

---

# Orthogonal Transformer: An Efficient Vision Transformer Backbone with Token Orthogonalization

---

Anonymous Author(s)

Affiliation

Address

email

## Abstract

1 We present a general vision transformer backbone, called as Orthogonal Trans-  
2 former, in pursuit of both efficiency and effectiveness. A major challenge for  
3 vision transformer is that self-attention, as the key element in capturing long-range  
4 dependency, is very computationally expensive for dense prediction tasks (e.g.,  
5 object detection). Coarse global self-attention and local self-attention are then  
6 designed to reduce the cost, but they suffer from either neglecting local correlations  
7 or hurting global modeling. We present an orthogonal self-attention mechanism  
8 to alleviate these issues. Specifically, self-attention is computed in the orthogonal  
9 space that is reversible to the spatial domain but has much lower resolution. The  
10 capabilities of learning global dependency and exploring local correlations are  
11 maintained because every orthogonal token in self-attention can attend to the entire  
12 visual tokens. Remarkably, orthogonality is realized by constructing an endoge-  
13 nously orthogonal matrix that is friendly to neural networks and can be optimized  
14 as arbitrary orthogonal matrices. We also introduce Positional MLP to incorporate  
15 position information for arbitrary input resolutions as well as enhance the capacity  
16 of MLP. Finally, we develop a hierarchical architecture for Orthogonal Transformer.  
17 Extensive experiments demonstrate its strong performance on a broad range of  
18 vision tasks, including image classification, object detection, instance segmentation  
19 and semantic segmentation.

## 20 1 Introduction

21 Recently, transformer [45] has made a tremendous success in the field of natural language processing  
22 (NLP). Benefiting from the self-attention (SA) mechanism, it has a strong capacity in building  
23 long-range dependencies in sequential data. Since long-range modeling is also essential for a wide  
24 range of vision tasks, transformer has been adapted into computer vision by converting an image into  
25 a sequence of patches [13] (called as tokens) and achieved competitive performance compared to  
26 CNN [42, 43]. Moreover, self-attention has a quadratic computational complexity to the number of  
27 tokens, resulting in intolerably expensive computational cost for dense prediction tasks, e.g., object  
28 detection and segmentation.

29 Many efforts have been made to design efficient self-attention mechanisms for vision transformer.  
30 PVT [49] employs a pyramid architecture with downsampled key and value tokens to reduce the  
31 computation of global attention. Swin [34] performs self-attention in a local region with shifted  
32 window to allow cross-region connection. GG-Transformer [62] and CrossFormer [51] present dilated  
33 self-attention to learn large-scale features efficiently. As shown in Fig. 1 (a-c), the aforementioned  
34 attention mechanisms can reduce the number of tokens conveniently. However, the penalty is losing  
35 fine-level details for coarse global self-attention (including downsampled and dialted ones) or hurting  
36 long-range modeling for local self-attention [40, 51, 41].

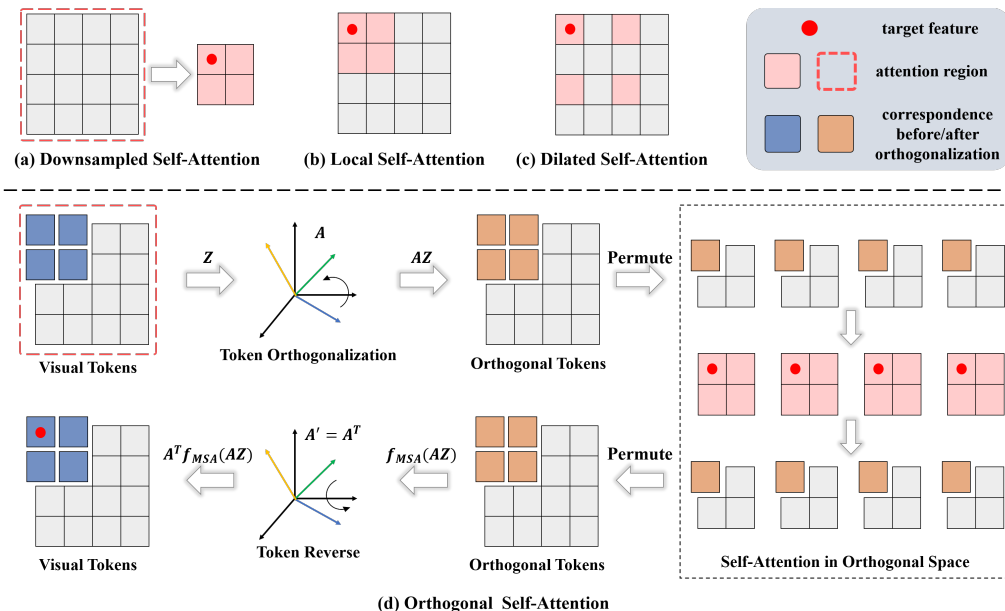


Figure 1: Illustration of different kinds of efficient self-attention mechanisms.

37 This paper presents an orthogonal self-attention (OSA) mechanism to capture global dependency  
 38 without losing fine-level details. As shown in Fig. 1 (d), we orthogonalize tokens within local  
 39 regions, permute them into token groups and perform group-wise self-attention in the orthogonal  
 40 space. We apply Orthogonal Transformation (OT) for the tokens for the following reasons. Firstly,  
 41 tokens has a lower resolution in the orthogonal space than in the original visual space, thus reducing  
 42 the computational cost. Secondly, the orthogonal space is reversible to the visual space without  
 43 information loss (thus outperforming downsampled self-attention), and can be easily reversed back  
 44 since the inverse matrix  $\mathbf{A}'$  is just the transpose of the orthogonal matrix  $\mathbf{A}$ , i.e.,  $\mathbf{A}' = \mathbf{A}^T$ . Thirdly,  
 45 OT can separate tokens into linearly independent groups. Thereby computing self-attention in such  
 46 groups helps to explore different properties in representation. Lastly, OT also builds connections  
 47 among adjacent tokens explicitly and is more capable of modeling local correlations than dilated  
 48 partition.

49 Despite of these benefits, it often needs complex optimization algorithms [30, 31] or imposing  
 50 extra penalties on the loss functions [57, 47] to enforce orthogonality. In this work, we construct  
 51 an endogenously orthogonal matrix that is friendly to neural networks with gradient optimizers.  
 52 We build an orthogonal matrix  $\mathbf{A} \in \mathbf{R}^{n \times n}$  by the product of  $n$  Householder transformations  
 53  $\{\mathbf{H}_i = \mathbf{I} - 2\mathbf{u}_i\mathbf{u}_i^T | i = 0, \dots, n - 1\}$ , where  $\mathbf{u}_i \in \mathbf{R}^n$  is a learnable unit vector. In this way,  $\mathbf{A}$   
 54 is naturally orthogonal and can be optimized as arbitrary orthogonal matrices. In particular, the dilated  
 55 self-attention can be viewed as a degeneration of OSA when  $\mathbf{A} = \mathbf{I}$ . The strong capacity in exploring  
 56 diverse transformations makes OSA more essential for the model to achieve competitive performance.

57 Based on the orthogonal self-attention mechanism, we propose a new Orthogonal Transformer in  
 58 this work. As shown in Fig. 2, it follows the hierarchical design [49, 34] and serves as a general-  
 59 purpose backbone for computer vision. We stack window self-attention (WSA) and orthogonal  
 60 self-attention (OSA) alternatively to capture both global and local dependencies. Note that although  
 61 our OSA is superior in preserving local details, combing it with WSA can further enhance Orthogonal  
 62 Transformer. To improve the flexibility for arbitrary resolutions, we adopt convolutional position  
 63 embedding to incorporate position information. Specifically, Positional MLP (PMLP) is introduced  
 64 by equipping MLP with a depth-wise convolution (DConv) after the non-linear activation. Such  
 65 a simple design not only enables the network to generate position information flexibly, but also  
 66 enhance the capacity of MLP. Besides, it also allows for downsampling within the transformer block  
 67 with strided DConv. We empirically show that downsampling within the block can achieve better  
 68 performance than outside.

69 Extensive experiments demonstrate the strong performance of Orthogonal Transformer on a wide  
 70 range of vision tasks. For example, without extra training data or supervision (such as token label-

71 ing [25] and distillation [42]), our large model Ortho-L achieves 85.4 top-1 accuracy on ImageNet-1K  
72 image classification, surpassing the previous state-of-the-art Swin Transformer [34] by **+1.2** with  
73 similar model size and FLOPs. Our base model Ortho-B achieves 53.0 box AP and 45.9 mask AP  
74 on the COCO detection task, 49.8 mIOU on the ADE20K semantic segmentation task, surpassing  
75 the Swin Transformer counterpart by **+1.2**, **+1.2** and **+2.2**, respectively. Under a smaller setting of  
76 FLOPs, our Ortho-S even obtain larger performance gains, i.e., **+2.1** on ImageNet classification with  
77  $224 \times 224$  resolution, **+1.8** box AP, **+1.6** mask AP on COCO detection and **+4.0** mIoU on ADE20K  
78 segmentation.

## 79 2 Related Work

80 **Vision Transformers.** Transformer network is firstly proposed in the natural language processing  
81 (NLP) field and achieves superior performance in many NLP tasks [45, 26]. Recently, there emerges  
82 a trend towards incorporating the transformer into computer vision tasks, which are previously  
83 dominated by CNNs [18, 39]. ViT [13] is the pioneering work of vision transformers. After that,  
84 there are a large number of works focusing on designing a general vision Transformer backbone [9,  
85 16, 34, 5, 49, 14, 64, 62, 60, 12, 51, 59, 25] for different vision tasks. To adapt transformer to the  
86 image inputs, hierarchical architectures [34, 49], efficient self-attentions [34, 62, 24] and diverse  
87 positional encodings [45, 37, 10, 12] are proposed. Vision transformers are also applied to different  
88 downstream vision tasks, such as object detection [4, 67], semantic segmentation [38, 65, 52], image  
89 restoration [6, 46], and video processing [32, 1]. We propose a new vision transformer backbone to  
90 tackle various vision tasks efficiently and effectively.

91 **Efficient Self-Attentions.** Many efficient self-attention mechanisms [29, 28, 8, 2, 48, 36] have been  
92 proposed in the field of NLP to efficiently handle long sequences. Since the image resolution is usually  
93 high in many vision tasks, the global self-attention can be applied only for a few times and in the  
94 low resolution feature space. Early approaches, such as ViT [13] and DeiT [42], employ large image  
95 patches to restrict the number of tokens. PVT [49] and Twins [9] use downsampled tokens to compute  
96 self-attention. Many works adopt local self-attention to limit the token number by attending only sub-  
97 regions of the input, such as the window attention [34], axis attention [20], and criss-cross attention  
98 [24]. Inspired by dilated convolution [61], GG Transformer [62] and CrossFormer [51] introduce  
99 dilated self-attention to capture long-distance dependency. Unlike existing methods, our orthogonal  
100 self-attention can attend to the entire tokens at a low computational cost. Global dependency can be  
101 captured without losing local details.

102 **Positional Encodings.** Recently, researchers propose to add different kinds of position information to  
103 the token feature or self-attention process. Absolute positional encoding (APE) [45] is the first work  
104 to add position encodings to the Transformer inputs. Relative positional encoding (RPE) [34, 37]  
105 considers the relative distance between two tokens instead of the absolute position. Conditional  
106 positional encoding (CPE) [10] takes feature as inputs and generates position information via a  
107 convolution layer. We follow CPE but compensate position information inside the MLP module. This  
108 design enables MLP to explore local correlations, thus enhancing the model performance.

109 **Orthogonality.** Orthogonality has been widely used in neural networks to regularize neurons [57,  
110 23, 47] or learning disentangled representations [53]. In contrast to existing works, we employ  
111 orthogonality to group visual tokens for efficient self-attention. Besides, most of previous works  
112 depend on complex optimization algorithms, such as singular value decomposition [30] and iterative  
113 approximation [31], or imposing extra regularization terms on the loss functions [57, 47]. Unlike  
114 them, we construct an endogenously orthogonal matrix that is more friendly to neural networks with  
115 gradient optimizers.

## 116 3 Method

### 117 3.1 Overall Architecture

118 An overview of the Orthogonal Transformer architecture is illustrated in Fig. 2. For an input image  
119  $x \in \mathbf{R}^{H \times W \times 3}$ , we follow [56] and adopt several stacked convolution layers to obtain  $\frac{H}{4} \times \frac{W}{4}$  patch  
120 tokens. Similar to Swin [34], the transformer blocks are divided into 4 stages to produce hierarchical  
121 representations. We stack blocks of window self-attention (WSA) and orthogonal self-attention  
122 (OSA) alternatively for better capacity of global and local modeling. For the Orthogonal Transformer

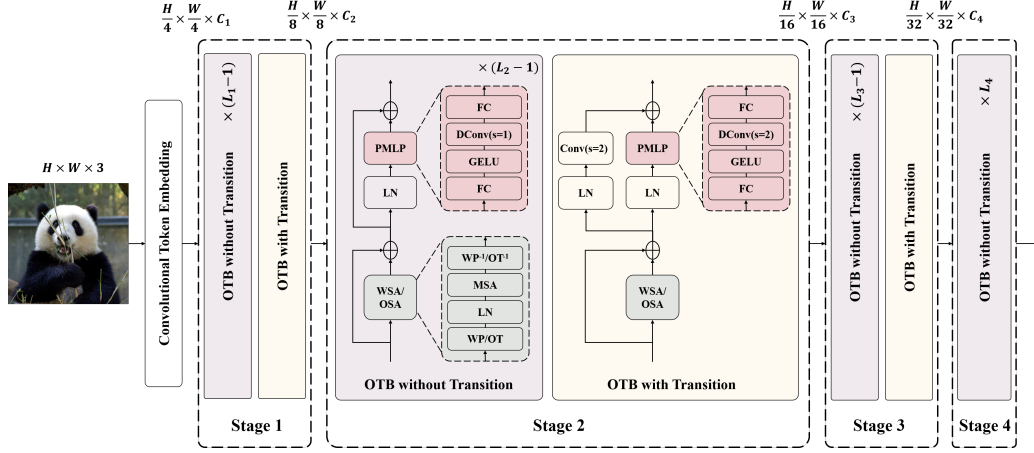


Figure 2: The architecture of Orthogonal Transformer. Orthogonal self-attention (OSA) and window self-attention (WSA) are used successively, where the tokens are grouped via window partition (WP) or orthogonal transformation (OT).  $WP^{-1}$  and  $OT^{-1}$  denote the respective reverse processes.

123 blocks (OTBs), we introduce Positional MLP (PMLP) with an additional depth-wise convolution  
 124 (DConv) after GELU in MLP to provide position information. Based on PMLP, we design two  
 125 kinds of OTBs, one with transition and one without. OTBs with transition aim to reduce the spatial  
 126 resolution by strided DConv and a residual strided convolution. The first three stages are composed  
 127 of stacked OTBs without transition following by a single OTB with transition, while the last stage  
 128 only consists of OTBs without transition.

### 129 3.2 Orthogonal Self-Attention

130 Orthogonal self-attention (OSA) is proposed to enable the transformer layers to encode high-resolution  
 131 images efficiently. Unlike previous works [34, 49, 9, 62], OSA performs self-attention in the  
 132 orthogonal space to capture global dependency without neglecting local correlations. It can cover the  
 133 same global receptive field with much less cost than standard self-attention.

#### 134 3.2.1 Orthogonal Transformation

135 Orthogonal transformation (OT) corresponds to an orthogonal matrix  $\mathbf{A}$  satisfying  $\mathbf{A}^T \mathbf{A} = \mathbf{I}$ , where  
 136  $\mathbf{I}$  is the identity matrix. As shown in Fig. 1, we employ orthogonal transformation  $\mathbf{A}$  to convert visual  
 137 tokens into orthogonal tokens for self-attention computation. Since it is hard to enforce orthogonality  
 138 in neural networks [30, 31], we first introduce a simple yet effective way to construct an endogenously  
 139 orthogonal matrix that can be optimized as arbitrary orthogonal matrices. It is based on the following  
 140 linear algebra theorem [44] (Proof is provided in the appendix).

141 **Theorem 1.** *Every real orthogonal  $n \times n$  matrix is the product of at most  $n$  real orthogonal*  
 142 *Householder transformations.*

143 Given  $n$  tokens, inspired from Theorem 1, we construct an orthogonal matrix  $\mathbf{A} \in \mathbf{R}^{n \times n}$  as:

$$\mathbf{A} = \mathbf{H}_0 \mathbf{H}_1 \cdots \mathbf{H}_{n-1}, \quad (1)$$

144 where  $\mathbf{H}_i = \mathbf{I} - 2\mathbf{u}_i \mathbf{u}_i^T$  is Householder transformation,  $\mathbf{u}_i$  is a unit column vector. Providing a set of  
 145  $n$  random-initialized vectors  $\{\mathbf{v}_i | i = 0, \dots, n-1\}$ , the orthogonal matrix can be reformulated as:

$$\mathbf{A} = (\mathbf{I} - 2 \frac{\mathbf{v}_0 \mathbf{v}_0^T}{\|\mathbf{v}_0\|^2}) (\mathbf{I} - 2 \frac{\mathbf{v}_1 \mathbf{v}_1^T}{\|\mathbf{v}_1\|^2}) \cdots (\mathbf{I} - 2 \frac{\mathbf{v}_{n-1} \mathbf{v}_{n-1}^T}{\|\mathbf{v}_{n-1}\|^2}). \quad (2)$$

146 We set  $\{\mathbf{v}_i | i = 0, \dots, n-1\}$  as learnable parameters, and hence  $\mathbf{A}$  can be optimized as arbitrary  
 147 orthogonal matrices in neural networks with gradient optimizers. Remarkably,  $\mathbf{A}$  can maintain  
 148 orthogonality endogenously throughout the training process without extra regularizer.

### 149 3.2.2 OSA Block

150 To clarify, we first define two terms:

- 151 • **Orthogonal window size**  $m_o$  is the size of the sub-window on which OT is performed. We  
152 perform OT separately for local windows to control its computation complexity.
- 153 • **Orthogonal groups**  $n_o$  is the number of groups into which the tokens are separated by OT. Since  
154 the orthogonal matrix is square,  $n_o = m_o^2$ .

155 Now we elaborate how OSA works in the following three steps, token orthogonalization, self-attention  
156 for orthogonal tokens, token reverse.

157 **Token Orthogonalization.** Given the input feature  $Z \in \mathbf{R}^{(h \times w) \times c}$  (viewed as  $h \times w$  tokens with  
158 the dimension  $c$ ), it is firstly divided into a grid of non-overlapped windows of size  $m_o \times m_o$ , i.e.,  
159  $Z \in \mathbf{R}^{n_w \times n_o \times c}$ , where  $n_w = \frac{h}{m_o} \times \frac{w}{m_o}$  and  $n_o = m_o^2$ .  $Z$  can be viewed as  $n_w$  sub-windows, of  
160 which each consists of  $n_o$  tokens, i.e.,  $Z = \{Z_i \in \mathbf{R}^{n_o \times c} | i = 0, \dots, n_w - 1\}$ . Then we perform  
161 OT separately for  $Z_i$  by multiplying the orthogonal matrix  $\mathbf{A} \in \mathbf{R}^{n_o \times n_o}$  and get orthogonalized  
162 tokens  $\hat{Z}_i = \mathbf{A}Z_i \in \mathbf{R}^{n_o \times c}$ . Finally, we combine the  $n_w$  sub-windows and permute the feature  
163  $\hat{Z} \in \mathbf{R}^{n_w \times n_o \times c}$  into  $n_o$  groups of orthogonal tokens  $\hat{Z}^j \in \mathbf{R}^{n_w \times c}$  (where  $j = 0, \dots, n_o - 1$ ).

164 **Self-attention for orthogonal tokens.** We perform standard multi-head self-attention (MSA) group-  
165 wily for  $\hat{Z}^j$ . Since  $\hat{Z}^j$  has a smaller token number than the input  $Z$  (actually  $\frac{h}{m_o} \times \frac{w}{m_o}$  against  
166  $h \times w$ ), with an appropriate  $m_o$ , performing self-attention in the orthogonal space can significantly  
167 reduce the computation cost. Besides, for OSA, local correlations can be modeled because  $\hat{Z}^j$  is  
168 computed from the adjacent tokens in  $Z$ , and local details are preserved due to the fact that  $\hat{Z}$  and  $Z$   
169 are reversible to each other.

170 **Token reverse.** After the self-attention computation, we employ  $\mathbf{A}^T$ , as the inverse of  $\mathbf{A}$ , to recover  
171 the visual tokens from the enhanced orthogonal representations. The token reverse process is not  
172 repeated herein because it is exactly the inverse process of token orthogonalization except that the  
173 orthogonal matrix is  $\mathbf{A}^T$ .

174 The complete process of OSA can be defined as:

$$f_{OSA}(Z) = \mathbf{A}^T f_{MSA}(f_{LN}(\mathbf{A}Z)) + Z. \quad (3)$$

175 We simplify the process of token orthogonalization as  $\mathbf{A}Z$  for description convenience. Note that  
176 when  $m_o = 1$ , OSA equals to the global self-attention.

177 **Complexity Analysis.** The computation complexity of the standard global self-attention is

$$\Omega(GSA) = 4hwC^2 + 2(hw)^2C. \quad (4)$$

178 For our OSA, the computational complexity is

$$\Omega(OSA) = 4hwC^2 + \frac{1}{n_o} 2(hw)^2C + 2n_o hwC. \quad (5)$$

179 The third term is produced by OT but can be ignored when  $n_o \ll \sqrt{hw}$  (this is usually true for  
180 high-resolution vision tasks). Compared to the standard global self-attention, OSA's computation  
181 cost is reduced significantly from  $O((hw)^2)$  to  $O(\frac{1}{n_o}(hw)^2)$  for high-resolution vision tasks.

### 182 3.3 Window Self-Attention

183 Since Orthogonal Transformer is composed of stacked window self-attentions (WSAs) and OSAs,  
184 we give a brief review of WSA here. WSA first employ window partition (WP) to split the input  
185 feature  $Z \in \mathbf{R}^{(h \times w) \times c}$  into  $\frac{h}{m_w} \times \frac{w}{m_w}$  non-overlapped windows of size  $m_w \times m_w$ , then performs  
186 MSA within each window, and finally reconstruct the tokens from the enhanced representation. The  
187 complete process of WSA is defined as:

$$f_{WSA}(Z) = f_{WP}^{-1}(f_{MSA}(f_{LN}(f_{WP}(Z)))) + Z, \quad (6)$$

188 where  $f_{WP}$  and  $f_{WP}^{-1}$  denote window partition and its reverse, respectively. WSA has a linear  
189 computational complexity to token number  $hw$ :

$$\Omega(WSA) = 4hwC^2 + 2m_w^2 hwC. \quad (7)$$

Table 1: Configurations of Orthogonal Transformer. The FLOPs are measured at resolution  $224 \times 224$ .

Model	Blocks	Channels	Heads	Ratio	Params (M)	FLOPs (G)
Ortho-T	[2, 2, 6, 2]	[32, 64, 160, 256]	[1, 2, 5, 8]	3	3.9	0.7
Ortho-S	[3, 5, 13, 3]	[64, 128, 256, 512]	[2, 4, 8, 16]	4	24	4.5
Ortho-B	[3, 5, 19, 4]	[80, 160, 320, 640]	[2, 4, 8, 16]	4	50	8.6
Ortho-L	[4, 6, 24, 5]	[96, 192, 384, 768]	[3, 6, 12, 24]	4	88	15.4

190 **3.4 Positional MLP**

191 Earlier ViTs adopt absolute position embedding (APE) [45] or relative position embedding (RPE)  
 192 [34, 37] to handle the position information but they cannot be well adapted for arbitrary input  
 193 resolutions [10]. Recently, convolution position embedding (CPE) [10] is proposed to generate  
 194 position information by convolutional layers. In detail, CPE adds position encodings into the input  
 195 tokens outside the transformer blocks. We follow CPE [10] but apply it in a different way. As shown  
 196 in Fig. 2, we introduce Positional MLP (PMLP) to tackle position information in the MLP module.  
 197 We equip PMLP with a depth-wise convolution (DConv) after the non-linear GELU. This design  
 198 can not only provide position information, but also enhance the capacity of MLP in exploring local  
 199 correlations. Based on it, we further develop PMLP with transition by setting the stride of DConv as  
 200 2 and adding a residual strided convolution. PMLP without transition is defined as:

$$f_{PMLP}(Z) = f_{FC}(f_{DConv}^{(s1)}(f_{GELU}(f_{FC}(f_{LN}(Z)))))) + Z, \quad (8)$$

201 while PMLP with transition is defined as

$$f_{PMLP}^{(t)}(Z) = f_{FC}(f_{DConv}^{(s2)}(f_{GELU}(f_{FC}(f_{LN}(Z)))))) + f_{Conv}^{(s2)}(f_{LN}(Z)), \quad (9)$$

202 where  $f_{DConv}^{(s1)}$  and  $f_{DConv}^{(s2)}$  denote depth-wise convolutions with stride 1 and stride 2, respectively,  
 203 and  $f_{Conv}^{(s2)}$  denote a convolution with stride 2.

204 **3.5 Architecture Variants**

205 As shown in Table 1, we build four different Orthogonal Transformer backbones by changing the  
 206 block number and channel number in each stage. Specifically, Ortho-S, Ortho-B and Ortho-L are  
 207 designed to have a similar setting of FLOPs and model size with their Swin Transformer counterparts.  
 208 For all the models, we set the orthogonal window sizes  $m_o$  for OSA in the four stages as 8, 4, 2,  
 209 1, respectively, and the window size for WSA as 7. The expansion ratios in MLP are set as 3 for  
 210 Ortho-T and 4 for the other three variants. For the convolutional patch embedding, we borrow early  
 211 convolutions from [56] and apply 5 convolutions with the same setting of [56]. For the transformer  
 212 blocks, we set the last one in each of the first 3 stages as OTB with transition. The second FC and  
 213 residual convolution change the feature channels. The convolutions reduce the spatial resolution with  
 214 the stride of 2. The kernel sizes of DConv and the residual convolution are set to be  $5 \times 5$  and  $3 \times 3$ ,  
 215 respectively. More details are in the Appendix.

216 **4 Experiments**

217 We conduct experiments on a wide range of vision tasks: image classification on ImageNet-1K [11],  
 218 object detection and instance segmentation on COCO 2017 [33], and semantic segmentation on  
 219 ADE20K [66]. We also take ablation studies to validate the importance of each component. More  
 220 details about experimental settings and extra ablation studies are in the Appendix.

221 **4.1 Image Classification on ImageNet-1K**

222 **Experimental Settings.** We benchmark our Orthogonal Transformer on ImageNet-1K [11] image  
 223 classification. We follow the same training strategy in DeiT [42] and adopt the strong data augmenta-  
 224 tion and regularization, except for repeated augmentation [21] that does not improve performance.  
 225 For a fair comparison, we do not use extra training data or extra supervision techniques, such as token  
 226 labeling [25] and distillation [42], because most of previous works didn't use them. All our models  
 227 are trained from scratch for 300 epochs with the input size of  $224 \times 224$ . An AdamW optimizer

Table 2: Comparison with the state-of-the-art on ImageNet-1K classification.

Model	#Param (M)	FLOPs (G)	Acc. (%)	Model	#Param (M)	FLOPs (G)	Acc. (%)
DeiT-Ti [42]	5.7	1.3	72.2	PVT-L [49]	61	9.8	81.7
PVTv2-b0 [50]	3.7	0.6	70.5	T2T-24 [63]	64	13.2	82.2
T2T-7 [63]	4.3	1.1	71.7	Swin-S [34]	50	8.7	83.0
QuadTree-B-b0 [40]	3.5	0.7	72.0	CvT-21 [54]	32	7.1	82.5
Ortho-T	3.9	0.7	<b>74.0</b>	CaiT-s24 [43]	47	9.4	82.7
RegNetY-4G [35]	21	4.0	80.0	Focal-S [60]	51	9.1	83.5
DeiT-S [42]	22	4.6	79.9	CrossFormer-B [51]	52	9.2	83.4
PVT-S [49]	25	3.8	79.8	RegionViT-M [5]	41	7.4	83.1
T2T-14 [63]	22	5.2	80.7	Ortho-B	50	8.6	<b>84.0</b>
Swin-T [34]	29	4.5	81.3	CvT-21 $\uparrow$ 384 [54]	32	25.0	84.9
Twins-SVT-S [9]	24	2.9	81.7	CaiT-s24 $\uparrow$ 384 [43]	47	32.2	84.3
CvT-13 [54]	20	4.5	81.6	Ortho-B $\uparrow$ 384	50	26.6	<b>85.2</b>
CaiT-xs24 [43]	27	5.4	81.8	DeiT-B [42]	86	17.5	81.8
Focal-T [60]	29	4.9	82.2	Swin-B [34]	88	15.4	83.3
CrossFormer-S [51]	31	4.9	82.5	CaiT-s48 [43]	90	18.6	83.5
RegionViT-S [5]	31	5.3	82.6	Focal-B [60]	90	16.0	83.8
Container [15]	22	8.1	82.7	CrossFormer-L [51]	92	16.1	84.0
QuadTree-B-b2 [40]	24	4.5	82.7	RegionViT-B [5]	73	13.0	83.2
Ortho-S	24	4.5	<b>83.4</b>	Ortho-L	88	15.4	<b>84.2</b>
CvT-13 $\uparrow$ 384 [54]	20	16.3	83.0	Swin-B $\uparrow$ 384 [34]	88	47.0	84.2
CaiT-xs24 $\uparrow$ 384 [43]	27	19.3	83.8	CaiT-s48 $\uparrow$ 384 [43]	90	63.8	85.1
Ortho-S $\uparrow$ 384	24	14.3	<b>84.8</b>	Ortho-L $\uparrow$ 384	88	47.4	<b>85.4</b>

Table 3: Object detection and instance segmentation with Mask R-CNN on COCO val2017. FLOPs are measured at resolution  $800 \times 1280$ . All the models are pre-trained on ImageNet-1K.

Backbone	#Param (M)	FLOPs (G)	Mask R-CNN $1 \times$ schedule							Mask R-CNN $3 \times +$ MS schedule				
			$AP^b$	$AP_{50}^b$	$AP_{75}^b$	$AP^m$	$AP_{50}^m$	$AP_{75}^m$	$AP^b$	$AP_{50}^b$	$AP_{75}^b$	$AP^m$	$AP_{50}^m$	$AP_{75}^m$
Res50 [19]	44	260	38.0	58.6	41.4	34.4	55.1	36.7	41.0	61.7	44.9	37.1	58.4	40.1
PVT-S [49]	44	245	40.4	62.9	43.8	37.8	60.1	40.3	43.0	65.3	46.9	39.9	62.5	42.8
Twins-S [9]	44	228	43.4	66.0	47.3	40.3	63.2	43.4	46.8	69.2	51.2	42.6	66.3	45.8
Swin-T [34]	48	264	42.2	64.6	46.2	39.1	61.6	42.0	46.0	68.2	50.2	41.6	65.1	44.8
ViL-S [64]	45	218	44.9	67.1	49.3	41.0	64.2	44.1	47.1	68.7	51.5	42.7	65.9	46.2
Focal-T [60]	49	291	44.8	67.7	49.2	41.0	64.7	44.2	47.2	69.4	51.9	42.7	66.5	45.9
RegionViT-S [5]	51	183	44.2	-	-	40.8	-	-	47.6	-	-	43.4	-	-
Ortho-S	44	277	<b>47.0</b>	<b>69.4</b>	<b>51.3</b>	<b>42.5</b>	<b>66.1</b>	<b>45.7</b>	<b>48.7</b>	<b>70.5</b>	<b>53.3</b>	<b>43.6</b>	<b>67.3</b>	<b>47.3</b>
Res101 [19]	63	336	40.4	61.1	44.2	36.4	57.7	38.8	42.8	63.2	47.1	38.5	60.1	41.3
X101-32 [58]	63	340	41.9	62.5	45.9	37.5	59.4	40.2	44.0	64.4	48.0	39.2	61.4	41.9
PVT-M [49]	64	302	42.0	64.4	45.6	39.0	61.6	42.1	44.2	66.0	48.2	40.5	63.1	43.5
Twins-B [9]	76	340	45.2	67.6	49.3	41.5	64.5	44.8	48.0	69.5	52.7	43.0	66.8	46.6
Swin-S [34]	69	354	44.8	66.6	48.9	40.9	63.4	44.2	48.5	70.2	53.5	43.3	67.3	46.6
Focal-S [60]	71	401	47.4	69.8	51.9	42.8	66.6	46.1	48.8	70.5	53.6	43.8	67.7	47.2
RegionViT-B [5]	93	307	45.4	-	-	41.6	-	-	48.3	-	-	43.5	-	-
Ortho-B	69	372	<b>48.3</b>	<b>70.5</b>	<b>53.0</b>	<b>43.3</b>	<b>67.3</b>	<b>46.5</b>	<b>49.9</b>	<b>71.4</b>	<b>54.8</b>	<b>44.3</b>	<b>68.6</b>	<b>47.9</b>

228 with a cosine decay learning rate scheduler and 5 epochs of linear warm-up is employed. The initial  
 229 learning rate, weight decay, and batch-size are 0.001, 0.05, and 1024, respectively. For Ortho-T, we  
 230 use a smaller weight decay of 0.01. The maximum rates of increasing stochastic depth [22] are set  
 231 as 0.1/0.2/0.4/0.5 for the models from tiny to large. For the results of  $384 \times 384$ , we fine-tune the  
 232 models for 30 epochs with learning rate of  $1e-5$ , weight decay of  $1e-8$  and batch-size of 512.

233 **Results.** Table 2 reports the image classification results on ImageNet-1K. It clearly shows that our  
 234 Orthogonal Transformer has a stronger performance than previous models under similar settings  
 235 of FLOPs and model size. Specifically, the tiny model Ortho-T achieves a 74.0% Top-1 accuracy  
 236 with only 0.7G FLOPs, surpassing QuadTree-B-b0, T2T-7 and PVTv2-b0 by 2%, 2.3% and 3.5%,  
 237 respectively. The small model Ortho-S achieves the same accuracy as base model CrossFormer-B  
 238 with 51% fewer FLOPs. As for  $384 \times 384$  input size, the large model Ortho-L achieves an accuracy  
 239 of 85.4%, surpassing Swin-B by 1.2% with similar model size, and outperforming CaiT-s48 with  
 240 26% fewer FLOPs. These quantitative comparisons with SOTA methods demonstrate the efficiency  
 241 and effectiveness of Orthogonal Transformer.

Table 4: Object detection and instance segmentation with Cascade Mask R-CNN on COCO val2017.

Method	#Params (M)	FLOPs (G)	3× + MS schedule					
			$AP^b$	$AP_{50}^b$	$AP_{75}^b$	$AP^m$	$AP_{50}^m$	$AP_{75}^m$
DeiT [42]	80	889	48.0	67.2	51.7	41.4	64.2	44.3
Swin-T [34]	86	745	50.5	69.3	54.9	43.7	66.6	47.1
Focal-T [60]	87	770	51.5	70.6	55.9	-	-	-
Ortho-S	81	755	<b>52.3</b>	<b>71.3</b>	<b>56.8</b>	<b>45.3</b>	<b>68.6</b>	<b>49.2</b>
X101-32 [58]	101	819	48.1	66.5	52.4	41.6	63.9	45.2
Swin-S [34]	107	838	51.8	70.4	56.3	44.7	67.9	48.5
Ortho-B	107	851	<b>53.0</b>	<b>72.0</b>	<b>57.4</b>	<b>45.9</b>	<b>69.4</b>	<b>49.9</b>

Table 5: Semantic segmentation with Semantic FPN and Upernet on ADE20K. The FLOPs are measured at resolution 512×2048.

Backbone	Semantic FPN 80k			Backbone	Upernet 160k			
	#Param (M)	FLOPs (G)	mIoU (%)		#Param (M)	FLOPs (G)	mIoU (%)	MS mIoU (%)
Res50 [19]	29	183	36.7	TwinsP-S [9]	55	919	46.2	47.5
PVT-S [49]	28	161	39.8	Twins-S [9]	54	901	46.2	47.1
TwinsP-S [9]	28	162	44.3	Swin-T [34]	60	945	44.5	45.8
Swin-T [34]	32	182	41.5	Focal-T [60]	62	998	45.8	47.0
Ortho-S	28	195	<b>48.2</b>	Ortho-S	54	956	<b>48.5</b>	<b>49.9</b>
Res101 [19]	48	260	38.8	Res101 [19]	86	1029	-	44.9
PVT-L [49]	65	283	42.1	Twins-B [9]	89	1020	47.7	48.9
TwinsP-L [9]	65	283	46.4	Swin-S [34]	81	1038	47.6	49.5
Swin-S [34]	53	274	45.2	Focal-T [60]	85	1130	48.0	50.0
Ortho-B	53	297	<b>49.0</b>	Ortho-B	81	1057	<b>49.8</b>	<b>51.2</b>

242 **4.2 Object Detection on COCO**

243 **Experimental Settings.** Experiments on object detection and instance segmentation are conducted  
 244 on COCO 2017 dataset [33]. Following [34], we use our transformer models as the backbone network,  
 245 and use Mask-RCNN [17] and Cascaded Mask R-CNN [3] as the detection and segmentation heads.  
 246 For both tasks, the backbones are pretrained on ImageNet-1K, and fine-tuned on the COCO training  
 247 set with AdamW optimizer. We take experiments on two common settings: "1 ×" (12 training epochs)  
 248 and "3 × +MS" (36 training epochs with multi-scale training). The configurations follow the setting  
 249 of [34] and are implemented with MMDetection [7].

250 **Results.** Table 3 and Table 4 show the results with Mask R-CNN and Cascade Mask R-CNN,  
 251 respectively. The results show that our method achieves the best performance in all comparisons. For  
 252 vision transformers with Mask R-CNN framework, our Ortho-S outperforms Focal-T (having similar  
 253 model size and computational cost) by **+2.2** box AP, **+1.5** mask AP with the 1× schedule and **+1.5**  
 254 box AP, **+0.9** mask AP with the 3× multi-scale learning schedule. The results with Cascade Mask  
 255 R-CNN also show that our Orthogonal Transformer exceeds the counterparts by evident margins.

256 **4.3 Semantic Segmentation on ADE20K**

257 **Experimental Settings.** Experiments on semantic segmentation are conducted on the ADE20K  
 258 dataset[66]. We employ Semantic FPN [27] and Upernet [55] as the segmentation heads and replace  
 259 the backbones with our Orthogonal Transformer. For a fair comparison, we follow the same setting  
 260 of PVT [49] to train Semantic FPN 80k iterations and apply the setting of Swin [34] to train Upernet  
 261 for 160k iterations. More details are provided in the Appendix.

262 **Results.** The results on semantic segmentation are listed in Table 5. Our method obtains large  
 263 performance gain over other vision transformers. For methods that using Semantic FPN as the  
 264 segmentation head, Orthogonal Transformer surpasses Twins by **+3.9**, **+2.6** on mIoU. The results  
 265 with Upernet indicate that Orthogonal Transformer achieves **+2.9**, **+1.2** absolutely higher Multi-



Table 6: Ablation Study on three benchmarks using the Ortho-S backbone.

	ImageNet-1K			COCO		ADE20K
	#Param (M)	FLOPs (G)	Acc. (%)	AP <sup>b</sup>	AP <sup>m</sup>	mIoU
window sa	24.0	4.5	82.6	46.2	42.0	47.5
dilated sa	24.0	4.5	82.9	46.2	41.7	46.6
ortho. sa	24.0	4.6	83.2	<b>47.0</b>	42.3	48.0
window/ortho. sa	24.0	4.5	<b>83.4</b>	<b>47.0</b>	<b>42.5</b>	<b>48.2</b>
no pos.	23.4	4.4	82.2	44.0	40.4	44.6
abs. pos.	23.4	4.4	82.2	44.1	40.6	44.4
rel. pos.	23.4	4.4	82.4	43.4	40.1	44.6
conv. pos.	24.0	4.5	<b>83.4</b>	<b>47.0</b>	<b>42.5</b>	<b>48.2</b>
w/o early convs	23.9	4.2	82.6	46.5	41.9	47.1
early convs	24.0	4.5	<b>83.4</b>	<b>47.0</b>	<b>42.5</b>	<b>48.2</b>
outside transition	24.0	4.6	83.1	46.7	42.4	48.0
inside transition	24.0	4.5	<b>83.4</b>	<b>47.0</b>	<b>42.5</b>	<b>48.2</b>

266 Scale(MS) mIOU than Focal Transformer, and reports **51.2** MS mIOU with 81M parameters. The  
 267 comparisons on downstream tasks further validate the superiority of our method.

#### 268 4.4 Ablation Study

269 In this section, we conduct ablation study to inspect the respective roles of each part in Orthogonal  
 270 Transformer, using ImageNet-1K image classification, Mask R-CNN (1×) on COCO object detection,  
 271 and Semantic FPN on ADE20K semantic segmentation.

272 **Self-Attention Mechanism.** Ablations of the orthogonal self-attention (SA) are reported in Table 6.  
 273 Orthogonal SA outperforms window SA and dilated SA consistently on the three tasks. Since window  
 274 SA and dilated SA belong to fine-grained local SAs and coarse-grained global SAs respectively,  
 275 the results may validate the superiority of orthogonal SA in capturing global dependency while  
 276 preserving local details. Besides, replacing half of orthogonal SA with window SA can lead to slight  
 277 performance gain with 0.1G decrease of FLOPs, justifying the design of Orthogonal Transformer.

278 **Position Encoding.** Comparisons of different position embedding methods are reported in Table 6.  
 279 Convolution position embedding (CPE) in PMLP performs better than those without position encoding  
 280 and those with APE or RPE. The gains are specifically large on downstream high-resolution tasks,  
 281 indicating that the design of our CPE is more suitable for arbitrary large input resolutions.

282 **Patch Embedding.** We borrow the idea of early convolutions from [56] to perform overlapped patch  
 283 embedding. The results in Table 6 indicate that Ortho-S with early convolutions outperforms the  
 284 counterpart with non-overlapped patch embedding.

285 **Transition.** Most of previous works using hierarchical architectures [34, 49] perform the transition  
 286 of token number and feature dimension between stages. In this work, we adopt a different manner  
 287 and employ the strided convolutions in PMLP to perform transition in the transformer block. The  
 288 results in Table 6 imply that Ortho-S with inside transition can yield performance gain as well as  
 289 computation reduction compared to the counterpart with outside transition.

## 290 5 Conclusion

291 We present a new vision transformer, called as Orthogonal Transformer, to serve as a strong general  
 292 backbone for computer vision. The orthogonal self-attention mechanism is introduced to capture  
 293 global information while exploring local correlations. Positional MLP is developed to handle position  
 294 information within the transformer block for arbitrary input resolutions. Orthogonal Transformer  
 295 achieves state-of-the-art performance in various vision tasks, including image classification, ob-  
 296 ject detection, instance segmentation and semantic segmentation. We expect to apply Orthogonal  
 297 Transformer for more vision tasks, such as image processing and video prediction.

298 **References**

- 299 [1] Anurag Arnab, Mostafa Dehghani, Georg Heigold, Chen Sun, Mario Lučić, and Cordelia Schmid. Vivit: A  
300 video vision transformer. In *ICCV*, pages 6836–6846, 2021.
- 301 [2] Iz Beltagy, Matthew E Peters, and Arman Cohan. Longformer: The long-document transformer. In  
302 *EMNLP*, 2020.
- 303 [3] Zhaowei Cai and Nuno Vasconcelos. Cascade r-cnn: Delving into high quality object detection. In *CVPR*,  
304 pages 6154–6162, 2018.
- 305 [4] Nicolas Carion, Francisco Massa, Gabriel Synnaeve, Nicolas Usunier, Alexander Kirillov, and Sergey  
306 Zagoruyko. End-to-end object detection with transformers. In *ECCV*, pages 213–229, 2020.
- 307 [5] Chun-Fu Chen, Rameswar Panda, and Quanfu Fan. Regionvit: Regional-to-local attention for vision  
308 transformers. In *ICLR*, 2022.
- 309 [6] Hanting Chen, Yunhe Wang, Tianyu Guo, Chang Xu, Yiping Deng, Zhenhua Liu, Siwei Ma, Chunjing Xu,  
310 Chao Xu, and Wen Gao. Pre-trained image processing transformer. In *CVPR*, pages 12299–12310, 2021.
- 311 [7] Kai Chen, Jiaqi Wang, Jiangmiao Pang, Yuhang Cao, Yu Xiong, Xiaoxiao Li, Shuyang Sun, Wansen Feng,  
312 Ziwei Liu, Jiarui Xu, et al. Mmdetection: Open mmlab detection toolbox and benchmark. *arXiv preprint*  
313 *arXiv:1906.07155*, 2019.
- 314 [8] Krzysztof Choromanski, Valerii Likhoshesterov, David Dohan, Xingyou Song, Andreea Gane, Tamas Sarlos,  
315 Peter Hawkins, Jared Davis, David Belanger, Lucy Colwell, et al. Masked language modeling for proteins  
316 via linearly scalable long-context transformers. In *ICLR*, 2020.
- 317 [9] Xiangxiang Chu, Zhi Tian, Yuqing Wang, Bo Zhang, Haibing Ren, Xiaolin Wei, Huaxia Xia, and Chunhua  
318 Shen. Twins: Revisiting the design of spatial attention in vision transformers. In *NeurIPS*, 2021.
- 319 [10] Xiangxiang Chu, Zhi Tian, Bo Zhang, Xinlong Wang, Xiaolin Wei, Huaxia Xia, and Chunhua Shen.  
320 Conditional positional encodings for vision transformers. *arXiv preprint arXiv:2102.10882*, 2021.
- 321 [11] Jia Deng, Wei Dong, Richard Socher, Li-Jia Li, Kai Li, and Li Fei-Fei. Imagenet: A large-scale hierarchical  
322 image database. In *CVPR*, pages 248–255, 2009.
- 323 [12] Xiaoyi Dong, Jianmin Bao, Dongdong Chen, Weiming Zhang, Nenghai Yu, Lu Yuan, Dong Chen, and  
324 Baining Guo. Cswin transformer: A general vision transformer backbone with cross-shaped windows. In  
325 *CVPR*, 2022.
- 326 [13] Alexey Dosovitskiy, Lucas Beyer, Alexander Kolesnikov, Dirk Weissenborn, Xiaohua Zhai, Thomas  
327 Unterthiner, Mostafa Dehghani, Matthias Minderer, Georg Heigold, Sylvain Gelly, et al. An image is worth  
328 16x16 words: Transformers for image recognition at scale. In *ICLR*, 2021.
- 329 [14] Haoqi Fan, Bo Xiong, Karttikeya Mangalam, Yanghao Li, Zhicheng Yan, Jitendra Malik, and Christoph  
330 Feichtenhofer. Multiscale vision transformers. In *ICCV*, pages 6824–6835, 2021.
- 331 [15] Peng Gao, Jiasen Lu, Hongsheng Li, Roozbeh Mottaghi, and Aniruddha Kembhavi. Container: Context  
332 aggregation networks. In *NeurIPS*, 2021.
- 333 [16] Kai Han, An Xiao, Enhua Wu, Jianyuan Guo, Chunjing Xu, and Yunhe Wang. Transformer in transformer.  
334 *NeurIPS*, 34, 2021.
- 335 [17] Kaiming He, Georgia Gkioxari, Piotr Dollár, and Ross Girshick. Mask r-cnn. In *ICCV*, pages 2961–2969,  
336 2017.
- 337 [18] Kaiming He, Xiangyu Zhang, Shaoqing Ren, and Jian Sun. Deep residual learning for image recognition.  
338 In *CVPR*, pages 770–778, 2016.
- 339 [19] Kaiming He, X. Zhang, Shaoqing Ren, and Jian Sun. Deep residual learning for image recognition. *CVPR*,  
340 pages 770–778, 2016.
- 341 [20] Jonathan Ho, Nal Kalchbrenner, Dirk Weissenborn, and Tim Salimans. Axial attention in multidimensional  
342 transformers. *arXiv preprint arXiv:1912.12180*, 2019.
- 343 [21] Elad Hoffer, Tal Ben-Nun, Itay Hubara, Niv Giladi, Torsten Hoefler, and Daniel Soudry. Augment your  
344 batch: Improving generalization through instance repetition. In *CVPR*, pages 8129–8138, 2020.
- 345 [22] Gao Huang, Yu Sun, Zhuang Liu, Daniel Sedra, and Kilian Q Weinberger. Deep networks with stochastic  
346 depth. In *ECCV*, pages 646–661, 2016.
- 347 [23] Lei Huang, Xianglong Liu, Bo Lang, Adams Wei Yu, Yongliang Wang, and Bo Li. Orthogonal weight  
348 normalization: Solution to optimization over multiple dependent stiefel manifolds in deep neural networks.  
349 In *AAAI*, 2018.
- 350 [24] Zilong Huang, Xinggang Wang, Lichao Huang, Chang Huang, Yunchao Wei, and Wenyu Liu. Ccnet:  
351 Criss-cross attention for semantic segmentation. In *ICCV*, pages 603–612, 2019.
- 352 [25] Zi-Hang Jiang, Qibin Hou, Li Yuan, Daquan Zhou, Yujun Shi, Xiaojie Jin, Anran Wang, and Jiashi Feng.  
353 All tokens matter: Token labeling for training better vision transformers. In *NeurIPS*.
- 354 [26] Jacob Devlin Ming-Wei Chang Kenton and Lee Kristina Toutanova. Bert: Pre-training of deep bidirectional  
355 transformers for language understanding. In *NAACL-HLT*, pages 4171–4186, 2019.
- 356 [27] Alexander Kirillov, Ross Girshick, Kaiming He, and Piotr Dollár. Panoptic feature pyramid networks. In  
357 *CVPR*, pages 6399–6408, 2019.
- 358 [28] Nikita Kitaev, Lukasz Kaiser, and Anselm Levskaya. Reformer: The efficient transformer. In *ICLR*, 2019.

- 359 [29] Juho Lee, Yoonho Lee, Jungtaek Kim, Adam Kosiosek, Seungjin Choi, and Yee Whye Teh. Set transformer:  
360 A framework for attention-based permutation-invariant neural networks. In *ICML*, pages 3744–3753, 2019.
- 361 [30] José Lezama, Qiang Qiu, Pablo Musé, and Guillermo Sapiro. Ole: Orthogonal low-rank embedding-a plug  
362 and play geometric loss for deep learning. In *CVPR*, pages 8109–8118, 2018.
- 363 [31] Jun Li, Fuxin Li, and Sinisa Todorovic. Efficient riemannian optimization on the stiefel manifold via the  
364 cayley transform. In *ICLR*, 2019.
- 365 [32] Kunchang Li, Yali Wang, Peng Gao, Guanglu Song, Yu Liu, Hongsheng Li, and Yu Qiao. Uniformer:  
366 Unified transformer for efficient spatiotemporal representation learning. In *ICLR*, 2022.
- 367 [33] Tsung-Yi Lin, Michael Maire, Serge Belongie, James Hays, Pietro Perona, Deva Ramanan, Piotr Dollár,  
368 and C Lawrence Zitnick. Microsoft coco: Common objects in context. In *ECCV*, pages 740–755, 2014.
- 369 [34] Ze Liu, Yutong Lin, Yue Cao, Han Hu, Yixuan Wei, Zheng Zhang, Stephen Lin, and Baining Guo. Swin  
370 transformer: Hierarchical vision transformer using shifted windows. In *ICCV*, pages 10012–10022, 2021.
- 371 [35] Ilija Radosavovic, Raj Prateek Kosaraju, Ross Girshick, Kaiming He, and Piotr Dollár. Designing network  
372 design spaces. In *CVPR*, pages 10428–10436, 2020.
- 373 [36] Aurko Roy, Mohammad Saffar, Ashish Vaswani, and David Grangier. Efficient content-based sparse  
374 attention with routing transformers. *TACL*, 9:53–68, 2021.
- 375 [37] Peter Shaw, Jakob Uszkoreit, and Ashish Vaswani. Self-attention with relative position representations. In  
376 *NAACL-HLT (2)*, 2018.
- 377 [38] Robin Strudel, Ricardo Garcia, Ivan Laptev, and Cordelia Schmid. Segmenter: Transformer for semantic  
378 segmentation. In *ICCV*, pages 7262–7272, 2021.
- 379 [39] Mingxing Tan and Quoc Le. Efficientnet: Rethinking model scaling for convolutional neural networks. In  
380 *ICML*, pages 6105–6114. PMLR, 2019.
- 381 [40] Shitao Tang, Jiahui Zhang, Siyu Zhu, and Ping Tan. Quadtree attention for vision transformers. In *ICLR*,  
382 2022.
- 383 [41] Yi Tay, Mostafa Dehghani, Dara Bahri, and Donald Metzler. Efficient transformers: A survey. *arXiv*  
384 *preprint arXiv:2009.06732*, 2020.
- 385 [42] Hugo Touvron, Matthieu Cord, Matthijs Douze, Francisco Massa, Alexandre Sablayrolles, and Hervé Jégou.  
386 Training data-efficient image transformers & distillation through attention. In *ICML*, pages 10347–10357,  
387 2021.
- 388 [43] Hugo Touvron, Matthieu Cord, Alexandre Sablayrolles, Gabriel Synnaeve, and Hervé Jégou. Going deeper  
389 with image transformers. In *ICCV*, pages 32–42, 2021.
- 390 [44] Frank Uhlig. Constructive ways for generating (generalized) real orthogonal matrices as products of  
391 (generalized) symmetries. *Linear Algebra and its Applications*, 332:459–467, 2001.
- 392 [45] Ashish Vaswani, Noam Shazeer, Niki Parmar, Jakob Uszkoreit, Llion Jones, Aidan N Gomez, Łukasz  
393 Kaiser, and Illia Polosukhin. Attention is all you need. In *NeurIPS*, pages 5998–6008, 2017.
- 394 [46] Ziyu Wan, Jingbo Zhang, Dongdong Chen, and Jing Liao. High-fidelity pluralistic image completion with  
395 transformers. In *ICCV*, pages 4692–4701, 2021.
- 396 [47] Jiayun Wang, Yubei Chen, Rudrasis Chakraborty, and Stella X Yu. Orthogonal convolutional neural  
397 networks. In *CVPR*, pages 11505–11515, 2020.
- 398 [48] Sinong Wang, Belinda Z Li, Madian Khabsa, Han Fang, and Hao Ma. Linformer: Self-attention with linear  
399 complexity. *arXiv preprint arXiv:2006.04768*, 2020.
- 400 [49] Wenhai Wang, Enze Xie, Xiang Li, Deng-Ping Fan, Kaitao Song, Ding Liang, Tong Lu, Ping Luo, and  
401 Ling Shao. Pyramid vision transformer: A versatile backbone for dense prediction without convolutions.  
402 In *ICCV*, pages 568–578, 2021.
- 403 [50] Wenhai Wang, Enze Xie, Xiang Li, Deng-Ping Fan, Kaitao Song, Ding Liang, Tong Lu, Ping Luo, and  
404 Ling Shao. Pvt v2: Improved baselines with pyramid vision transformer. *CVM*, pages 1–10, 2022.
- 405 [51] Wenxiao Wang, Lu Yao, Long Chen, Binbin Lin, Deng Cai, Xiaofei He, and Wei Liu. Crossformer: A  
406 versatile vision transformer hinging on cross-scale attention. In *ICLR*, 2022.
- 407 [52] Yuqing Wang, Zhaoliang Xu, Xinlong Wang, Chunhua Shen, Baoshan Cheng, Hao Shen, and Huaxia Xia.  
408 End-to-end video instance segmentation with transformers. In *CVPR*, pages 8741–8750, 2021.
- 409 [53] Yuxiang Wei, Yupeng Shi, Xiao Liu, Zhilong Ji, Yuan Gao, Zhongqin Wu, and Wangmeng Zuo. Orthogonal  
410 jacobian regularization for unsupervised disentanglement in image generation. In *ICCV*, pages 6721–6730,  
411 2021.
- 412 [54] Haiping Wu, Bin Xiao, Noel Codella, Mengchen Liu, Xiyang Dai, Lu Yuan, and Lei Zhang. Cvt:  
413 Introducing convolutions to vision transformers. In *ICCV*, pages 22–31, 2021.
- 414 [55] Tete Xiao, Yingcheng Liu, Bolei Zhou, Yuning Jiang, and Jian Sun. Unified perceptual parsing for scene  
415 understanding. In *ECCV*, pages 418–434, 2018.
- 416 [56] Tete Xiao, Mannat Singh, Eric Mintun, Trevor Darrell, Piotr Dollár, and Ross Girshick. Early convolutions  
417 help transformers see better. *NeurIPS*, 34:30392–30400, 2021.
- 418 [57] Di Xie, Jiang Xiong, and Shiliang Pu. All you need is beyond a good init: Exploring better solution for  
419 training extremely deep convolutional neural networks with orthonormality and modulation. In *CVPR*,  
420 pages 6176–6185, 2017.

- 421 [58] Saining Xie, Ross B. Girshick, Piotr Dollár, Zhuowen Tu, and Kaiming He. Aggregated residual transfor-  
422 mations for deep neural networks. *CVPR*, pages 5987–5995, 2017.
- 423 [59] Weijian Xu, Yifan Xu, Tyler Chang, and Zhuowen Tu. Co-scale conv-attentional image transformers. In  
424 *ICCV*, pages 9981–9990, 2021.
- 425 [60] Jianwei Yang, Chunyuan Li, Pengchuan Zhang, Xiyang Dai, Bin Xiao, Lu Yuan, and Jianfeng Gao. Focal  
426 self-attention for local-global interactions in vision transformers. In *NeurIPS*, 2021.
- 427 [61] Fisher Yu and Vladlen Koltun. Multi-scale context aggregation by dilated convolutions. *arXiv preprint*  
428 *arXiv:1511.07122*, 2015.
- 429 [62] Qihang Yu, Yingda Xia, Yutong Bai, Yongyi Lu, Alan L Yuille, and Wei Shen. Glimpse-and-gaze vision  
430 transformer. *NeurIPS*, 34, 2021.
- 431 [63] Li Yuan, Yungpeng Chen, Tao Wang, Weihao Yu, Yujun Shi, Zi-Hang Jiang, Francis EH Tay, Jiashi Feng,  
432 and Shuicheng Yan. Tokens-to-token vit: Training vision transformers from scratch on imagenet. In *ICCV*,  
433 pages 558–567, 2021.
- 434 [64] Pengchuan Zhang, Xiyang Dai, Jianwei Yang, Bin Xiao, Lu Yuan, Lei Zhang, and Jianfeng Gao. Multi-  
435 scale vision longformer: A new vision transformer for high-resolution image encoding. In *ICCV*, pages  
436 2998–3008, 2021.
- 437 [65] Sixiao Zheng, Jiachen Lu, Hengshuang Zhao, Xiatian Zhu, Zekun Luo, Yabiao Wang, Yanwei Fu, Jianfeng  
438 Feng, Tao Xiang, Philip HS Torr, et al. Rethinking semantic segmentation from a sequence-to-sequence  
439 perspective with transformers. In *CVPR*, pages 6881–6890, 2021.
- 440 [66] Bolei Zhou, Hang Zhao, Xavier Puig, Tete Xiao, Sanja Fidler, Adela Barriuso, and Antonio Torralba.  
441 Semantic understanding of scenes through the ade20k dataset. *IJCV*, 127(3):302–321, 2019.
- 442 [67] Xizhou Zhu, Weijie Su, Lewei Lu, Bin Li, Xiaogang Wang, and Jifeng Dai. Deformable detr: Deformable  
443 transformers for end-to-end object detection. In *ICLR*, 2020.

## 444 Checklist

- 445 1. For all authors...
- 446 (a) Do the main claims made in the abstract and introduction accurately reflect the paper’s contribu-  
447 tions and scope? [Yes]
- 448 (b) Did you describe the limitations of your work? [Yes]
- 449 (c) Did you discuss any potential negative societal impacts of your work? [Yes]
- 450 (d) Have you read the ethics review guidelines and ensured that your paper conforms to them? [Yes]
- 451 2. If you are including theoretical results...
- 452 (a) Did you state the full set of assumptions of all theoretical results? [Yes]
- 453 (b) Did you include complete proofs of all theoretical results? [Yes]
- 454 3. If you ran experiments...
- 455 (a) Did you include the code, data, and instructions needed to reproduce the main experimental  
456 results (either in the supplemental material or as a URL)? [Yes]
- 457 (b) Did you specify all the training details (e.g., data splits, hyperparameters, how they were chosen)?  
458 [Yes]
- 459 (c) Did you report error bars (e.g., with respect to the random seed after running experiments  
460 multiple times)? [N/A]
- 461 (d) Did you include the total amount of compute and the type of resources used (e.g., type of GPUs,  
462 internal cluster, or cloud provider)? [N/A]
- 463 4. If you are using existing assets (e.g., code, data, models) or curating/releasing new assets...
- 464 (a) If your work uses existing assets, did you cite the creators? [Yes]
- 465 (b) Did you mention the license of the assets? [N/A]
- 466 (c) Did you include any new assets either in the supplemental material or as a URL? [No]
- 467 (d) Did you discuss whether and how consent was obtained from people whose data you’re us-  
468 ing/curating? [N/A]
- 469 (e) Did you discuss whether the data you are using/curating contains personally identifiable informa-  
470 tion or offensive content? [N/A]
- 471 5. If you used crowdsourcing or conducted research with human subjects...
- 472 (a) Did you include the full text of instructions given to participants and screenshots, if applicable?  
473 [N/A]
- 474 (b) Did you describe any potential participant risks, with links to Institutional Review Board (IRB)  
475 approvals, if applicable? [N/A]
- 476 (c) Did you include the estimated hourly wage paid to participants and the total amount spent on  
477 participant compensation? [N/A]