relevant regions and subsequently attends to the union of these regions. (f) In our approach, we realize bi-level routing attention, where the initial step involves searching for top-k (k = 1 in this case) relevant regions. Subsequently, attention is directed to the union of these regions by deforming regular grid attendance via top-k relevant regions.

To improve attention, numerous pieces of research have used thoughtfully crafted, efficient attention patterns in which each query selectively focused by a smaller portion of key-value pairs. As shown in Figure 1, among the various representation approaches, some include local windows [50] and dilated windows [45,40,24]. In addition, some research has taken a different path through sparsity adaptation to data in their methodology, as demonstrated in the works of [5,47]. However, despite the varying strategies for merging or selecting key and value tokens. These tokens are not semantic for queries. With this approach, when applied to other downstream tasks for pretrained ViT [41] and DETR [1]. Queries do not originate from semantic-region key-value pairs. Consequently, compelling all queries to focus on insufficient sets of tokens may not yield the most optimal results. Recently, with the dynamic query-aware sparsity attention mechanism, queries are focused by the most dynamic semantically key-value pairs, which is referred to as bi-level routing attention [56]. However, in this approach, queries

are handled by semantic key-value pairs instead of originating from detailed regions, which may not yield the most optimal results in all cases. In addition, when calculating the attention, these keys and values selected for all queries are influenced by too many less relevant queries, resulting in a decrease of attention for important queries, which has a significant impact when performing segmentations.[13,25].

To make the attention for queries more efficient, we propose the Deformable Bi-level Routing Attention (DBRA), an attention-in-attention architecture for visual recognition. During the process in DBRA, the first problem is how to locate deformable points. We use the observation in [47] that attention has an offset network that takes as input query features and generates corresponding offsets for all reference points. Thus, candidate deformable points are shifted towards important regions with high flexibility and efficiency to capture more informative features. The second problem is how to aggregate information from semantically relevant key-value pairs, and then broadcast the information back to queries. Therefore, we propose an attention-in-attention architecture that shifted toward deformable points as shown above acts as the agent for the queries. As the keyvalue pairs are selected for deformable points, we use the observation in [56] to select a small portion of the most semantically relevant key-value pairs that a region only needs by focusing on the top-k routed regions. Then, with the semantically relevant key-value pairs selected, we first apply a token-to-token attention with deformable points queries. And then, we apply a second token-to-token attention to broadcast the information back to queries, in which deformable points as key-value pairs are designed to represent the most important points in a portion of semantic regions.

To summarize, our contributions are as follows:

- 1. We propose Deformable Bi-level Routing Attention (DBRA), an attention-in-attention architecture for visual recognition, where data-dependent attention patterns are obtained flexibly and semantically.
- 2. By utilizing the DBRA module, we propose a novel backbone, called DeBiFormer, which has a stronger recognition ability based on the visualization results of the attention heat map.
- 3. Extensive experiments on ImageNet [35], ADE20K [55], and COCO [17] demonstrate that our model consistently outperforms other competitive baselines.

2 Related Work

2.1 Vision Transformers

The Transformer-based backbone incorporates channel-wise MLP [38] blocks to embed per-location features through channel mixing. Additionally, attention [41] blocks are used for cross-location relation modeling and facilitating spatial mixing. Initially devised for natural language processing [41,11], Transformers were subsequently introduced to the domain of computer vision through works like

4 NguyenHuu B. et al.

DETR [1] and ViT [41]. Compared with CNNs, the primary distinction lies in the fact that transformers use attention as a substitute for convolution, thereby facilitating global context modeling. Nevertheless, vanilla attention, which calculates pairwise feature affinity across all spatial locations, imposes a significant computational burden and leads to substantial memory footprints, particularly when dealing with high-resolution inputs. Thus, a key research focus is to devise more efficient attention mechanisms, crucial for mitigating computational demands, especially with high-resolution inputs.

2.2 Attention mechanisms

Numerous studies have aimed to alleviate the computational and memory complexities associated with vanilla attention. Approaches include sparse connection patterns [6], low-rank approximations [42], and recurrent operations [10]. In the context of Vision Transformers, sparse attention has gained popularity, particularly following the remarkable success of the Swin Transformer [29]. Within the Swin Transformer framework, attention is constrained to non-overlapping local windows, and an innovative shift window operation is introduced. This operation facilitates communication between adjacent windows, contributing to its unique approach to handling attention mechanisms. To attain larger or approximate global receptive fields without exceeding computational constraints, recent studies have incorporated diverse manually designed sparse patterns. These include the integration of dilated windows [45,40,24] and cross-shaped windows [14]. Moreover, certain studies endeavor to make sparse patterns adaptable to data, as demonstrated by works like DAT [47], TCFormer [53], and DPT [5]. Despite their efforts to decrease the number of key-value tokens using diverse merging or selection strategies, it is crucial to recognize that these tokens lack semantic specificity. Instead, we reinforce query-aware key-value token selection.

Our work is motivated by an observation: semantically attentive regions for important queries can exhibit significant differences, as illustrated by visualizations from pre-trained models like ViT [41] and DETR [1]. In achieving query-adaptive sparsity through a coarse-to-fine approach, we propose an attention-in-attention architecture which utilizes the deformable attention [47] with bi-level routing attention [56]. Diverging from deformable attention [47] and bi-level routing attention [56], our deformable bi-level routing attention aims to reinforce the most semantic and flexible key-value pairs. In contrast, bi-level routing attention only focuses on locating a few highly relevant key-value pairs, while deformable attention prioritizes identifying a few of the most flexible key-value pairs.

3 Our Approach: DeBiFormer

3.1 Preliminaries

Initially, we revisit the attention mechanism used in recent Vision Transformers. Taking a flattened feature map $x \in \mathbb{R}^{N \times C}$ as the input, a multi-head self-

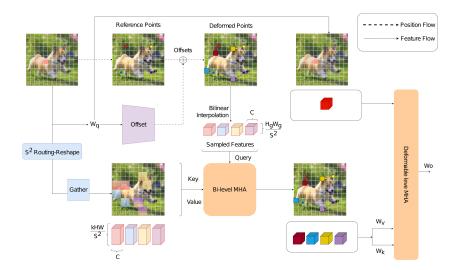


Fig. 2. Detailed architecture of Deformable Bi-level Routing Attention. In the top-left part, the set of reference points is uniformly distributed across the feature map. Offsets for these points are learned from queries through the offset network. Then, in the top-middle part, deformed features are projected from sampled features based on the locations of deformed points. In the bottom-left-middle part, we attend to projected deformed features by utilizing gathered key-value pairs in top-k-related windows.

attention (MHSA) block with M heads is formulated as

$$q = xW_q, k = xW_k, v = xW_v, \tag{1}$$

$$z^{(m)} = \sigma(q^{(m)k^{(m)\top}/\sqrt{d}})v^{(m)}, m = 1, ..., M,$$
(2)

$$z = Concat(z^{(1)}, ..., z^{(M)})W_o$$
(3)

where $\sigma(\cdot)$ denotes the softmax function, and d=C/M is the dimension of each head. z(m) denotes the embedding output from the m-th attention head, and $q(m), k(m), v(m) \in \mathbb{R}^{N \times d}$ denote the query, key, and value embeddings, respectively. $W_q, W_k, W_v, W_o \in \mathbb{R}^{C \times C}$ are the projection matrices. With normalization layers and identity shortcuts, the l-th Transformer block, for which LN means layer normalization, is formulated as

$$z'_{l} = MHSA(LN(z_{l-1})) + z_{l-1}, (4)$$

$$z_l = MLP(LN(z_l')) + z_l' \tag{5}$$

3.2 Deformable bi-level routing attention (DBRA)

The architecture of the proposed Deformable Bi-level Routing Attention (DBRA) is illustrated in Figure 2. We first employ a deformable attention module, which

includes an offset network that generates offsets for reference points based on the query features, creating deformable points. However, these points tend to cluster in important regions, leading to an over-concentration in certain areas.

To address this, we introduce deformable points-aware region partitioning, ensuring that each deformable point interacts with only a small subset of key-value pairs. Yet, solely relying on region partitioning can result in an imbalance between important and less important regions. To tackle this, the DBRA module is designed to distribute attention more effectively. In DBRA, each deformable point acts as an agent query, computing attention with semantic region key-value pairs. This approach ensures that only a few deformable points are assigned to each important region, allowing attention to be spread across all critical areas of the image rather than clustered in one spot.

By employing the DBRA module, attention is reduced in less important regions and increased in more important ones, ensuring a balanced distribution of attention throughout the image.

Deformable attention module and input projection. As illustrated in Fig.2, given the input feature map $x \in \mathbb{R}^{H \times W \times C}$, a uniform grid of points $p \in \mathbb{R}^{H_G \times W_G \times 2}$ is generated by downsampling the input feature map by factor r, $H_G = H/r$, $W_G = W/r$ as a reference. To obtain the offset for each reference point, the features are linearly projected to generate query tokens $q = xW_q$, which are then input into the $\theta_{offset}(\cdot)$ subnetwork to produce the offsets $\Delta p = \theta_{offset}(q)$. Subsequently, the features are sampled at the locations of deformed points as keys and values and further processed by projection matrices:

$$q = xW_q, \Delta p = \theta_{offset}(q), \bar{x} = \varphi(x; p + \Delta p),$$
 (6)

where \bar{x} represent the deformed key \bar{k} and value \bar{v} embeddings, respectively. Specifically, we set the sampling function $\varphi(\cdot;\cdot)$ to a bilinear interpolation to make it differentiable:

$$\varphi(z;(p_x, p_y)) = \sum_{r_x, r_y} g(p_x, r_x) g(p_y, r_y) z[r_y, r_x, :],$$
(7)

where the function g(a,b) = max(0,1-|a-b|) and (r_x,r_y) represent indices for all locations on $z \in \mathbb{R}^{H \times W \times C}$. In a similar setup as deformable attention, where g is nonzero on the four integral points closest to (px,py), Equation 7 is simplified to a weighted average across these four locations.

Region partition and region-to-region routing. With the deformable attention feature map input $\bar{x} \in \mathbb{R}^{H_G \times W_G \times C}$ and feature map $x \in \mathbb{R}^{H \times W \times C}$, the process begins by dividing it into regions of size $S \times S$ non-overlapped regions such that each region contains $\frac{H_G W_G}{S^2}$ feature vectors with reshaped \bar{x} as $\bar{x}^r \in \mathbb{R}^{S^2 \times \frac{H_G W_G}{S^2}} \times C$ and x as $x^r \in \mathbb{R}^{S^2 \times \frac{H_W}{S^2}} \times C$. Then, we derive the query, key, and value with linear projections:

$$\hat{q} = \bar{x}^r W_q, \hat{k} = x^r W_k, \hat{v} = x^r W_v. \tag{8}$$

Next, we use the region-to-region method, as introduced in BiFormer [56], which is applied to establish the attending relationship by constructing a directed

graph. To initiate the process, region queries and keys $\hat{q}^r, \hat{k}^r \in S^{S^2 \times C}$ are derived through the application of per-region averaging. Subsequently, the adjacency matrix $A^r \in S^2 \times S^2$ for the region-to-region affinity graph is derived through Q^r and K^r^{\top} matrix multiplication:

$$A^r = \hat{q}^r (\hat{k}^r)^\top, \tag{9}$$

where adjacency matrix A^r quantifies the semantic relationship between two regions. The crucial step in this method involves pruning the affinity graph by retaining only the topk connections for each region with a routing index matrix $I^r \in \mathbb{N}^{S^2 \times k}$ through the use of the topk operator:

$$I^r = topk(A^r). (10)$$

Bi-level token to deformable-level token attention. Utilizing the region routing matrix I^r , we can then apply token attention. For each deformable query token within region i, its attention spans all key-value pairs located in the topk routed regions, that is, those indexed by $I_{i,1}^r, I_{i,2}^r, ..., I_{i,k}^r$. Hence, we continue the process of gathering the key and value:

$$\hat{k}^g = gather(\hat{k}, I^r), \hat{v}^g = gather(\hat{v}, I^r), \tag{11}$$

where $\hat{k}^g, \hat{v}^g \in \mathbb{R}^{S^2 \times \frac{kHW}{S^2} \times C}$ are the gathered key and value. Then, we apply attention on \hat{k}^g, \hat{v}^g as:

$$\hat{O} = \hat{x} + W_{o'}(Attention(\hat{q}, \hat{k}^g, \hat{v}^g) + LCE(\hat{v})), \tag{12}$$

$$O = MLP(LN(\hat{O})) + \hat{O}, \tag{13}$$

where $W_{o'}$ is a projection weight for output features, and $LCE(\cdot)$ uses a kernel size 5 depth-wise convolution.

Deformable-level token to token attention. Following that, the deformable features that are semantically attended to via [56] are reshaped O as $O^r \in \mathbb{R}^{H_G \times W_G \times C}$ and parameterized at the locations of keys and values:

$$k = O^r W_k, v = O^r W_v, \tag{14}$$

k and v represent the embeddings of semantically deformed keys and values, respectively. Using existing approaches, we perform self-attention on q, k, v, and relative position offsets R. The output of attention is formulated as follows:

$$z^{m} = W_{\bar{o}}(\sigma(q^{m}k^{(m)\top}/\sqrt{d} + \phi(\hat{B};R))v^{m}). \tag{15}$$

Here, $\phi(\hat{B}; R) \in \mathbb{R}^{HW \times H_GW_G}$ corresponds to the position embedding, following the approach of previous work [29]. Then z^m is projected though W_o to get the final output z as Equation 3.

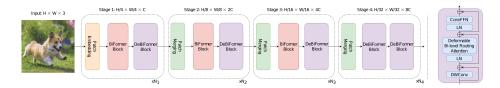


Fig. 3. Overall model architecture of our DeBiFormer. Left: Network architecture of DeBiFormer. N_1 to N_4 represent numbers of stacked successive local and Deformable Bi-level Routing Attention blocks. Please consult Table 1 for specific configurations. Right: Details on DeBiFormer Block.

	Variant Arch	itectures of DeBiFo	rmer
	DeBi-T	DeBi-S	DeBi-B
Store 1	$N_1=2, C=64$	$N_1=4, C=64$	$N_1=4, C=96$
Stage 1 56×56	$r=8, M=2, D_r=3$	$r=8, M=2, D_r=3$	$r=8, M=3, D_r=3$
	G=1, K=9, B_r =3	G=1, K=9, B_r =3	G=1, K=9, B_r =3
Store 2	$N_2=2, C=128$	$N_2=4, C=128$	$N_2=4, C=192$
Stage 2 28×28	$r=4, M=4, D_r=3$	$r=4, M=4, D_r=3$	$r=4, M=6, D_r=3$
	G=2, K=7, B_r =3	G=2, K=7, B_r =3	G=2, K=7, B_r =3
Store 2	$N_3=8, C=256$	$N_3=18, C=256$	$N_3=18, C=384$
Stage 3 14×14	$r=2, M=8, D_r=3$	$r=2, M=8, D_r=3$	$r=2, M=12, D_r=3$
14 / 14	$G=4, K=5, B_r=3$	$G=4, K=5, B_r=3$	G=4, K=5, B_r =3
Stage 4 7×7	$N_4=2, C=512$	$N_4=6, C=512$	$N_4=4, C=768$
	$r=1, M=16, D_r=3$	$r=1, M=16, D_r=1$	$r=1, M=24, D_r=3$
	G=8, K=3, B_r =3	G=8, K=3, B_r =2	G=8, K=3, B_r =3

Table 1. DeBiFormer model architecture specifications. N_i : Number of blocks at stage i. C: Base channels in each block. r: Downsample ratio of deformed points. M: Number of attention heads in DBRMHA. G: Number of offset groups in DBRMHA. D_r : Deformable level MLP expansion ratio. B_r : Bi level MLP expansion ratio. K: Kernel size of offset module.

3.3 Model architectures

Leveraging DBRA as a fundamental building block, we introduce a novel vision transformer called DeBiFormer. As depicted in Figure3, we adhere to the recent state-of-the-art Vision Transformers [14,29,56,47], using a four-stage pyramid structure. In stage i, we utilize an overlapped patch embedding in the first stage and a patch merging module [26,34] in the second to fourth stages. This is done to decrease the input spatial resolution while increasing the number of channels. Subsequently, N_i consecutive DeBiFormer blocks are used to transform the features. Within each DeBiFormer block, we adhere to recent methodologies [26,40,56] by using a 3×3 depthwise convolution at the outset. This is done to implicitly encode relative position information. Following that, we sequentially use a DBRA module with a 2-ConvFFN module with an expansion ratio e for cross-location relation modeling and per-location embedding, respectively. DeBiFormer is instantiated in three distinct model sizes, achieved by scaling the

M - 1-1	FLOPs	Params	Top-1	M - 1-1	FLOPs Params Top-1		
Model	(G)	(M)	(%)	Model	(G)	(M)	(%)
ResNet-18 [20]	1.8	11.7	69.8	ConvNeXt-B [30]	15.4	89	83.8
PVTv2-b1 [42]	2.1	13.1	78.7	SLaK-S [27]	9.8	55	83.8
Shunted-T [34]	2.1	11.5	79.8	Twins-SVT-L [7]	14.8	99	83.7
QuadTree-B-b1 [37]	2.3	13.6	80.0	PVTv2-B5 [44]	11.8	82	83.8
BiFormer-T [56]	2.2	13.1	81.4	Swin-B [29]	15.4	88	83.5
Conv2Former-N [22]	2.2	15.0	81.5	Focal-B [49]	16.4	90	84.0
DeBiFormer-T	2.6	21.4	81.9	CSWin-B [14]	15.0	78	84.2
PVTv2-B3 [44]	6.9	45	81.2	Shunted-B [34]	8.1	39	84.0
Swin-T [29]	4.5	29	81.3	UniFormer-B [26]	8.3	50.0	83.9
CSWin-T [14]	4.5	23	82.7	ScalableViT-B [51]	8.6	81	84.1
DAT-T [47]	4.6	29	82.0	Slide-Swin-B [33]	15.5	89.0	84.2
CrossFormer-S [45]	5.3	31	82.5	DAT-S [47]	9.0	50	83.7
RegionViT-S+ $[2]$	5.7	31	83.3	QuadTree-B-b4[37]	11.5	64	84.1
QuadTree-B-b3 [37]	7.8	46	83.7	CrossFormer-L [45]	16.1	92	84.0
MaxViT-T [40]	5.6	31	83.6	RegionViT-B+ [2]	13.6	74	83.8
InternImage-T [43]	5.0	30	83.5	InternImage-S [43]	8.0	50	84.2
MixFormer-B4 [4]	3.6	35	83.0	MixFormer-B6 [4]	12.7	119	83.8
BiFormer-S [56]	4.5	26	83.8	BiFormer-B [56]	9.8	57	84.3
UniRepLKNet-T [12]	4.9	31	83.2	UniRepLKNet-S [12]	9.1	56	83.9
DeBiFormer-S	5.4	44	83.9	DeBiFormer-B	11.8	77	84.4
TD 11 0 TO 1 11			1:00	. 1 11 T	37 . 477		

Table 2. Evaluating and comparing different backbones on ImageNet-1K on images with resolution of 224×224 .

network width and depth as outlined in Table 1. Each attention head comprises 32 channels, and we use a bi-level ConvFFN and deformable-level ConvFFN with an MLP expansion ratio of e=3. For the BRA, we use $topk=1,4,16,S^2$, and for the DBRA, we use $topk=4,8,16,S^2$ for the four stages. Moreover, we set the region partition factor S to specific values: S=7 for classification, S=8 for semantic segmentation, and S=20 for object detection tasks.

4 Experiments

We experimentally evaluated the effectiveness of our proposed DeBiFormer on various mainstream computer vision tasks, including image classification (Section 4.1), semantic segmentation (Section 4.2) and object detection and instance segmentation (Section 4.3). In our approach, we commence training from scratch on ImageNet-1K [35] for image classification. Subsequently, we fine-tune the pretrained backbones on ADE20K [55] for semantic segmentation and on COCO [17] for object detection and instance segmentation. Furthermore, we perform an ablation study to confirm the efficacy of the proposed Deformable Bi-level Routing Attention and top-k choices of DeBiFormer (Section 4.4). Finally, in order to validate that the recognition ability and interpretability of our DeBiFormer, we visualize the attention map (Section 5).

Backbone	Semantic-FPN	UperNet
Dackbone	mIoU(%)	mIoU(%)
Swin-T [29]	41.5	44.5
DAT-T [47]	42.6	45.5
CSWin-T [14]	48.2	49.3
RegionViT-S+ $[2]$	45.3	-
InternImage-T [43]	-	47.9
CrossFormer-S [45]	46.0	47.6
Uniformer-S [26]	46.6	47.6
Shunted-S [34]	48.2	48.9
BiFormer-S [56]	48.9	49.8
DeBiFormer-S	49.2	50.0
Swin-S [29]	-	47.6
DAT-S [47]	46.1	48.3
CSWin-S [14]	49.2	50.4
RegionViT-B+ [2]	47.5	-
InternImage-S [43]	-	50.1
CrossFormer-B [45]	47.7	49.7
Uniformer-B [26]	48.0	50.0
BiFormer-B [56]	49.9	51.0
DeBiFormer-B	50.6	51.4

Table 3. Evaluating DeBiFormer on semantic segmentation with two segmentation heads (Semantic FPN and UpperNet) on ADE20K dataset.

4.1 Image classification on ImageNet-1K

Settings. We conducted image classification experiments on the ImageNet-1K [35] dataset, following the experimental settings of DeiT [39] for a fair comparison. Specifically, each model was trained for 300 epochs on 8 V100 GPUs with an input size of 224×224 . We used AdamW as the optimizer with a weight decay of 0.05 and used a cosine decay learning rate schedule with an initial learning rate of 0.001, while the first five epochs were used for linear warm-up. The batch size was set to 1024. To avoid overfitting, we used regularization techniques including RandAugment [9] (rand-m9-mstd0.5-inc1), MixUp [54] (prob = (0.8), CutMix [52] (prob = 1.0), Random Erasing (prob = 0.25), and increasing stochastic depth [23] (prob = 0.1/0.2/0.4 for DeBiFormer-T/S/B, respectively). Results. We report our results in Table 2 showing the top-1 accuracy with similar computational complexities. Our DeBiFormer outperformed the Swin Transformer [29], PVT [44], DeiT [39], DAT[47], and Biformer [56] in all three scales. Without inserting convolutions in Transformer blocks or using overlapped convolutions in patch embeddings, DeBiFormer achieved gains of 0.5pt, 0.1pt and 0.1pt over BiFormer [56] counterparts.

4.2 Semantic segmentation on ADE20K

Settings. The same as existing works, we used our DeBiFormer on SemanticFPN [46] and UperNet [48]. In both cases, the backbone was initialized with

								DC	DII OI	inci	11	
Backbone]	Retina	Net 1	\times sch	edule		1	Mask I	R-CNI	$1 \times sc$	hedule	
Dackbone	mAP	AP_{50}	AP_{75}	AP_S	AP_{M}	AP_L	mAP^b	AP_{50}^b	AP_{75}^b	mAP^m	AP_{50}^m	AP_{75}^{m}
Swin-T [29]	41.5	62.1	44.2	25.1	44.9	55.5	42.2	64.6	46.2	39.1	61.6	42.0
DAT-T [47]	42.8	64.4	45.2	28.0	45.8	57.8	44.4	67.6	48.5	40.4	64.2	43.1
CSWin-T [14]	-	-		-	-	-	46.7	68.6	51.3	42.2	65.6	45.4
RegionViT-S+ [2]	43.9	65.5	47.3	28.5	47.3	58.0	44.2	67.3	48.2	40.9	64.1	44.0
MixFormer-B4 [4]	-	-		-	-	-	45.1	67.1	49.2	41.2	64.3	44.4
CrossFormer-S [45]	44.4	55.3	38.6	19.3	40.0	48.8	45.4	68.0	49.7	41.4	64.8	44.6
QuadTree-B2 [37]	46.2	67.2	49.5	29.0	50.1	61.8	-	-	-	-	-	-
InternImage-T [43]	-	-		-	-	-	47.2	69.0	52.1	42.5	66.1	45.8
Agent-Swin-T [18]	-	-		-	-	-	44.6	67.5	48.7	40.7	64.4	43.4
BiFormer-S [56]	45.9	66.9	49.4	30.2	49.6	61.7	47.8	69.8	52.3	43.2	66.8	46.5
DeBiFormer-S	45.6	66.6	48.9	28.7	49.3	61.6	47.5	69.7	52.1	42.5	66.2	45.7
Swin-S [29]	44.5	65.7	47.5	27.4	48.0	59.9	44.8	66.6	48.9	40.9	63.4	44.2
DAT-S [47]	45.7	67.7	48.5	30.5	49.3	61.3	47.1	69.9	51.5	42.5	66.7	45.4
CSWin-S [14]	-	-	-	-	-	-	47.9	70.1	52.6	43.2	67.1	46.2
RegionViT-B+ [2]	44.6	66.4	47.6	29.6	47.6	59.0	45.4	68.4	49.6	41.6	65.2	44.8
CrossFormer-B [45]	46.2	67.8	49.5	30.1	49.9	61.8	47.2	69.9	51.8	42.7	66.6	46.2
QuadTree-B3 [37]	47.3	68.2	50.6	30.4	51.3	62.9	-	-	-	-	-	-
InternImage-S [43]	-	-		-	-	-	47.8	69.8	52.8	43.3	67.1	46.7
Agent-Swin-S [18]	-	-		-	-	-	47.2	69.6	52.3	42.7	66.6	45.8
BiFormer-B [56]	47.1	68.5	50.4	31.3	50.8	62.6	48.6	70.5	53.8	43.7	67.6	47.1
DeBiFormer-B	47.1	68.2	50.2	30.3	51.1	63.0	48.5	70.2	53.3	43.2	67.2	46.4

DeBiFormer

11

Table 4. Results on object detection (left group) and instance segmentation (right group) tasks, performed on COCO 2017 dataset.

ImageNet-1K pretrained weights. The optimizer was AdamW [31], and the batch size was 32. For a fair comparison, we followed the same setting as PVT [44] to train the model with 80k steps and Swin Transformer [29] to train the model with 160k steps.

Results. Table 8 shows the results of the two different frameworks. It shows that with the Semantic FPN framework, our DeBiFormer-S/B achieved 49.2/50.6 mIoU, respectively, improving BiFormer by 0.3pt./0.7pt. A similar performance gain for the UperNet framework was also observed. By utilizing the DBRA module, our DeBiFormer could caputure the most semantic key-value pairs, which makes the attention selection more reasonable and achieve higher performance on downstream semantic tasks.

4.3 Object detection and instance segmentation

Settings. We used our DeBiFormer as the backbone in the Mask RCNN [19] and RetinaNet [16] frameworks to evaluate the effectiveness of models for object detection and instance segmentation on COCO 2017 [17]. The experiments were conducted with the MMDetection [3] toolbox. Before training on COCO, we initialized the backbone with weights pre-trained on ImageNet-1K and followed the same training strategies as BiFormer [56] to compare our methods fairly. Note that due to device limitations, we set mini batch size as 4 for these experiments,

Sparse Attention	IN1K	ADE20K
Sparse Attention	Top1(%)	mIoU(%)
Shifted window [29]	81.3	41.5
Spatially sep [7]	81.5	42.9
Sequential axial [21]	81.5	39.8
Criss-cross [45]	81.7	43.0
Cross-shaped window [14]	82.2	43.4
Deformable [47]	82.0	42.6
Block-grid [40]	81.8	42.8
Bi-level routing [56]	82.7	44.8
Deformable bi-level routing	82.9	48.0

Table 5. Ablation study on different attention mechanisms. All models follow the architecture design of the Swin-T model.

while in BiFormer this value is 16. For details on the specific settings of the experiment, please refer to the supplementary paper.

Results. We list the results in Table 4.2. For object detection with RetinaNet, we report the mean average precision (mAP) and the average precision (AP) at different IoU thresholds (50%, 75%) for three object sizes (i.e., small, medium, and large (S/M/L)). From the results, we can see that while the overall performance of DeBiFormer was only comparable to some of the most competitive existing methods, the performance on large objects (AP_L) outperformed these methods although we use a limited resources. This may be because the DBRA allocates deformable points more reasonably. These points are not to focus only on small things, but to focus on important things in the image. Therefore the attention is not limited to a small area, which improves the detection accuracy of large objects. For instance segmentation with Mask R-CNN, we report the bounding box and mask the average precision $(AP_b \text{ and } AP_m)$ at different IoU thresholds (50%, 75%). Note that our DeBiFormer still achieved great performance under the device limitation of mini batch size. We believe that we could achieve better results if the mini batch size could be the same to other methods since it has been proved on semantic segmentation tasks.

4.4 Ablation study

Effectiveness of DBRA. We compared DBRA with several existing sparse attention mechanisms. Following CSWIN [14], we aligned macro architecture designs with Swin-T [29] for fair comparison. Specifically, we used 2, 2, 6, 2 blocks for the four stages and non-overlapped patch embedding, and we set the initial patch embedding dimension to C=96 and MLP expansion ratio to e=4. The results are reported in Table 5. Our Deformable Bi-level Routing Attention had significantly better performance than the existing sparse attention mechanisms, in terms of both image classification and semantic segmentation.

Partition factor S. Similar to BiFormer, we opted to use S as a divisor of the training size to prevent padding. We used an image classification with a resolu-

#partition factor	k	Bi-level Routing tokens to attend	CoreML Latency(ms)	Top-1 Acc (%)
7	1,4,16,49	64, 64, 64, 49	291	81.81
7	2,8,32,49	128, 128, 128, 49	459	81.74
7	4,8,16,32	256, 128, 64, 32	276	81.47
7	8,16,32,49	512, 256, 128, 49	498	81.60
7	4,8,16,49	256, 128, 64, 49	382	81.90

Table 6. Ablation study on top-k selection. The four numbers in the k column represent the top-k values of the four stages and the same to the tokens to attend column. The CoreML Latency is conducted by MacBook Pro M1 using CPU and Nerual Engine.

tion of $224 = 7 \times 32$, and we set S = 7 to ensure that the size of the feature maps was divided at each stage. This choice aligns with the strategy used for the Swin Transformer [29], where a window size of 7 was used.

Top-k Choices. We systematically adjusted k to ensure a reasonable number of tokens attended to deformable queries as the region size diminished in later stages. Exploring various combinations of k is a viable option. In Table 9, we present ablation results on IN-1K, following DeBiFormer-STL ("STL" denotes Swin-T Layout). A crucial observation from these experiments is that augmenting the number of tokens paid attention to the deformable queries had a detrimental effect on accuracy and latency, and increasing the number of tokens paid attention in stages 1 and 2 had an effect on accuracy.

Deformable Bi-level Routing Multi-Head Attention (DBRMHA) at different stages. To evaluate the impact of design choices, we systematically replaced bi-level routing attention blocks with DBRMHA blocks across different stages, as shown in Table7. Initially, all stages used bi-level routing attention, similar to BiFormer-T [56], achieving 81.3% accuracy in image classification. Replacing just one block in the 4th stage with DBRMHA immediately boosted accuracy by +0.21. Replacing all blocks in the 4th stage added another +0.05. Further DBRMHA replacements in the 3rd stage continued to improve performance across tasks. While gains tapered off with earlier stage replacements, we settled on a final version—DeBiFormer—where all stages use Deformable Bi-level Routing Attention for simplicity.

5 Grad-CAM Visualization

To further illustrate the ability of the proposed DeBiFormer to recognize the attention in important regions, we used Grad-CAM [36] to visualize the areas of greatest concern of BiFormer-Base and DeBiFormer-Base. As shown in Figure 4, by utilizing DBRA module, our DeBiFormer-Base model performed better in target objects locating in which more regions have been focused on. In addition, our model scales down the attention in the unnecessary regions and pays more attention to the necessary regions. Depending on the attention of more necessary regions, our DeBiFormer model focused on semantic areas more continuously and

14 NguyenHuu B. et al.

Stage Config	Stage Configurations			IN-1K
1 2 3	4	FLOPs	#Param	Acc.
BB BB BB	BB	2.2	13.1	81.37
BB BB BB	\mathbf{B}^*	2.4	15.8	81.58
BB BB BB	**	2.6	18.5	81.63
BB BB B*	**	2.7	21.2	81.84
BB $\mathbf{B^*}$ $\mathbf{B^*}$	**	2.7	21.4	81.87
B* B* B*	**	2.6	21.4	81.90

Table 7. Configuration **B***: this stage is constructed with successive bi-level routing attention and DBRMHA blocks, while BB and ****** denotes that stage uses the same blocks as bi-level routing attention and DBRMHA. respectively.

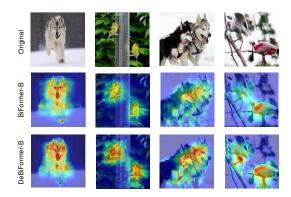


Fig. 4. Grad-CAM Visualization of BiFormer-Base and DeBiFormer-Base. These images are sampled from the validation set of ImageNet-1K.

completely, suggesting the stronger recognition ability of our model. Such ability yields better classification, and semantic segmentation performance compared with BiFormer-Base.

6 Conclusion

The paper introduces the Deformable Bi-level Routing Attention Transformer, a novel hierarchical vision transformer designed for both image classification and dense prediction tasks. Through Deformable Bi-level Routing Attention, our model optimizes query-key-value interactions while adaptively selecting semantically relevant regions. This leads to more efficient and meaningful attention. Extensive experiments show our model's effectiveness compared to strong baselines. We hope this work offers insights into designing flexible and semantically aware attention mechanisms.

7 Supplementary Material

7.1 Offset groups

Like [47], to foster diversity among deformed points, we adhere to a similar paradigm as in MHSA, where the channels are split into multiple heads to compute a variety of attentions. Hence, we partition the channels into G groups to generate diverse offsets. The offset generation network shares weights for features from different groups.

7.2 Deformable relative position bias

Certainly, incorporating position information into attention mechanisms has proven beneficial for model performance. Various approaches, such as APE [15], RPE [29], CPE [8], LogCPB [28], and others have demonstrated improved results. The relative position embedding (RPE) introduced in the Swin Transformer specifically encodes the relative positions between every pair of query and key, thereby enhancing vanilla attention with spatial inductive bias [29]. The explicit modeling of relative locations is especially well-suited for attention heads at the deformable level. In this context, deformed keys can assume arbitrary continuous locations, as opposed to being confined to fixed discrete grids.

In accordance with [47], the relative coordinate displacements are restricted to the range of [-H, +H] and [-W, +W] along the spatial dimension with a relative position bias (RPB), denoted as \hat{B} with dimensions $(2H-1)\times(2W-1)$.

Then, the relative locations within the range of [-1, +1] are sampled using bilinear interpolation $\varphi(\hat{B}; R)$ with a parameterized bias. This is done by considering continuous relative displacements, ensuring coverage of all conceivable offset values.

7.3 Computational complexity

Deformable bi-level routing attention (DBRA) incurs a comparable computation cost to its counterpart in the Swin Transformer. The computation of DBRA consists of two parts: token-to-token attention and offset & sampling. Therefore, the computation of this part is:

$$FLOPs_{def} = FLOPs_{attn} + FLOPs_{offset\&sampling}$$

= $2HWN_sC + 2HWC^2 + 2N_sC^2 + (k^2 + 6)N_sC$ (16)

where $N_s = HW/r^2$ is the number of sampled points, and C is the token embedding dimension. The computation of bi-level routing multi-head attention consists of three parts: linear projection, region-to-region routing, and token-to-token attention. Hence, the computation of this part is:

$$FLOPs_{bi} = FLOPs_{proj} + FLOPs_{routing} + FLOPs_{attn}$$

$$= 2HWC^{2} + 2N_{s}C^{2} + 2(S^{2})^{2}C + 2HWk\frac{N_{s}}{S^{2}}C$$

$$= 2HWC^{2} + 2N_{s}C^{2} + C\{2S^{4} + 2HWk\frac{N_{s}}{S^{2}}\}$$

$$= 2HWC^{2} + 2N_{s}C^{2} + C\{2S^{4} + \frac{kHWN_{s}}{S^{2}} + \frac{kHWN_{s}}{S^{2}}\}$$

$$\geq 2HWC^{2} + 2N_{s}C^{2} + 3C\{2S^{4} \cdot \frac{kHWN_{s}}{S^{2}} \cdot \frac{kHWN_{s}}{S^{2}}\}^{\frac{1}{3}}$$

$$= 2HWC^{2} + 2N_{s}C^{2} + 3Ck^{\frac{2}{3}}\{2HWN_{s}\}^{\frac{2}{3}}$$

$$= 2HWC^{2} + 2N_{s}C^{2} + 3Ck^{\frac{2}{3}}\{2HWN_{s}\}^{\frac{2}{3}}$$

$$(17)$$

where k is the number of regions to attend to, and S is the region partition factor. Last, the total computation of DBRA consists of two parts: deformable level multi-head attention and bi-level routing multi-head attention. The total amount of computations is therefore:

$$FLOPs = FLOPs_{bi} + FLOPs_{def}$$

$$= 2HWN_sC + 2HWC^2 + 2N_sC^2 + (k^2 + 6)N_sC$$

$$+2HWC^2 + 2N_sC^2 + 3Ck^{\frac{2}{3}}\{2HWN_s\}^{\frac{2}{3}}$$

$$= 2HWN_sC + 4HWC^2 + 4N_sC^2$$

$$+(k^2 + 6)N_sC + 3Ck^{\frac{2}{3}}\{2HWN_s\}^{\frac{2}{3}}$$
(18)

In other words, DBRA achieves $O((HWN_s)^{\frac{2}{3}})$ complexity. For example, the 3^{rd} stage of a hierarchical model with 224×224 input for image classification usually has computation sizes of $H = W = 14 \, S^2 = 49 \, N_s = 1 \, C = 384$ and thus computational complexity for multi-head self-attention. Additionally, the complexity could be further reduced by enlarging the downsampling factor r and scale with region partition factor S, making it friendly to tasks with much higher resolution inputs such as object detection and instance segmentation.

7.4 Token to Attend

In Table 9, we present the token to attend to the query and the token to attend to the deformable point. Compared with other methods, DeBiFormer has the fewest tokens to attend for each query but has a high performance in Imagenet1K, ADE20K(S-FPN head) and COCO(Retina head).

7.5 More visualization results

Effective Receptive Field Analysis To assess the effective receptive fields (ERFs) [32] of the central pixel with an input size of 224×224 across various

Backbone	mIoU	MS mIoU
Swin-S [29]	47.6	49.5
DAT-S [47]	48.3	49.8
CSWin-S [14]	50.4	51.5
InternImage-S [43]	50.1	50.9
CrossFormer-B [45]	49.7	50.6
Uniformer-B [26]	50.0	50.8
BiFormer-B [56]	51.0	51.7
DeBiFormer-B	51.4	52.0

Table 8. Evaluating DeBiFormer on semantic segmentation with two segmentation heads (UpperNet With MS IoU Scale) on ADE20K dataset.

// D = -1-1	Token to attend	IN1K	ADE20K	COCO
#Backbone	for each query	Top1(%)	mIoU(%)	mAP(%)
QuadTree-B3 [37]	2048, 512, 128, 32	83.8	50.0	47.3
BiFormer-B [56]	64, 64, 64, 49	84.3	49.9	47.1
Debiformer-B	1, 1, 1, 1	84.4	50.6	47.1

Table 9. Comparing different backbones Token to attend for each query in on ImageNet-1K, ADE20K, COCO.

models, we present a comparative analysis in Figure 5. To demonstrate the powerful representation capacity of our DeBiFormer, we also compare the ERF of several SOTA methods with similar computational costs. As shown in Fig. 5, our DeBiFormer has the largest and most consistent ERF among these methods while maintaining strong local sensitivity, which is challenging to achieve.

Grad-CAM Analysis To further show how DBRA works, we demonstrate more visualization results in Figure 6. Thanks to the flexible key-value pairs selection, in most cases, our DeBiFormer focuses on important objects at an earlier stage. Meanwhile, due to the reasonable allocation of deformable points, it also focuses on different important areas earlier in multi-object scenarios. With the powerful DBRA module, our DeBiFormer has a larger heat map area in the last two stages, which represents a stronger ability of recognition.

7.6 Detailed Experimental Settings

Image classification on ImageNet-1K As introduced in the main paper, each model was trained for 300 epochs on 8 V100 GPU with an input size of 224×224 . The experimental settings are strictly follow the DeiT[39] for a fair comparison. For more details, please refer to the Table 10 provided.

Object Detection and Instance Segmentation When fine-tuning our DeBiFormer to object detection and instance segmentation on COCO [17], we have considered two common frameworks: Mask R-CNN [19], RetinaNet [16]. For

optimization, we adopt the AdamW optimizer with an initial learning rate of 0.0002 and a **mini batch size of 4** due to the limitation of devices. When training models of different sizes, we adjust the training settings according to the settings used in image classification. The detailed hyper-parameters used in training models are presented in Table 11.

Semantic Segmentation For ADE20K, we utilized the AdamW optimizer with an initial learning rate of 0.00006, a weight decay of 0.01, and a mini batch size of 16 for all models trained for 160K iterations. In terms of testing, we reported the results using both single-scale (SS) and multi-scale (MS) testing in the main comparisons. For multi-scale testing, we experimented with resolutions ranging from 0.5 to 1.75 times that of the training resolution. To set the path drop rates in different models, we used the same hyper-parameters as those used for object detection and instance segmentation. Table 8 shows the results of the Upernet frameworks with the single and Multi-Scale IoU.

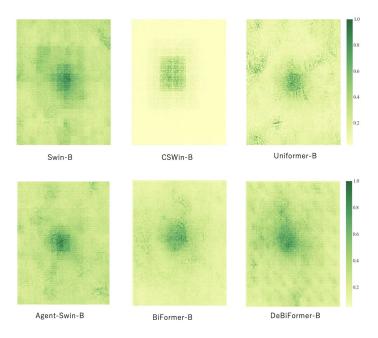


Fig. 5. ERF visualization of models incorporating various local operators and SOTA methods. The results are obtained by averaging over 100 images (resized to 224×224) from ImageNet.

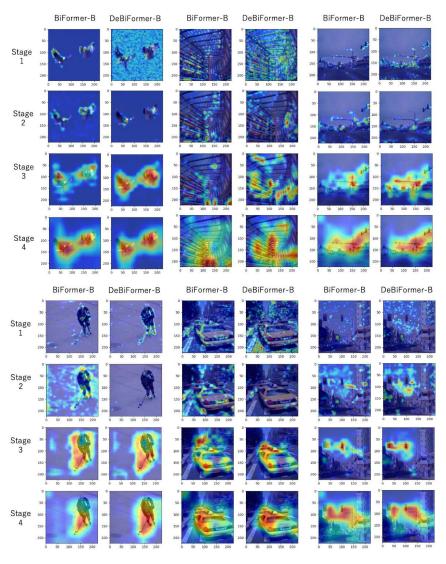


Fig. 6. More Visualization of attention maps for different stages of DeBiFormer-Base. The attention maps show the focused regions has extracted large to small in each stage. Meanwhile the model fitted out the unessential regions to extract and enhance the important regions features in each stage.

7.7 Limitation and Future Work

In contrast to sparse attention with simple static patterns, we propose a novel attention method that consists of two components. First, we prune a region-level graph and gather key-value pairs for important regions, which are focused by highly flexible key-value pairs. Then, we apply token-to-token attention. While

Settings	DeBi-T	DeBi-S	DeBi-B
Input resolution	224	224	224
Batch size	1024	512	512
Optimizer	AdamW	AdamW	AdamW
Learning rate	1×10^{-3}	5×10^{-4}	5×10^{-4}
LR schedule	cosine	cosine	cosine
Weight decay	5×10^{-2}	5×10^{-2}	5×10^{-2}
Warmup epochs	20	20	20
Epochs	300	300	300
Horizontal flip	√	√	√
Random resize	✓	\checkmark	\checkmark
${ m AutoAugment}$	✓	\checkmark	\checkmark
Mixup alpha	0.8	0.8	0.8
Cutmix alpha	1.0	1.0	1.0
Random erasing prob.	0.25	0.25	0.25
Color jitter	0.4	0.4	0.4
Label smoothing	0.1	0.1	0.1
Droppath rate	0.	0.2	0.4
Grad. clipping	5.0	5.0	5.0

Table 10. Image Classification Training Settings

Config	Value
Optimizer	AdamW
LR	0.0002
weight decay	0.05
batch size	4
LR schedule	steps:[8, 11]
training	epochs 12
scales	(800, 1333)
drop path	0.2 (Small), 0.3 (Base)

Table 11. Object Detection and Instance Segmentation Training Settings

this method does not incur much computation as it operates at a top-k routed semantically relevant region level and deformable important regions, it inevitably involves additional parameter capacity transactions during linear projection. In our future endeavors, we plan to investigate efficient sparse attention mechanisms and enhance Vision Transformers with parameter capacity awareness.

References

- 1. Carion, N., Massa, F., Synnaeve, G., Usunier, N., Kirillov, A., Zagoruyko, S.: End-to-end object detection with transformers. In: Proceedings of the IEEE/CVF European Conference on Computer Vision (ECCV) (2020) 1, 2, 4
- 2. Chen, C.F., Panda, R., Fan, Q.: Regionvit: Regional-to-local attention for vision transformers. arXiv preprint arXiv:2106.02689 (2021) 9, 10, 11

- Chen, K., Wang, J., Pang, J., Cao, Y., Xiong, Y., Li, X., Sun, S., Feng, W., Liu, Z., Xu, J., Zhang, Z., Cheng, D., Zhu, C., Cheng, T., Zhao, Q., Li, B., Lu, X., Zhu, R., Wu, Y., Dai, J., Wang, J., Shi, J., Ouyang, W., Loy, C.C., Lin, D.: MMDetection: Open mmlab detection toolbox and benchmark. arXiv preprint arXiv:1906.07155 (2019) 11
- Chen, Q., Wu, Q., Wang, J., Hu, Q., Hu, T., Ding, E., Cheng, J., Wang, J.: Mixformer: Mixing features across windows and dimensions (2022) 9, 11
- Chen, Z., Zhu, Y., Zhao, C., Hu, G., Zeng, W., Wang, J., Tang, M.: Dpt: Deformable patch-based transformer for visual recognition. In: Proceedings of the ACM International Conference on Multimedia (2021) 2, 4
- Child, R., Gray, S., Radford, A., Sutskever, I.: Generating long sequences with sparse transformers. URL https://openai.com/blog/sparse-transformers (2019) 4
- 7. Chu, X., Tian, Z., Wang, Y., Zhang, B., Ren, H., Wei, X., Xia, H., Shen, C.: Twins: Revisiting the design of spatial attention in vision transformers. In: NeurIPS 2021 (2021), https://openreview.net/forum?id=5kTlVBkzSRx 9, 12
- Chu, X., Tian, Z., Zhang, B., Wang, X., Shen, C.: Conditional positional encodings for vision transformers. In: ICLR 2023 (2023), https://openreview.net/forum?id=3KWnuT-R1bh 15
- Cubuk, E.D., Zoph, B., Shlens, J., Le, Q.: Randaugment: Practical automated data augmentation with a reduced search space. In: Larochelle, H., Ranzato, M., Hadsell, R., Balcan, M., Lin, H. (eds.) Advances in Neural Information Processing Systems. vol. 33, pp. 18613–18624. Curran Associates, Inc. (2020) 10
- Dai, Z., Yang, Z., Yang, Y., Carbonell, J., Le, Q.V., Salakhutdinov, R.: Transformer-xl: Attentive language models beyond a fixed-length context (2019) 4
- 11. Devlin, J., Chang, M.W., Lee, K., Toutanova, K.: Bert: Pre-training of deep bidirectional transformers for language understanding. arXiv preprint arXiv:1810.04805 (2018) 1, 3
- Ding, X., Zhang, Y., Ge, Y., Zhao, S., Song, L., Yue, X., Shan, Y.: Unireplknet: A universal perception large-kernel convnet for audio, video, point cloud, time-series and image recognition. arXiv preprint arXiv:2311.15599 (2023)
- 13. Dong, K., Xue, J., Lan, X., Lu, K.: Biunet: Towards more effective unet with bi-level routing attention. In: The British Machine Vision Conference (November 2023) 3
- Dong, X., Bao, J., Chen, D., Zhang, W., Yu, N., Yuan, L., Chen, D., Guo, B.: Cswin transformer: A general vision transformer backbone with cross-shaped windows. In: Proceedings of the IEEE/CVF Computer Vision and Pattern Recognition Conference (CVPR). pp. 12124–12134 (2022) 1, 4, 8, 9, 10, 11, 12, 17
- Dosovitskiy, A., Beyer, L., Kolesnikov, A., Weissenborn, D., Zhai, X., Unterthiner, T., Dehghani, M., Minderer, M., Heigold, G., Gelly, S., Uszkoreit, J., Houlsby, N.: An image is worth 16x16 words: Transformers for image recognition at scale. ICLR (2021) 1, 15
- Girshick, R., Radosavovic, I., Gkioxari, G., Dollár, P., He, K.: Focal loss for dense object detection. https://github.com/facebookresearch/detectron (2018) 11, 17
- 17. Gupta, A., Dollar, P., Girshick, R.: LVIS: A dataset for large vocabulary instance segmentation. In: Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition (2019) 3, 9, 11, 17
- Han, D., Ye, T., Han, Y., Xia, Z., Song, S., Huang, G.: Agent attention: On the integration of softmax and linear attention. arXiv preprint arXiv:2312.08874 (2023)

- 19. He, K., Gkioxari, G., Dollár, P., Girshick, R.: Mask r-cnn (2017) 11, 17
- 20. He, K., Zhang, X., Ren, S., Sun, J.: Deep residual learning for image recognition (2015) 9
- Ho, J., Kalchbrenner, N., Weissenborn, D., Salimans, T.: Axial attention in multidimensional transformers (2019) 12
- 22. Hou, Q., Lu, C.Z., Cheng, M.M., Feng, J.: Conv2former: A simple transformer-style convnet for visual recognition. arXiv preprint arXiv:2211.11943 (2022) 9
- Huang, G., Sun, Y., Liu, Z., Sedra, D., Weinberger, K.: Deep networks with stochastic depth (2016) 10
- 24. Jiao, J., Tang, Y.M., Lin, K.Y., Gao, Y., Ma, J., Wang, Y., Zheng, W.S.: Dilate-former: Multi-scale dilated transformer for visual recognition. IEEE Transactions on Multimedia (2023) 2, 4
- 25. Lan, L., Cai, P., Jiang, L., Liu, X., Li, Y., Zhang, Y.: Brau-net++: U-shaped hybrid cnn-transformer network for medical image segmentation. arXiv preprint arXiv:2401.00722 (2024) 3
- 26. Li, K., Wang, Y., Gao, P., Song, G., Liu, Y., Li, H., Qiao, Y.: Uniformer: Unified transformer for efficient spatiotemporal representation learning (2022) 8, 9, 10, 17
- Liu, S., Chen, T., Chen, X., Chen, X., Xiao, Q., Wu, B., Pechenizkiy, M., Mocanu, D., Wang, Z.: More convnets in the 2020s: Scaling up kernels beyond 51x51 using sparsity. arXiv preprint arXiv:2207.03620 (2022) 9
- Liu, Z., Hu, H., Lin, Y., Yao, Z., Xie, Z., Wei, Y., Ning, J., Cao, Y., Zhang, Z., Dong, L., Wei, F., Guo, B.: Swin transformer v2: Scaling up capacity and resolution. In: International Conference on Computer Vision and Pattern Recognition (CVPR) (2022) 15
- 29. Liu, Z., Lin, Y., Cao, Y., Hu, H., Wei, Y., Zhang, Z., Lin, S., Guo, B.: Swin transformer: Hierarchical vision transformer using shifted windows. In: Proceedings of the IEEE/CVF International Conference on Computer Vision (ICCV) (2021) 1, 4, 7, 8, 9, 10, 11, 12, 13, 15, 17
- 30. Liu, Z., Mao, H., Wu, C.Y., Feichtenhofer, C., Darrell, T., Xie, S.: A convnet for the 2020s. In: Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR). pp. 11976–11986 (June 2022) 9
- 31. Loshchilov, I., Hutter, F.: Decoupled weight decay regularization (2017) 11
- 32. Luo, W., Li, Y., Urtasun, R., Zemel, R.: Understanding the effective receptive field in deep convolutional neural networks. In: Lee, D., Sugiyama, M., Luxburg, U., Guyon, I., Garnett, R. (eds.) Advances in Neural Information Processing Systems. vol. 29. Curran Associates, Inc. (2016), https://proceedings.neurips.cc/paper_files/paper/2016/file/c8067ad1937f728f51288b3eb986afaa-Paper
- 33. Pan, X., Ye, T., Xia, Z., Song, S., Huang, G.: Slide-transformer: Hierarchical vision transformer with local self-attention. In: Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR). pp. 2082–2091 (June 2023)
- 34. Ren, S., Zhou, D., He, S., Feng, J., Wang, X.: Shunted self-attention via multi-scale token aggregation (2021) 8, 9, 10
- Russakovsky, O., Deng, J., Su, H., Krause, J., Satheesh, S., Ma, S., Huang, Z., Karpathy, A., Khosla, A., Bernstein, M., Berg, A.C., Fei-Fei, L.: ImageNet Large Scale Visual Recognition Challenge. International Journal of Computer Vision (IJCV) 115(3), 211–252 (2015). https://doi.org/10.1007/s11263-015-0816-y 3, 9, 10

- 36. Selvaraju, R.R., Cogswell, M., Das, A., Vedantam, R., Parikh, D., Batra, D.: Grad-cam: Visual explanations from deep networks via gradient-based localization (2016). https://doi.org/10.1007/s11263-019-01228-7 13
- Tang, S., Zhang, J., Zhu, S., Tan, P.: Quadtree attention for vision transformers.
 ICLR (2022) 9, 11, 17
- 38. Tolstikhin, İ., Houlsby, N., Kolesnikov, A., Beyer, L., Zhai, X., Unterthiner, T., Yung, J., Steiner, A., Keysers, D., Uszkoreit, J., Lucic, M., Dosovitskiy, A.: Mlpmixer: An all-mlp architecture for vision. arXiv preprint arXiv:2105.01601 (2021) 3
- Touvron, H., Cord, M., Douze, M., Massa, F., Sablayrolles, A., Jegou, H.: Training data-efficient image transformers distillation through attention. In: International Conference on Machine Learning. vol. 139, pp. 10347–10357 (July 2021) 10, 17
- 40. Tu, Z., Talebi, H., Zhang, H., Yang, F., Milanfar, P., Bovik, A., Li, Y.: Maxvit: Multi-axis vision transformer. ECCV (2022) 2, 4, 8, 9, 12
- Vaswani, A., Shazeer, N., Parmar, N., Uszkoreit, J., Jones, L., Gomez, A.N., Kaiser, L.u., Polosukhin, I.: Attention is all you need. In: Guyon, I., Von Luxburg, U., Bengio, S., Wallach, H., Fergus, R., Vishwanathan, S., Garnett, R. (eds.) Advances in Neural Information Processing Systems. vol. 30. Curran Associates, Inc. (2017) 1, 2, 3, 4
- 42. Wang, S., Li, B.Z., Khabsa, M., Fang, H., Ma, H.: Linformer: Self-attention with linear complexity (2020) 4, 9
- 43. Wang, W., Dai, J., Chen, Z., Huang, Z., Li, Z., Zhu, X., Hu, X., Lu, T., Lu, L., Li, H., et al.: Internimage: Exploring large-scale vision foundation models with deformable convolutions. arXiv preprint arXiv:2211.05778 (2022) 9, 10, 11, 17
- 44. Wang, W., Xie, E., Li, X., Fan, D.P., Song, K., Liang, D., Lu, T., Luo, P., Shao, L.: Pyramid vision transformer: A versatile backbone for dense prediction without convolutions. In: Proceedings of the IEEE/CVF International Conference on Computer Vision (ICCV). pp. 568–578 (2021) 1, 9, 10, 11
- 45. Wang, W., Yao, L., Chen, L., Lin, B., Cai, D., He, X., Liu, W.: Crossformer: A versatile vision transformer hinging on cross-scale attention. In: Proceedings of the International Conference on Learning Representations (ICLR) (2022), https://openreview.net/forum?id=_PHymLIxuI 1, 2, 4, 9, 10, 11, 12, 17
- 46. Wu, Y., Kirillov, A., Massa, F., Lo, W.Y., Girshick, R.: Detectron 2. https://github.com/facebookresearch/detectron 2 (2019) 10
- 47. Xia, Z., Pan, X., Song, S., Li, L.E., Huang, G.: Vision transformer with deformable attention. In: Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR). pp. 4794–4803 (June 2022) 2, 3, 4, 8, 9, 10, 11, 12, 15, 17
- 48. Xiao, T., Liu, Y., Zhou, B., Jiang, Y., Sun, J.: Unified perceptual parsing for scene understanding. In: European Conference on Computer Vision. Springer (2018) 10
- Yang, J., Li, C., Zhang, P., Dai, X., Xiao, B., Yuan, L., Gao, J.: Focal attention for long-range interactions in vision transformers. In: Ranzato, M., Beygelzimer, A., Dauphin, Y., Liang, P., Vaughan, J.W. (eds.) Advances in Neural Information Processing Systems. vol. 34, pp. 30008–30022. Curran Associates, Inc. (2021)
- 50. Yang, J., Li, C., Zhang, P., Dai, X., Xiao, B., Yuan, L., Gao, J.: Focal self-attention for local-global interactions in vision transformers (2021) 2
- 51. Yang, R., Ma, H., Wu, J., Tang, Y., Xiao, X., Zheng, M., Li, X.: Scalablevit: Rethinking the context-oriented generalization of vision transformer (2022) 9
- 52. Yun, S., Han, D., Oh, S.J., Chun, S., Choe, J., Yoo, Y.: Cutmix: Regularization strategy to train strong classifiers with localizable features. In: International Conference on Computer Vision (ICCV) (2019) 10

- 53. Zeng, W., Jin, S., Liu, W., Qian, C., Luo, P., Ouyang, W., Wang, X.: Not all tokens are equal: Human-centric visual analysis via token clustering transformer. In: Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition. pp. 11101–11111 (2022) 4
- 54. Zhang, H., Cisse, M., Dauphin, Y.N., Lopez-Paz, D.: mixup: Beyond empirical risk minimization. International Conference on Learning Representations (2018), https://openreview.net/forum?id=r1Ddp1-Rb 10
- 55. Zhou, B., Zhao, H., Puig, X., Xiao, T., Fidler, S., Barriuso, A., Torralba, A.: Semantic understanding of scenes through the ade20k dataset. International Journal of Computer Vision 127(3), 302–321 (2019) 3, 9
- Zhu, L., Wang, X., Ke, Z., Zhang, W., Lau, R.: Biformer: Vision transformer with bi-level routing attention. Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR) (2023) 2, 3, 4, 6, 7, 8, 9, 10, 11, 12, 13, 17