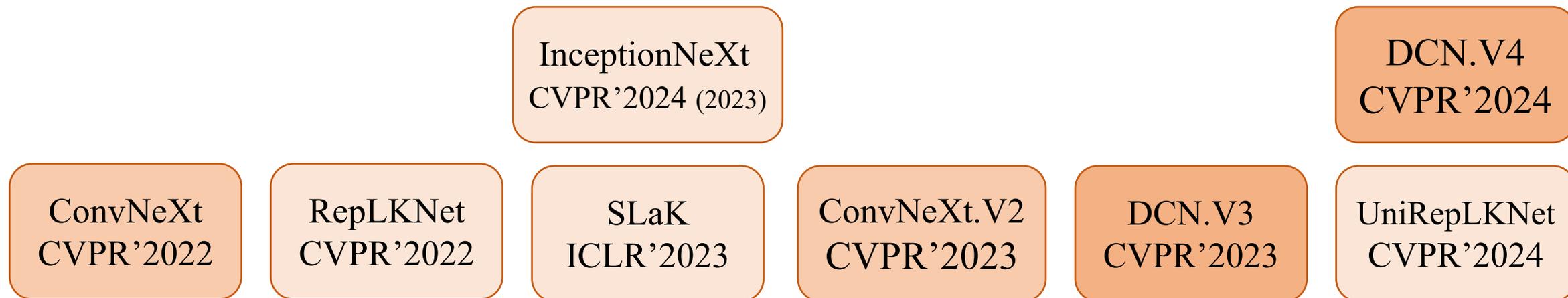


# Modern Convolutional Neural Networks

Siyuan Li

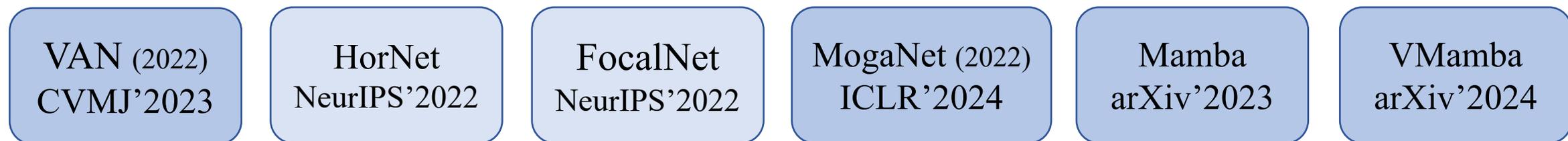
Westlake University, Zhejiang University  
March, 2024

# Timeline of Modern CNNs



## Convolution Kernel Designs

## Large-Kernel Conv + Gated Attentions



# Content

---

## 1. Modern CNNs: Macro Design and Pre-training

MetaFormer, ConvNeXt, ConvNeXt.V2 (Spark, A2MIM)

## 2. Design of Convolution Kernels

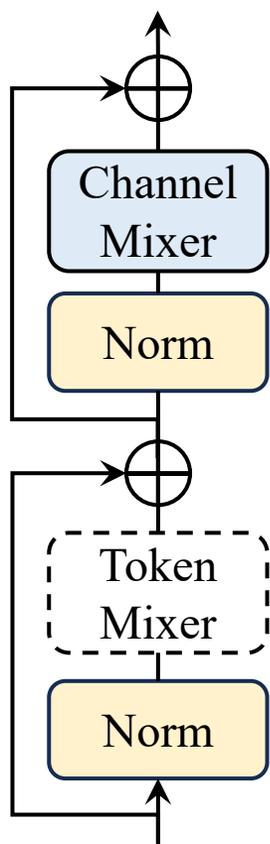
RepLKNet, SLaK, InceptionNext, DCN.V3/V4, UniRepLKNet

## 3. Combining Large Kernel with Gated Attention

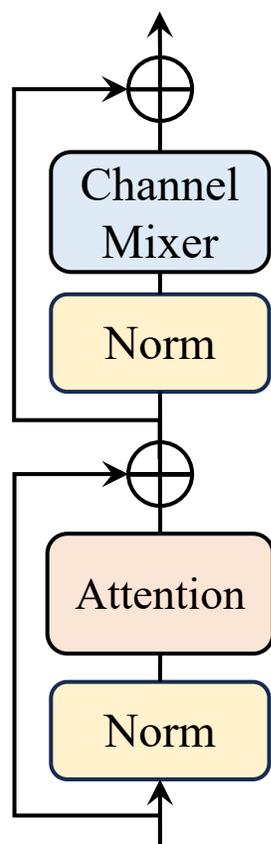
VAN, HorNet, FocalNet, MogaNet, Mamba, VMamba

# Modern CNNs: Macro Design

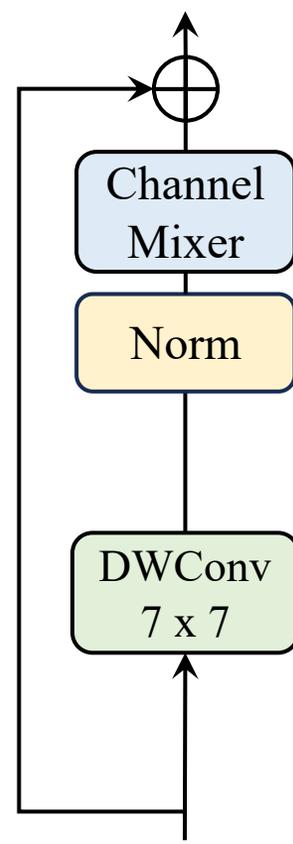
- Macro Design: Patch Embedding + Token Mixer + Channel Mixer + Pre-Norm & Short-cut.



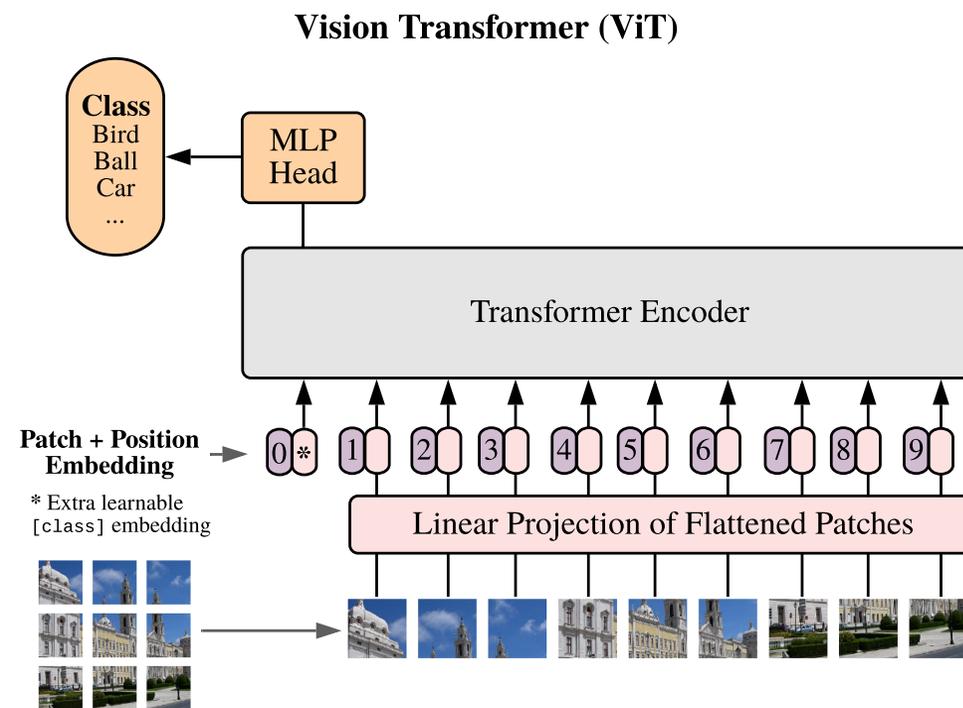
MetaFormer



TransFormer

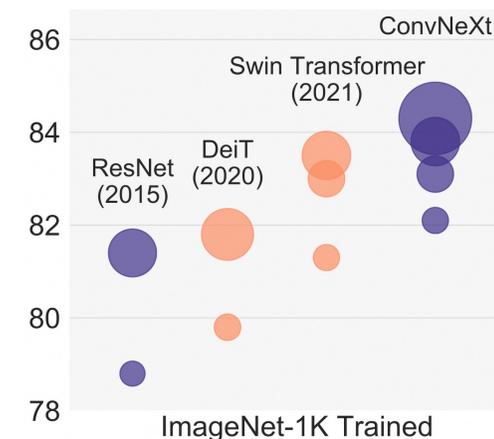
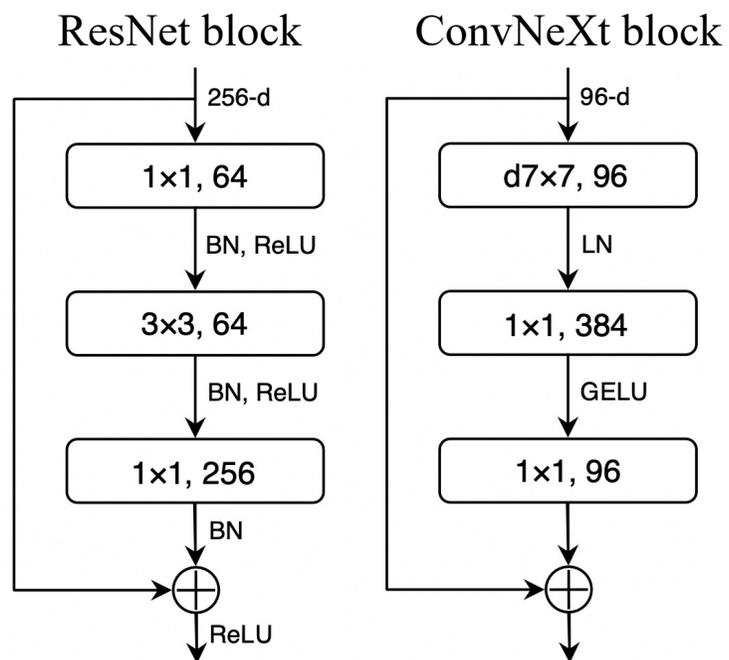
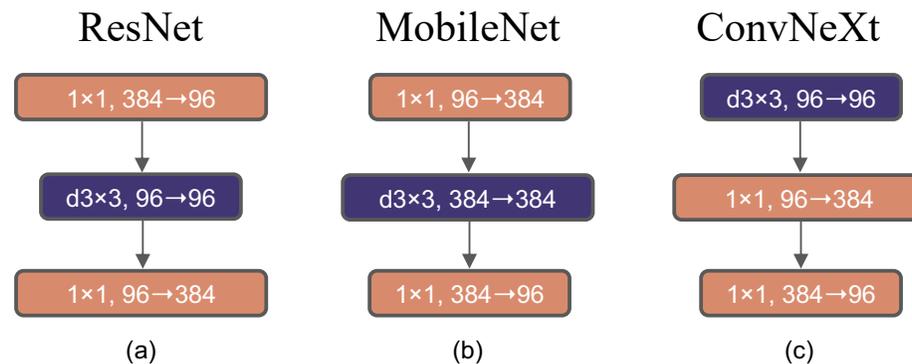
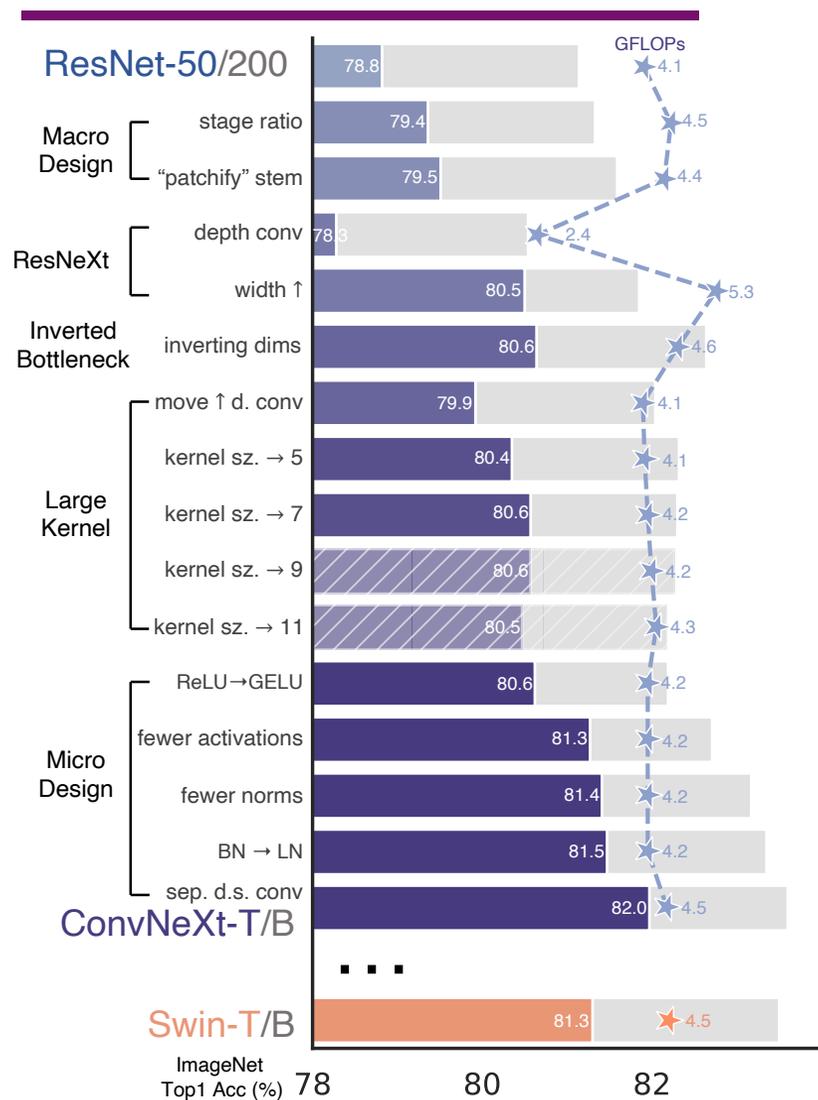


ConvNeXt



- [1] An Image is Worth 16x16 Words: Transformers for Image Recognition at Scale. ICLR, 2021  
 [2] PoolFormer: MetaFormer Is Actually What You Need for Vision. CVPR, 2022.  
 [3] A ConvNet for the 2020s. CVPR, 2022.

# Modern CNNs: ConvNeXt

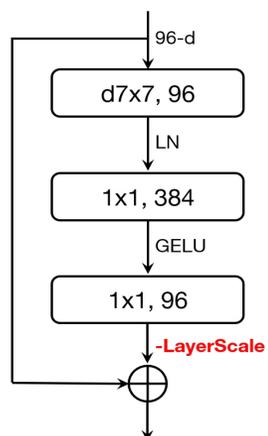


model	image size	#param.	FLOPs	throughput (image / s)	IN-1K top-1 acc.
ImageNet-1K trained models					
● RegNetY-16G [54]	224 <sup>2</sup>	84M	16.0G	334.7	82.9
● EffNet-B7 [71]	600 <sup>2</sup>	66M	37.0G	55.1	84.3
● EffNetV2-L [72]	480 <sup>2</sup>	120M	53.0G	83.7	85.7
○ DeiT-S [73]	224 <sup>2</sup>	22M	4.6G	978.5	79.8
○ DeiT-B [73]	224 <sup>2</sup>	87M	17.6G	302.1	81.8
○ Swin-T	224 <sup>2</sup>	28M	4.5G	757.9	81.3
● ConvNeXt-T	224 <sup>2</sup>	29M	4.5G	774.7	<b>82.1</b>
○ Swin-S	224 <sup>2</sup>	50M	8.7G	436.7	83.0
● ConvNeXt-S	224 <sup>2</sup>	50M	8.7G	447.1	<b>83.1</b>
○ Swin-B	224 <sup>2</sup>	88M	15.4G	286.6	83.5
● ConvNeXt-B	224 <sup>2</sup>	89M	15.4G	292.1	<b>83.8</b>
○ Swin-B	384 <sup>2</sup>	88M	47.1G	85.1	84.5
● ConvNeXt-B	384 <sup>2</sup>	89M	45.0G	95.7	<b>85.1</b>
● ConvNeXt-L	224 <sup>2</sup>	198M	34.4G	146.8	<b>84.3</b>
● ConvNeXt-L	384 <sup>2</sup>	198M	101.0G	50.4	<b>85.5</b>

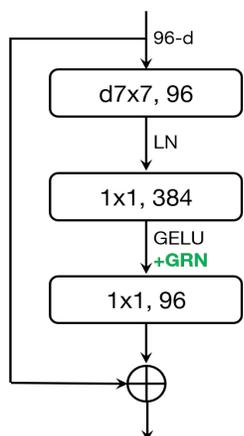
[1] A ConvNet for the 2020s. CVPR, 2022.

# Modern CNNs: ConvNeXt.V2

- CNNs benefit from Masked Image Modeling (MIM) Pre-training.



ConvNeXt.V1



ConvNeXt.V2

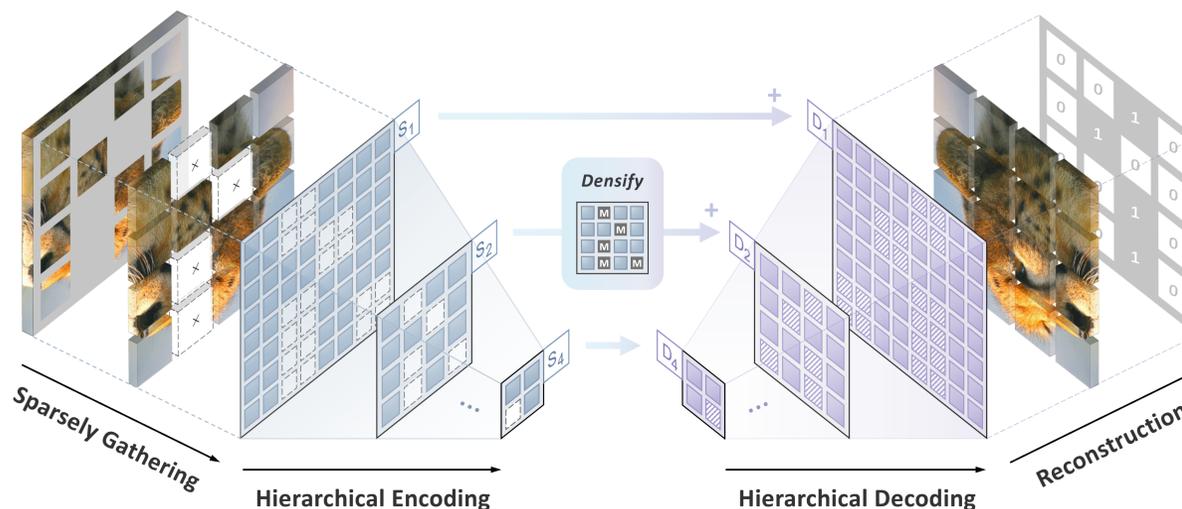
Global Response Normalization (GRN)

```
# gamma, beta: learnable affine transform parameters
# X: input of shape (N,H,W,C)
```

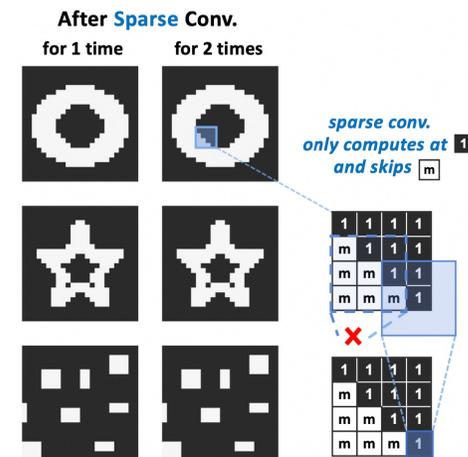
```
gx = torch.norm(X, p=2, dim=(1,2), keepdim=True)
nx = gx / (gx.mean(dim=-1, keepdim=True)+1e-6)
return gamma * (X * nx) + beta + X
```

$$\mathcal{G}(X) := X \in \mathcal{R}^{H \times W \times C} \rightarrow gx \in \mathcal{R}^C$$

$$\mathcal{N}(\|X_i\|) := \|X_i\| \in \mathcal{R} \rightarrow \frac{\|X_i\|}{\sum_{j=1, \dots, C} \|X_j\|} \in \mathcal{R}$$



MIM pre-training with SparK (or FCMAE in ConvNeXt.V2)



Sparse Conv for Masking

Backbone	Method	#param	FLOPs	Val acc.
ConvNeXt V1-B	Supervised	89M	15.4G	83.8
ConvNeXt V1-B	FCMAE	89M	15.4G	83.7
ConvNeXt V2-B	Supervised	89M	15.4G	84.3 (+0.5)
ConvNeXt V2-B	FCMAE	89M	15.4G	<b>84.6 (+0.8)</b>
ConvNeXt V1-L	Supervised	198M	34.4G	84.3
ConvNeXt V1-L	FCMAE	198M	34.4G	84.4
ConvNeXt V2-L	Supervised	198M	34.4G	84.5 (+0.2)
ConvNeXt V2-L	FCMAE	198M	34.4G	<b>85.6 (+1.3)</b>

Methods	#Para.	Sup.	MoCoV3 <sup>‡</sup>	SimMIM <sup>‡</sup>	SparK	A <sup>2</sup> MIM
Target	(M)	Label	CL	RGB	RGB	RGB
ResNet-50	25.6	79.8	80.1	79.9	80.6	<b>80.4</b>
ResNet-101	44.5	81.3	81.6	81.3	82.2	<b>81.9</b>
ResNet-152	60.2	81.8	82.0	81.9	82.7	<b>82.5</b>
ResNet-200	64.7	82.1	82.5	82.2	83.1	<b>83.0</b>
ConvNeXt-T	28.6	82.1	82.3	82.1	82.7	<b>82.5</b>
ConvNeXt-S	50.2	83.1	83.3	83.2	84.1	<b>83.7</b>
ConvNeXt-B	88.6	83.5	83.7	83.6	84.8	<b>84.1</b>

# Content

---

## 1. Modern CNNs: Macro Design and Pre-training

MetaFormer, ConvNeXt, ConvNeXt.V2 (SparK, A2MIM)

## 2. Design of Convolution Kernels

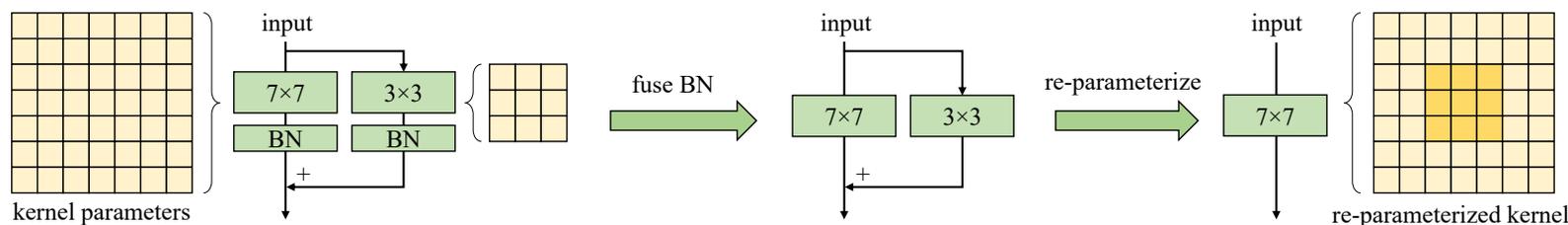
RepLKNet, SLaK, InceptionNext, DCN.V3/V4, UniRepLKNet

## 3. Combining Large Kernel with Gated Attention

VAN, HorNet, FocalNet, MogaNet, Mamba, VMamba

# Large Kernels: RepLKNet

- Large-Kernel (LK) Convolutions are **efficient** and **competitive** as Self-attention.
- Training extremely large convolutions with **Structural Re-parameterization**.



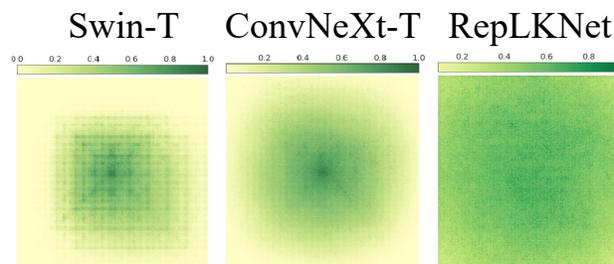
$$DW7 \times 7 = DW3 \times 3 (BN) + DW7 \times 7 (BN) + \text{Short-cut.}$$

Resolution $R$	Impl	Latency (ms) @ Kernel size									
		3	5	7	9	13	17	21	27	29	31
$16 \times 16$	Pytorch	5.6	11.0	14.4	17.6	36.0	57.2	83.4	133.5	150.7	171.4
	Ours	5.6	6.5	6.4	6.9	7.5	8.4	8.4	8.4	8.3	8.4
$32 \times 32$	Pytorch	21.9	34.1	54.8	76.1	141.2	230.5	342.3	557.8	638.6	734.8
	Ours	21.9	28.7	34.6	40.6	52.5	64.5	73.9	87.9	92.7	96.7
$64 \times 64$	Pytorch	69.6	141.2	228.6	319.8	600.0	977.7	1454.4	2371.1	2698.4	3090.4
	Ours	69.6	112.6	130.7	152.6	199.7	251.5	301.0	378.2	406.0	431.7

Kernel size	Architecture	ImageNet			ADE20K		
		Top-1	Params	FLOPs	mIoU	Params	FLOPs
7-7-7-7	ConvNeXt-Tiny	81.0	29M	4.5G	44.6	60M	939G
7-7-7-7	ConvNeXt-Small	82.1	50M	8.7G	45.9	82M	1027G
7-7-7-7	ConvNeXt-Base	82.8	89M	15.4G	47.2	122M	1170G
31-29-27-13	ConvNeXt-Tiny	81.6	32M	6.1G	<b>46.2</b>	64M	973G
31-29-27-13	ConvNeXt-Small	82.5	58M	11.3G	<b>48.2</b>	90M	1081G

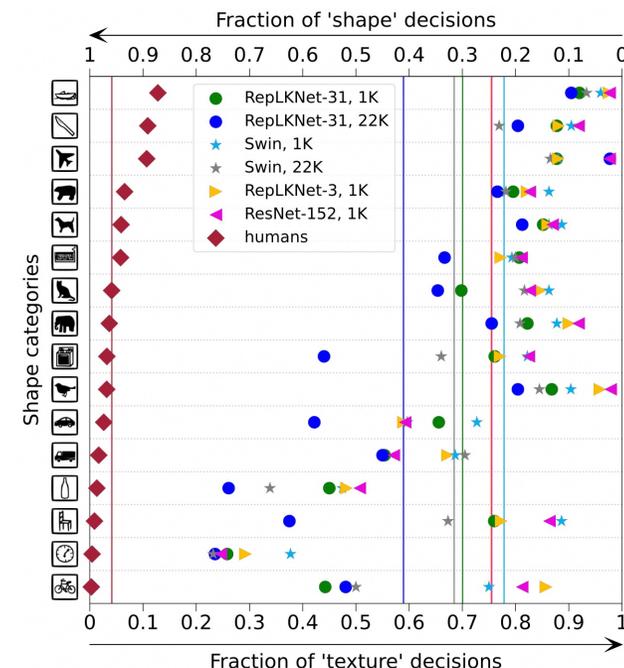
**Extremely large kernels** benefit both classification and downstream tasks and outperforms ViTs.

Large kernels are **memory bound** instead of compute bound.



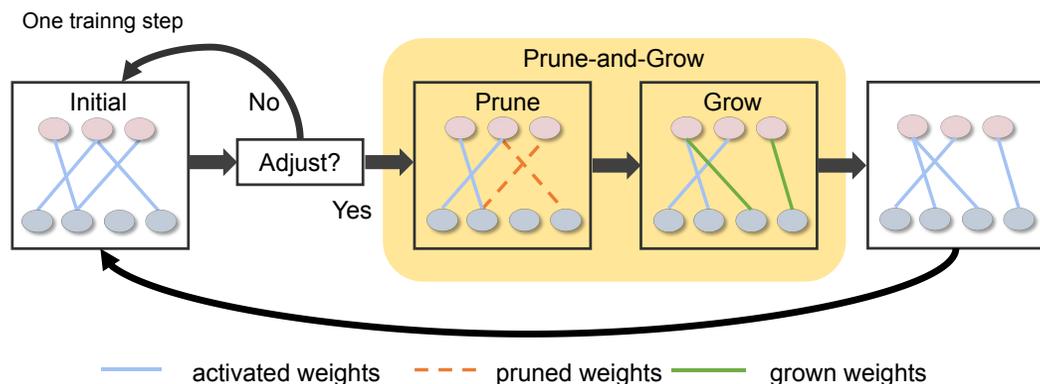
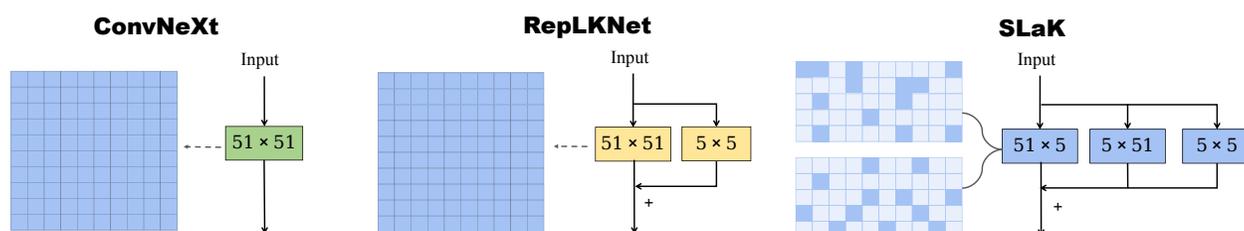
Effective receptive field

Large kernels are **shape biased** as ViTs.



# Large Kernels: SLaK

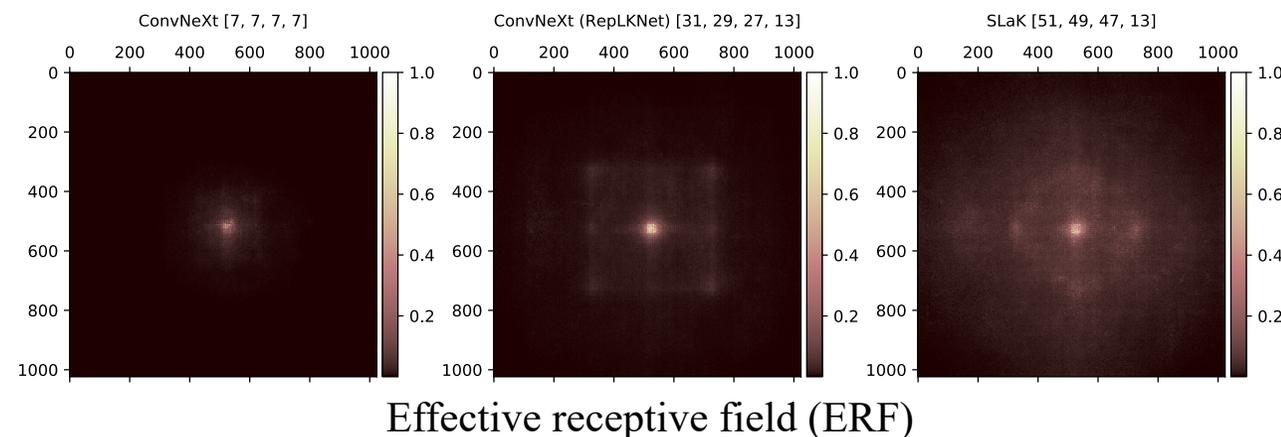
- Step 1: Decomposing a large kernel (61x61) into two rectangular, parallel kernels.
- Step 2: Using sparse groups training (speedup), expanding more width.



- (1) Initialization: Constructing Sparse Convolution based on SNIP<sup>[2]</sup>
- (2) Dynamic sparsity: Pruning (the lowest magnitude) and growing

Kernel Size	Top-1 Acc	#Params	FLOPs	Decomposed			Sparse groups			Sparse groups, expand more width		
				Top-1 Acc	#Params	FLOPs	Top-1 Acc	#Params	FLOPs	Top-1 Acc	#Params	FLOPs
7-7-7-7	81.0	29M	4.5G	80.0	17M	2.6G	81.1	29M	4.5G			
31-29-37-13	81.3	30M	5.0G	80.4	18M	2.9G	81.5	30M	4.8G			
51-49-47-13	81.5	31M	5.4G	80.5	18M	3.1G	81.6	30M	5.0G			
61-59-57-13	81.4	31M	5.6G	80.4	19M	3.2G	81.5	31M	5.2G			

Model	Kernel Size	AP <sup>box</sup>	AP <sub>50</sub> <sup>box</sup>	AP <sub>75</sub> <sup>box</sup>	AP <sup>mask</sup>	AP <sub>50</sub> <sup>mask</sup>	AP <sub>75</sub> <sup>mask</sup>
pre-trained for 120 epochs, finetuned for 1 × (12 epochs)							
ConvNeXt-T (Liu et al., 2022b)	7-7-7-7	47.3	65.9	51.5	41.1	63.2	44.4
ConvNeXt-T (RepLkNET)* (Ding et al., 2022)	31-29-27-13	47.8	66.7	52.0	41.4	63.9	44.7
SLaK-T	51-49-47-13	<b>48.4</b>	<b>67.2</b>	<b>52.5</b>	<b>41.8</b>	<b>64.4</b>	<b>45.2</b>
pre-trained for 300 epochs, finetuned for 3 × (36 epochs)							
ConvNeXt-T (Liu et al., 2022b)	7-7-7-7	50.4	69.1	54.8	43.7	66.5	47.3
SLaK-T	51-49-47-13	<b>51.3</b>	<b>70.0</b>	<b>55.7</b>	<b>44.3</b>	<b>67.2</b>	<b>48.1</b>

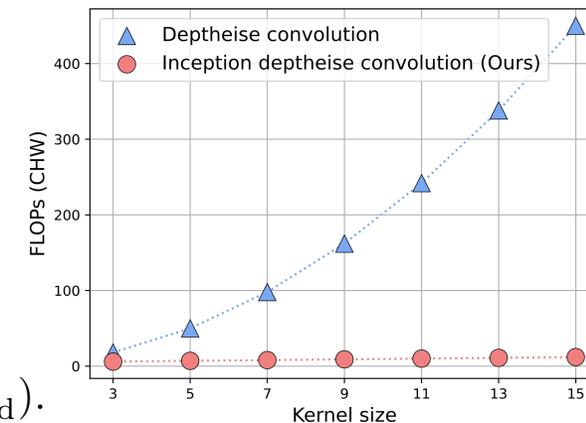
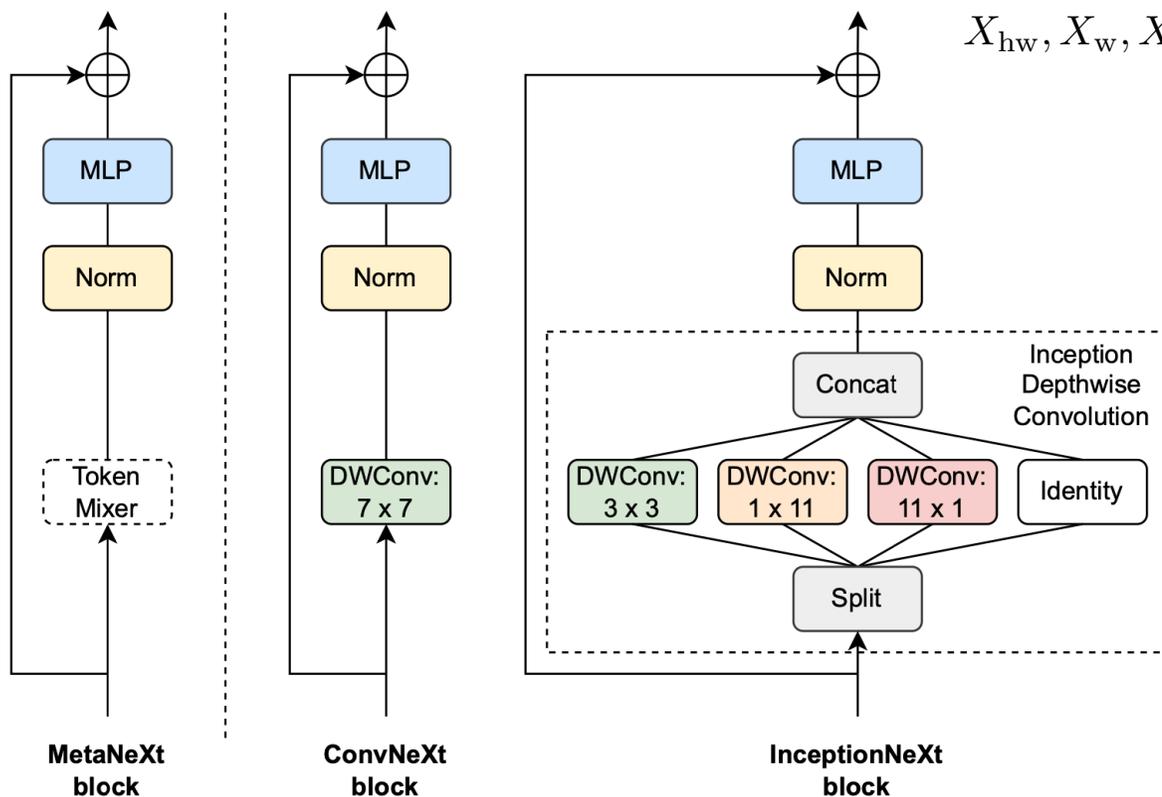


[1] More ConvNets in the 2020s: Scaling up Kernels Beyond 51x51 using Sparsity. ICLR, 2023.

[2] SNIP: Single-shot Network Pruning based on Connection Sensitivity. ICLR, 2019.

# Large Kernels: InceptionNeXt

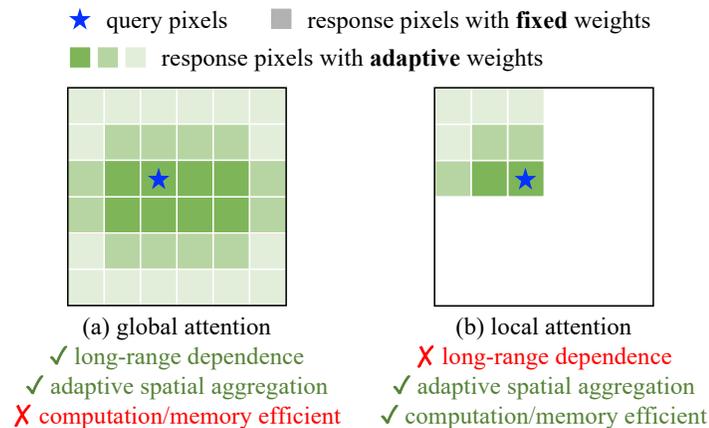
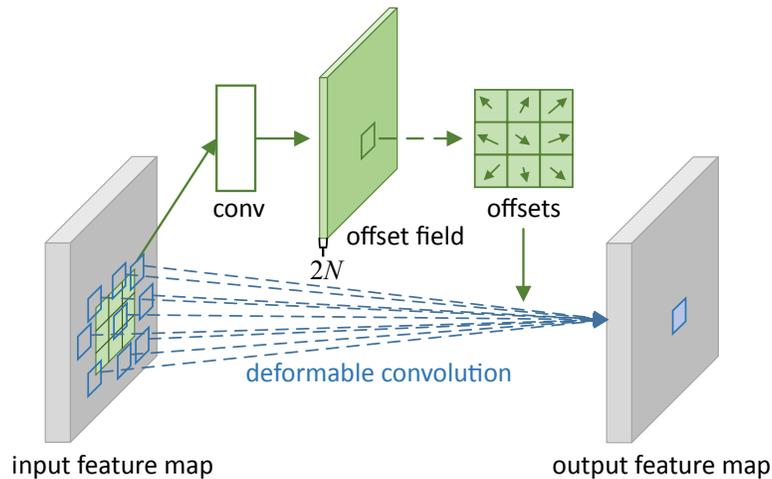
- MetaNeXt: Fusing Token Mixer with Channel Mixer + PreNorm + ShortCut.
- Inception Kernels: Better performance and throughputs than Depth-wise Conv 7x7.



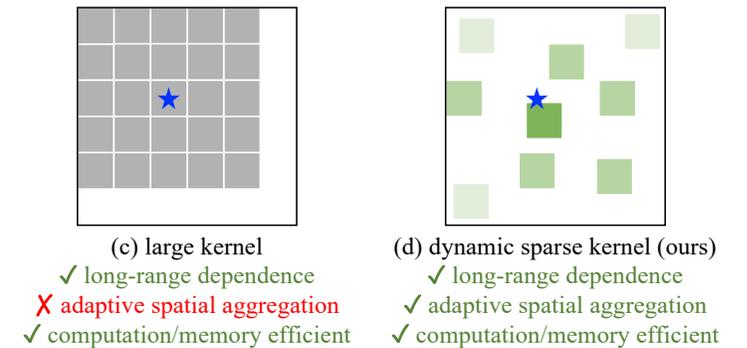
Model	Mixing Type	Image (size)	Params (M)	MACs (G)	Throughput (img/second)		Top-1 (%)
					Train	Inference	
DeiT-S [61]	Attn	224 <sup>2</sup>	22	4.6	1227	3781	79.8
T2T-ViT-14 [76]	Attn	224 <sup>2</sup>	22	4.8	–	–	81.5
TNT-S [18]	Attn	224 <sup>2</sup>	24	5.2	–	–	81.5
Swin-T [37]	Attn	224 <sup>2</sup>	29	4.5	564	1768	81.3
Focal-T [73]	Attn	224 <sup>2</sup>	29	4.9	–	–	82.2
ResNet-50 [20, 69]	Conv	224 <sup>2</sup>	26	4.1	969	3149	78.4
RSB-ResNet-50 [20, 69]	Conv	224 <sup>2</sup>	26	4.1	969	3149	79.8
RegNetY-4G [46, 69]	Conv	224 <sup>2</sup>	21	4.0	670	2694	81.3
FocalNet-T [72]	Conv	224 <sup>2</sup>	29	4.5	–	–	82.3
ConvNeXt-T [38]	Conv	224 <sup>2</sup>	29	4.5	575	2413 (1943)	82.1
InceptionNeXt-T (Ours)	Conv	224 <sup>2</sup>	28	4.2	901 (+57%)	2900 (+20%)	82.3 (+0.2)

# Kernel Designs: DCN.V3 (InternImage)

- DCN.V3: Learnable offsets (V1) + Softmax-normalized modulation (V2) + Grouping.



## Self-Attention vs. Conv vs. DCN



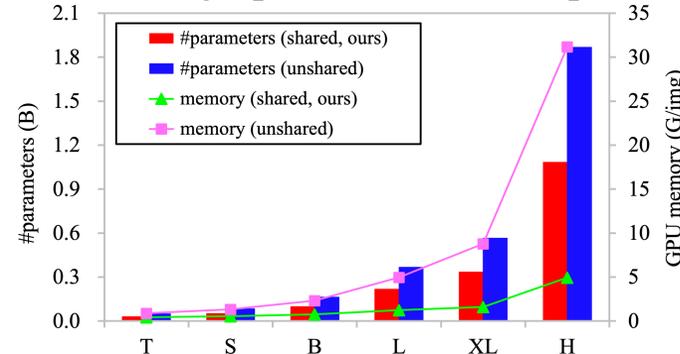
$$\text{DCN.V1: } \mathbf{y}(p_0) = \sum_{\mathbf{p}_n \in \mathcal{R}} \mathbf{w}(\mathbf{p}_n) \cdot \mathbf{x}(p_0 + \mathbf{p}_n + \Delta \mathbf{p}_n)$$

$$\text{DCN.V2: } \mathbf{y}(p_0) = \sum_{k=1}^K \mathbf{w}_k \mathbf{m}_k \mathbf{x}(p_0 + p_k + \Delta p_k)$$

$$\text{DCN.V3: } \mathbf{y}(p_0) = \sum_{g=1}^G \sum_{k=1}^K \mathbf{w}_g \mathbf{m}_{gk} \mathbf{x}_g(p_0 + p_k + \Delta p_{gk})$$

Offsets  $\Delta p_n$ , Regular grids  $p_n$ , Modulation  $m_k$ , weights  $w$

## Scaling-up with efficient impl.



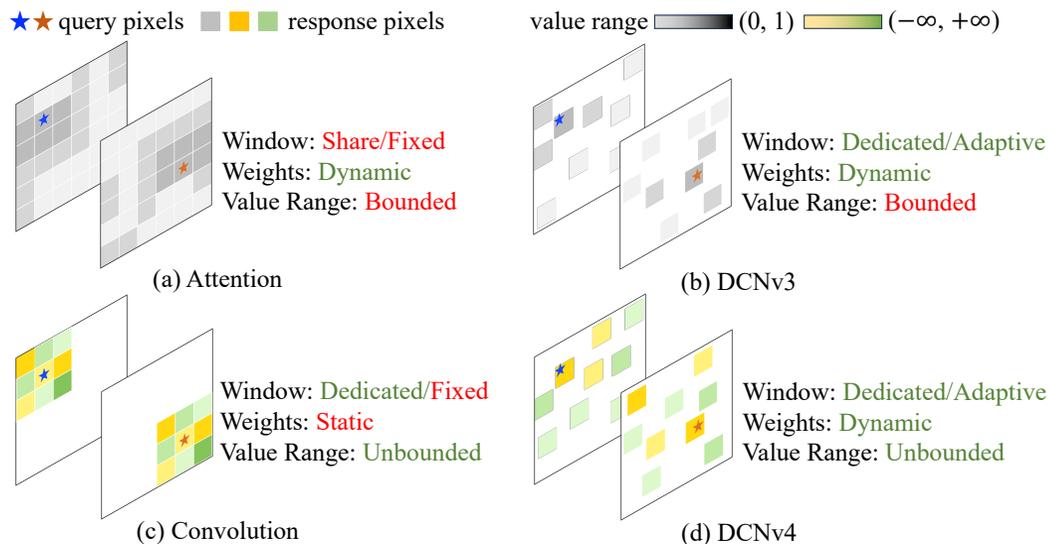
method	type	scale	#params	#FLOPs	acc (%)
SwinV2-L/24 <sup>‡</sup> [16]	T	384 <sup>2</sup>	197M	115G	87.6
RepLKNet-31L <sup>‡</sup> [22]	C	384 <sup>2</sup>	172M	96G	86.6
HorNet-L <sup>‡</sup> [43]	C	384 <sup>2</sup>	202M	102G	87.7
ConvNeXt-L <sup>‡</sup> [21]	C	384 <sup>2</sup>	198M	101G	87.5
ConvNeXt-XL <sup>‡</sup> [21]	C	384 <sup>2</sup>	350M	179G	87.8
InternImage-L <sup>‡</sup> (ours)	C	384 <sup>2</sup>	223M	108G	87.7
InternImage-XL <sup>‡</sup> (ours)	C	384 <sup>2</sup>	335M	163G	88.0
ViT-G/14 <sup>#</sup> [30]	T	518 <sup>2</sup>	1.84B	5160G	90.5
CoAtNet-6 <sup>#</sup> [20]	T	512 <sup>2</sup>	1.47B	1521G	90.5
CoAtNet-7 <sup>#</sup> [20]	T	512 <sup>2</sup>	2.44B	2586G	90.9
Florence-CoSwin-H <sup>#</sup> [59]	T	—	893M	—	90.0
SwinV2-G <sup>#</sup> [16]	T	640 <sup>2</sup>	3.00B	—	90.2
RepLKNet-XL <sup>#</sup> [22]	C	384 <sup>2</sup>	335M	129G	87.8
BiT-L-ResNet152x4 <sup>#</sup> [67]	C	480 <sup>2</sup>	928M	—	87.5
InternImage-H <sup>#</sup> (ours)	C	224 <sup>2</sup>	1.08B	188G	88.9
InternImage-H <sup>#</sup> (ours)	C	640 <sup>2</sup>	1.08B	1478G	89.6

[1] Deformable Convolutional Networks. ICCV, 2017. [2] Deformable ConvNets v2: More Deformable, Better Results. CVPR, 2018.

[3] InternImage: Exploring Large-Scale Vision Foundation Models with Deformable Convolutions. CVPR, 2023.

# Kernel Designs: DCN.V4 (FlashInternImage)

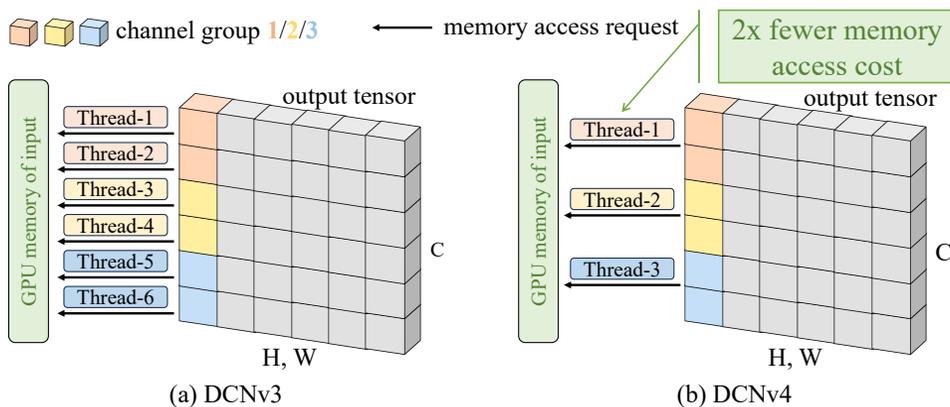
- DCN.V4: No Softmax normalization + Speed-up (reducing HRM as Flash-Attention).



Model	5th EP	10th Ep	20th Ep	50th Ep	300th Ep
ConvNeXt	29.9	53.5	66.1	74.8	83.8
ConvNeXt + softmax	<b>(-21.4)</b>	<b>(-28.2)</b>	<b>(-15.0)</b>	<b>(-5.7)</b>	<b>(-2.2)</b>

Using Softmax in DWConv7×7 degenerating performance

Operator	Runtime (ms)				
	56 × 56 × 128	28 × 28 × 256	14 × 14 × 512	7 × 7 × 1024	14 × 14 × 768
Attention (torch)	30.8 / 19.3	3.35 / 2.12	0.539 / 0.448	0.446 / 0.121	0.779 / 0.654
FlashAttention-2	N/A / 2.46	N/A / 0.451	N/A / <b>0.123</b>	N/A / 0.0901	N/A / 0.163
Window Attn (7 × 7)	4.05 / 1.46	2.07 / 0.770	1.08 / 0.422	0.577 / 0.239	1.58 / 0.604
DWConv (7 × 7, torch)	2.02 / 1.98	1.03 / 1.00	0.515 / 0.523	0.269 / 0.261	0.779 / 0.773
DWConv (7 × 7, cuDNN)	0.981 / 0.438	0.522 / 0.267	0.287 / 0.153	0.199 / 0.102	0.413 / 0.210
DCNv3	1.45 / 1.52	0.688 / 0.711	0.294 / 0.298	0.125 / 0.126	0.528 / 0.548
DCNv4	<b>0.606 / 0.404</b>	<b>0.303 / 0.230</b>	<b>0.145 / 0.123</b>	<b>0.0730 / 0.0680</b>	<b>0.224 / 0.147</b>



## ImageNet-1K Classification

Model	Size	Scale	Acc	Throughput
Swin-T	29M	224 <sup>2</sup>	81.3	1989 / 3619
ConvNeXt-T	29M	224 <sup>2</sup>	82.1	2485 / 4305
InternImage-T	30M	224 <sup>2</sup>	83.5	1409 / 1746
FlashInternImage-T	30M	224 <sup>2</sup>	<b>83.6</b>	2316 / 3154 (+64% / +80%)
Swin-S	50M	224 <sup>2</sup>	83.0	1167/2000
ConvNeXt-S	50M	224 <sup>2</sup>	83.1	1645/2538
InternImage-S	50M	224 <sup>2</sup>	84.2	1044/1321
FlashInternImage-S	50M	224 <sup>2</sup>	<b>84.4</b>	1625 / 2396

## COCO2017 Det. and Seg.

Model	#param	FPS	Cascade Mask R-CNN			
			1×		3×+MS	
			AP <sup>b</sup>	AP <sup>m</sup>	AP <sup>b</sup>	AP <sup>m</sup>
Swin-L	253M	20 / 26	51.8	44.9	53.9	46.7
ConvNeXt-L	255M	26 / 40	53.5	46.4	54.8	47.6
InternImage-L	277M	20 / 26	54.9	47.7	56.1	48.5
ConvNeXt-XL	407M	21 / 32	53.6	46.5	55.2	47.7
InternImage-XL	387M	16 / 23	55.3	48.1	56.2	48.8
FlashInternImage-L	277M	26 / 39	<b>55.6</b>	<b>48.2</b>	<b>56.7</b>	<b>48.9</b>

# Content

---

## 1. Modern CNNs: Macro Design and Pre-training

MetaFormer, ConvNeXt, ConvNeXt.V2 (Spark, A2MIM)

## 2. Design of Convolution Kernels

RepLKNet, SLaK, InceptionNext, DCN.V3/V4, UniRepLKNet

## 3. Combining Large Kernel with Gated Attention

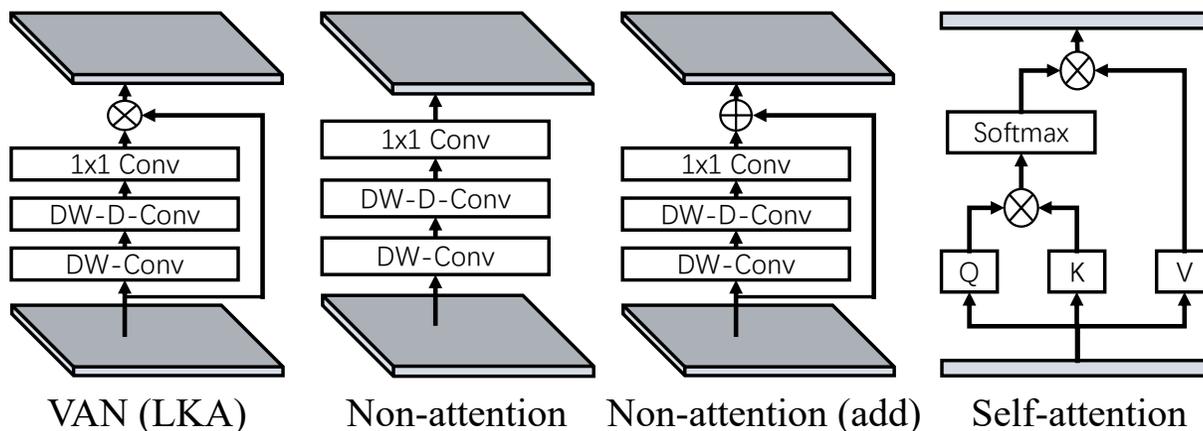
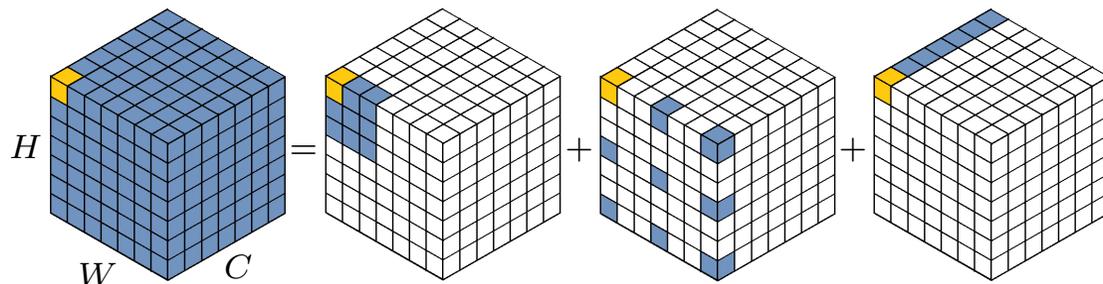
VAN, HorNet, FocalNet, MogaNet, Mamba, VMamba

# Gating & Large-kernel: VAN

- Decomposed large kernel + Gating.

$$\text{Conv}9 \times 9 = \text{DWConv}3 \times 3 + \text{DWConv}3 \times 3 + \text{PWConv}1 \times 1$$

(Dilation=3)



Properties	Convolution	Self-Attention	LKA
Local Receptive Field	✓	✗	✓
Long-range Dependence	✗	✓	✓
Spatial Adaptability	✗	✓	✓
Channel Adaptability	✗	✗	✓
Computational complexity	$\mathcal{O}(n)$	$\mathcal{O}(n^2)$	$\mathcal{O}(n)$

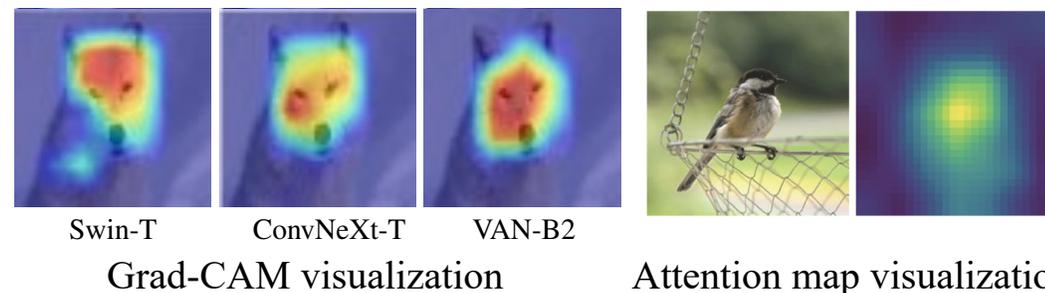
Properties of DWConv vs. MHSA vs. Large-kernel Attention

Method	$K$	Dilation	Params. (M)	GFLOPs	Acc(%)
VAN-B0	7	2	4.03	0.85	74.8
VAN-B0	14	3	4.07	0.87	75.3
VAN-B0	21	3	4.11	0.88	75.4
VAN-B0	28	4	4.14	0.90	75.4

Kernel size vs. Dilation vs. ImageNet Acc (%)

$$\text{Conv}21 \times 21 = \text{DWConv}5 \times 5 + \text{DWConv}7 \times 7 + \text{PWConv}1 \times 1$$

(Dilation=3)



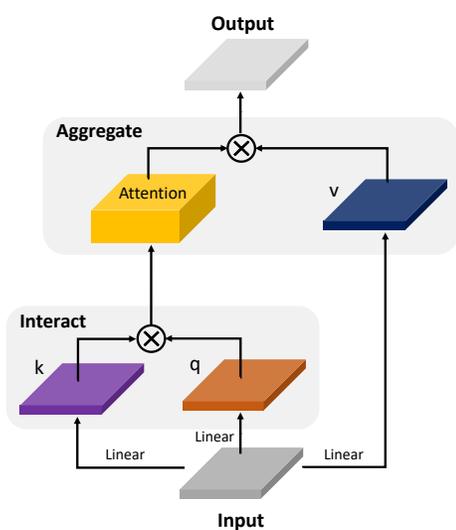
# Gating & Hierarchical Kernel: FocalNet

- Hierarchical Contextualization + Gated Aggregation.

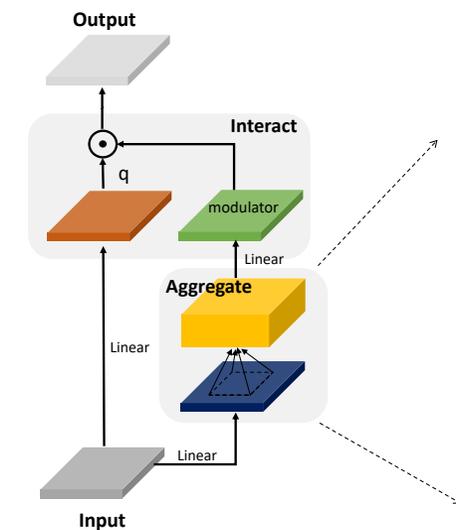


```

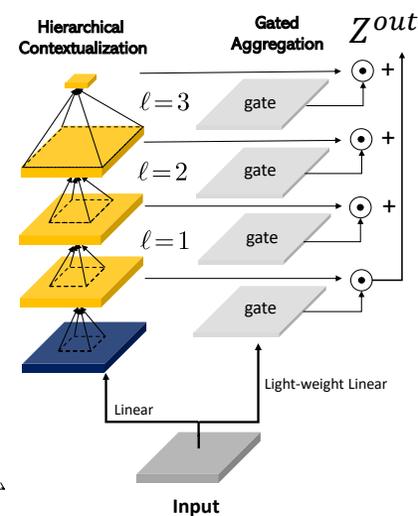
5 def forward(x, m=0):
6     x = pj_in(x).permute(0, 3, 1, 2)
7     q, z, gate = split(x, (C, C, L+1), 1)
8     for l in range(L):
9         z = hc_layers[l](z) # Eq.(4), hierarchical contextualization
10        m = m + z * gate[:, l:l+1] # Eq.(5), gated aggregation
11    m = m + GeLU(z.mean(dim=(2,3))) * gate[:, L:]
12    x = q * pj_cxt(m) # Eq.(6), Focal Modulation
13    return pj_out(x.permute(0, 2, 3, 1))
    
```



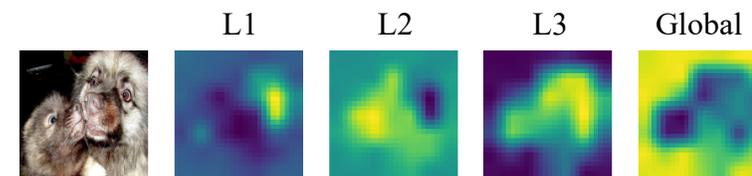
(a) Self-Attention



(b) Focal-Modulation

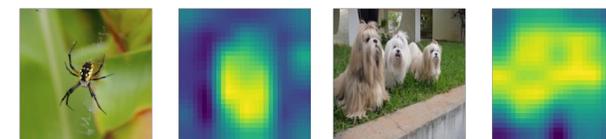


(c) Context Aggregation



$$\mathbf{Z}^l = f_a^l(\mathbf{Z}^{l-1}) \triangleq \text{GeLU}(\text{DWConv}(\mathbf{Z}^{l-1})) \in \mathbb{R}^{H \times W \times C} \quad \text{Eq. (4)}$$

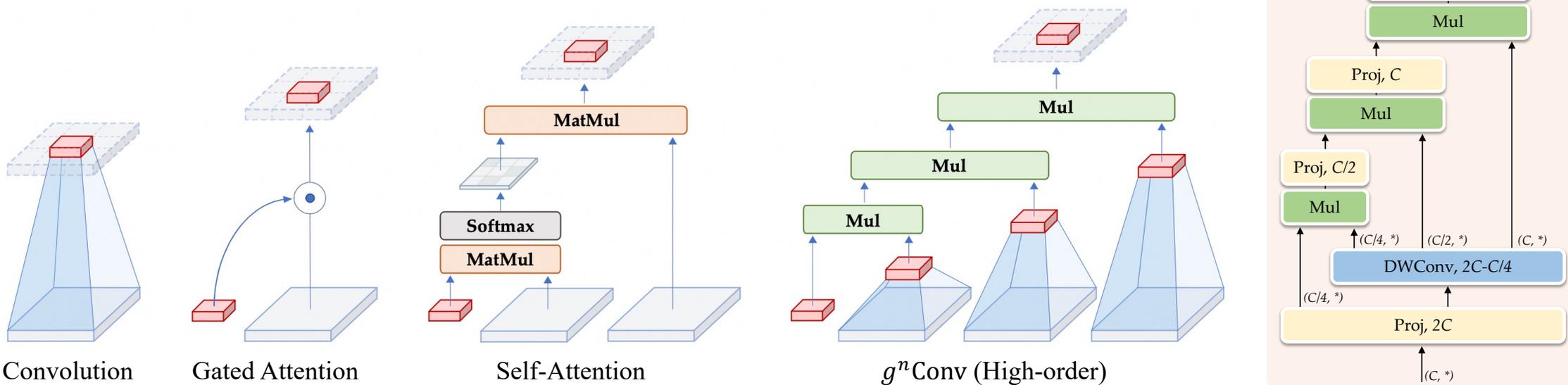
$$\mathbf{Z}^{\text{out}} = \sum_{\ell=1}^{L+1} \mathbf{G}^{\ell} \odot \mathbf{Z}^{\ell} \in \mathbb{R}^{H \times W \times C} \quad \text{Eq. (5)}$$



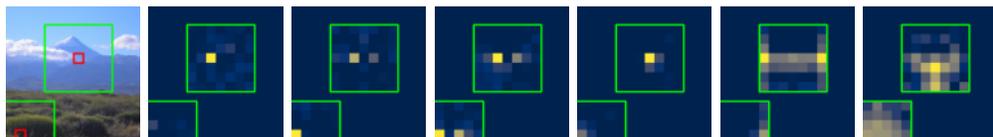
$$\mathbf{y}_i = q(\mathbf{x}_i) \odot h\left(\sum_{\ell=1}^{L+1} \mathbf{g}_i^{\ell} \cdot \mathbf{z}_i^{\ell}\right) \quad \text{Eq. (6)}$$

# Gating & Hierarchical Kernel: HorNet

- High-order Interactions: Recursive DWConv + Gating.



$$x_{g^n \text{Conv}}^{(i,c)} = p_n^{(i,c)} = \sum_{j \in \Omega_i} \sum_{c'=1}^C \frac{w_{n-1,i \rightarrow j}^c \mathbf{g}_{n-1}^{(i,c)} w_{\phi_{in}}^{(c',c)}}{w_{\phi_{in}}^{(c',c)}} x^{(j,c')} \triangleq \sum_{j \in \Omega_i} \sum_{c'=1}^C \frac{h_{ij}^c w_{\phi_{in}}^{(c',c)}}{w_{\phi_{in}}^{(c',c)}} x^{(j,c')} \quad \text{Eq. (3.8)}$$



Adaptive weights generated by  $g^n$  Conv, i.e.,  $\frac{1}{C} \sum_{c=1}^C h_{ij}^c$  in Eq. (3.8)

```
def forward(self, x):
    x = self.proj_in(x)
    y, x = torch.split(x, (self.dims[0], sum(self.dims)), dim=1)
    x = self.dwconv(x)
    x_list = torch.split(x, self.dims, dim=1)
    x = y * x_list[0]
    for i in range(self.order - 1):
        x = self.projs[i](x) * x_list[i+1]
    return self.proj_out(x)
```

```
self.projs = nn.ModuleList(
    [nn.Conv2d(self.dims[i], self.dims[i+1], 1)
     for i in range(order-1)])
self.proj_out = nn.Conv2d(dim, dim, 1)
```

# Multi-order Interaction: MogaNet

- Representation Bottleneck<sup>[1]</sup>: Loss in the middle-order interactions.

Multi-order Interactions

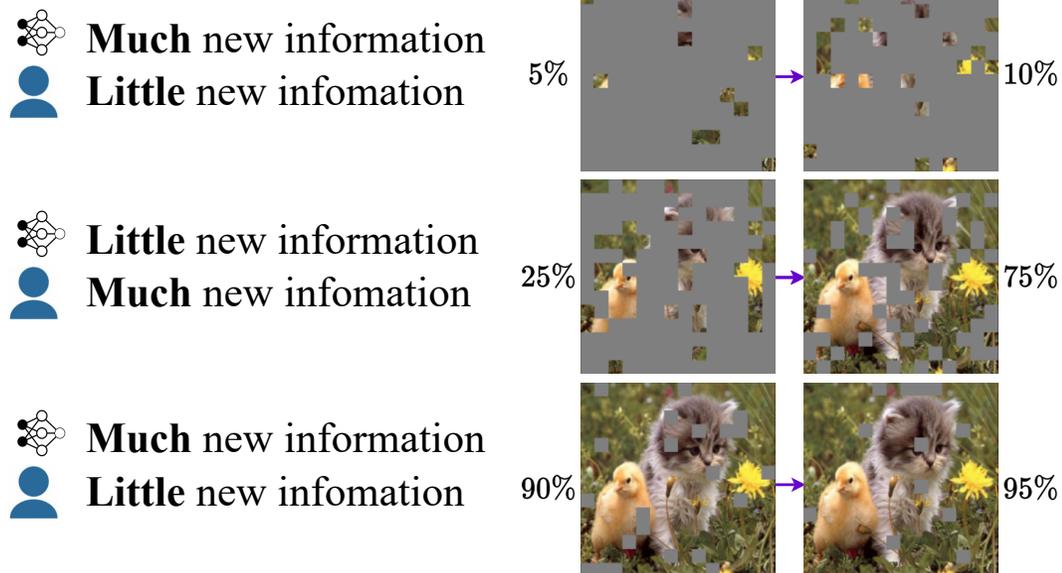
$$I^{(m)}(i, j) = \mathbb{E}_{S \subseteq N \setminus \{i, j\}, |S|=m} [\Delta f(i, j, S)]$$

$$N = \{1, \dots, n\} \quad 0 \leq m \leq n - 2$$

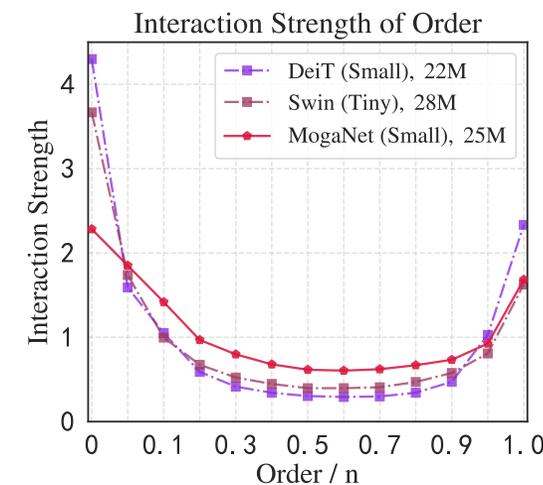
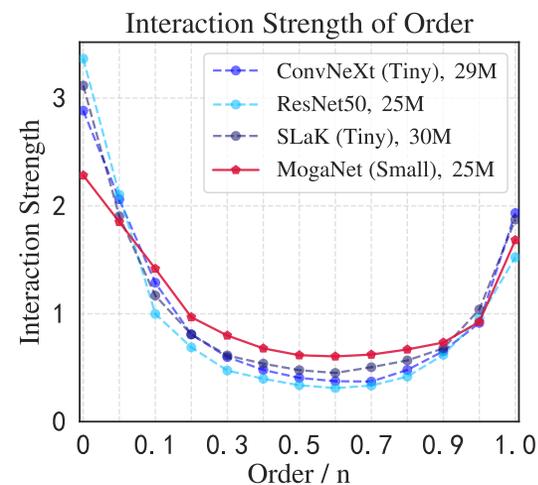
$$\Delta f(i, j, S) = f(S \cup \{i, j\}) - f(S \cup \{i\}) - f(S \cup \{j\}) + f(S)$$

Interaction Strengths

$$J^{(m)} = \frac{\mathbb{E}_{x \in \Omega} \mathbb{E}_{i, j} |I^{(m)}(i, j|x)|}{\mathbb{E}_{m'} \mathbb{E}_{x \in \Omega} \mathbb{E}_{i, j} |I^{(m')}(i, j|x)|}$$

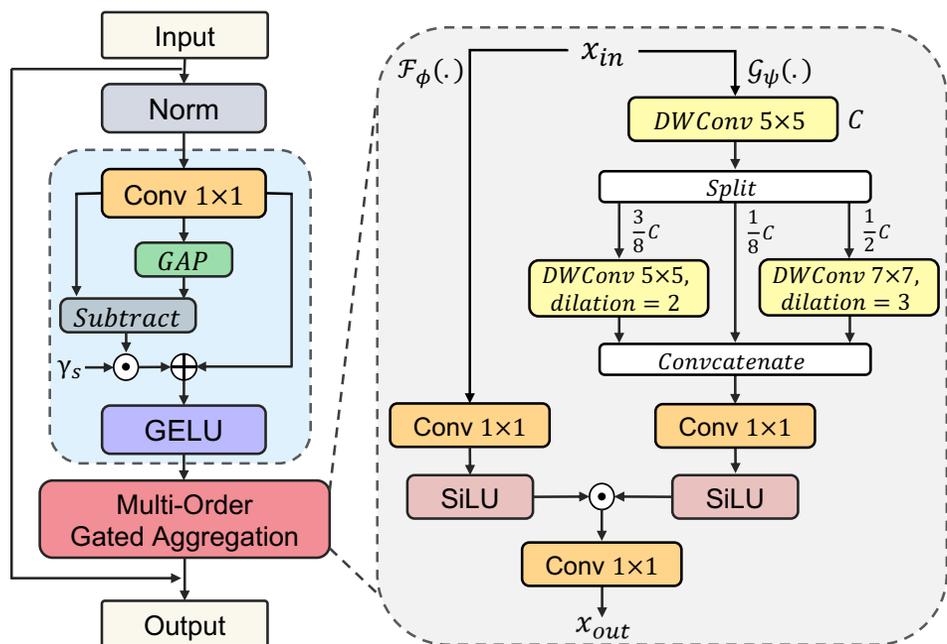


Both ViTs and modern CNN architectures fail to explore middle-order interactions, which are informative to humans.



# Multi-order Interaction: MogaNet

- Spatial Aggregation (SA): Multi-order context extraction + Gated aggregation.



$$Z = X + \text{Moga}\left(\text{FD}(\text{Norm}(X))\right)$$

Feature decomposition:  $Y = \text{Conv}_{1 \times 1}(X),$   
 $Z = \text{GELU}\left(Y + \gamma_s \odot (Y - \text{GAP}(Y))\right)$

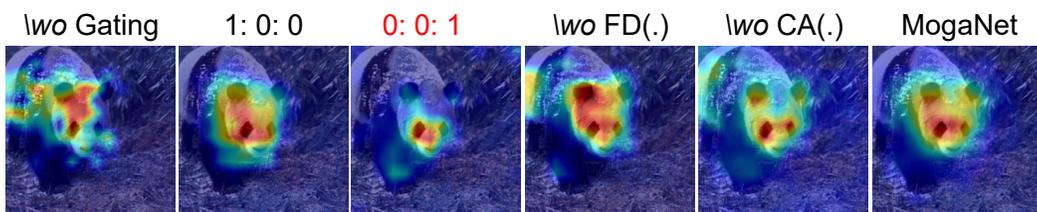
Gated aggregation branch:  $Z = \underbrace{\text{SiLU}(\text{Conv}_{1 \times 1}(X))}_{\mathcal{F}_\phi} \odot \underbrace{\text{SiLU}(\text{Conv}_{1 \times 1}(Y_C))}_{\mathcal{G}_\psi}$

Multi-order DWConvs: DW5×5, DW5×5 (d=2), DW7×7 (d=3)

$$C_l + C_m + C_h = C, Y_C = \text{Concat}(Y_{l,1:C_l}, Y_m, Y_h)$$

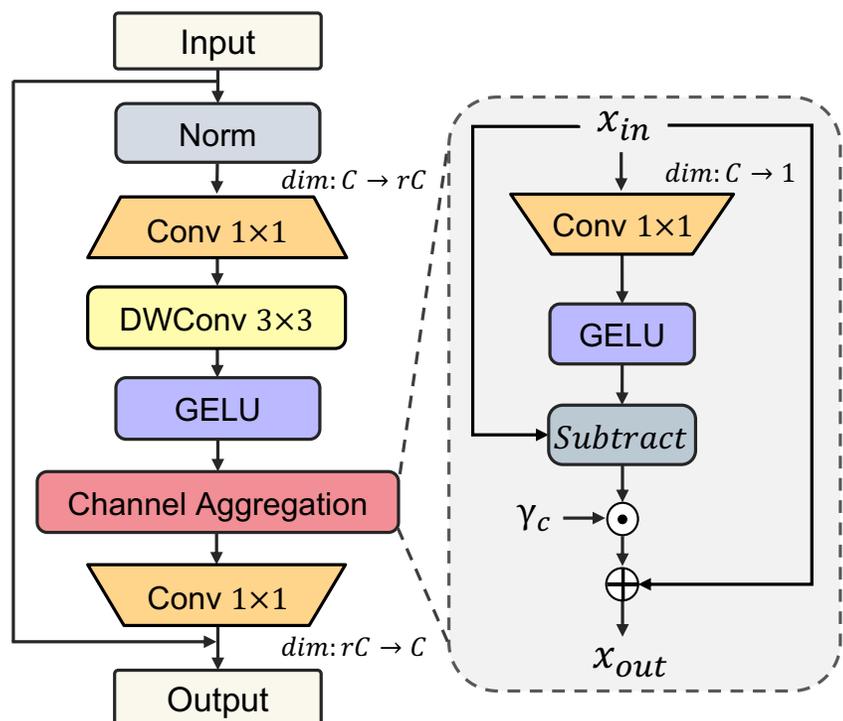
Modules	Top-1	Params.	FLOPs	Context branch				
	Acc (%)			(M)	(G)	None	GELU	SiLU
Baseline (+Gating branch)	77.2	5.09	1.070	None	76.3	76.7	76.7	
DW <sub>7×7</sub>	77.4	5.14	1.094	Gating branch	Sigmoid	76.8	77.0	76.9
DW <sub>5×5,d=1</sub> + DW <sub>7×7,d=3</sub>	77.5	5.15	1.112		GELU	76.7	76.8	77.0
DW <sub>5×5,d=1</sub> + DW <sub>5×5,d=2</sub> + DW <sub>7×7,d=3</sub>	77.5	5.17	1.185		SiLU	76.9	77.1	<b>77.2</b>
+Multi-order, $C_l : C_m : C_h = 1 : 0 : 3$	77.5	5.17	1.099					
+Multi-order, $C_l : C_m : C_h = 0 : 1 : 1$	77.6	5.17	1.103					
+Multi-order, $C_l : C_m : C_h = 1 : 6 : 9$	77.7	5.17	1.104					
+Multi-order, $C_l : C_m : C_h = 1 : 3 : 4$	<b>77.8</b>	5.17	1.102					

Ablation of SA module with MogaNet-T on ImageNet



# Multi-order Interaction: MogaNet

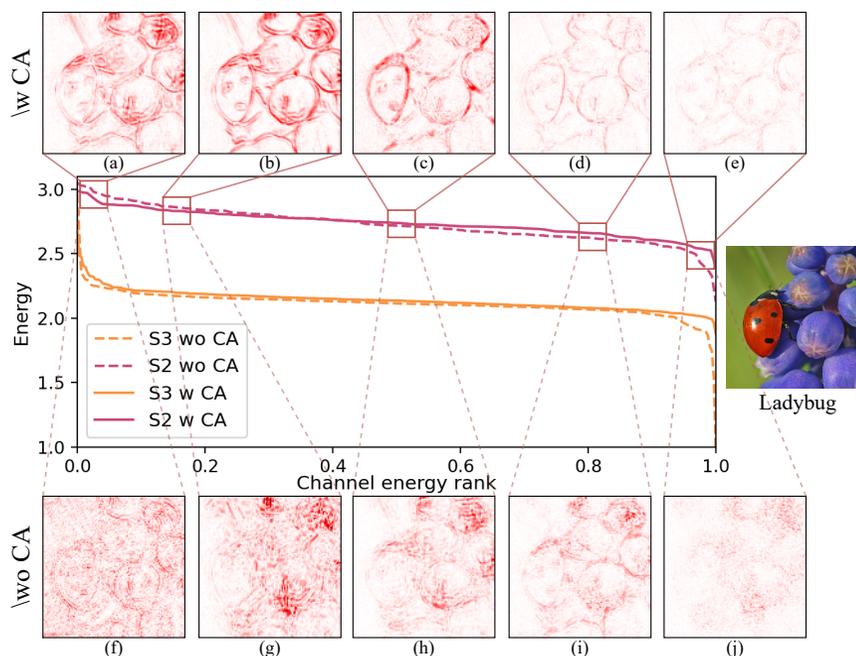
- Channel Aggregation (CA): Multi-order Channel Reallocation.



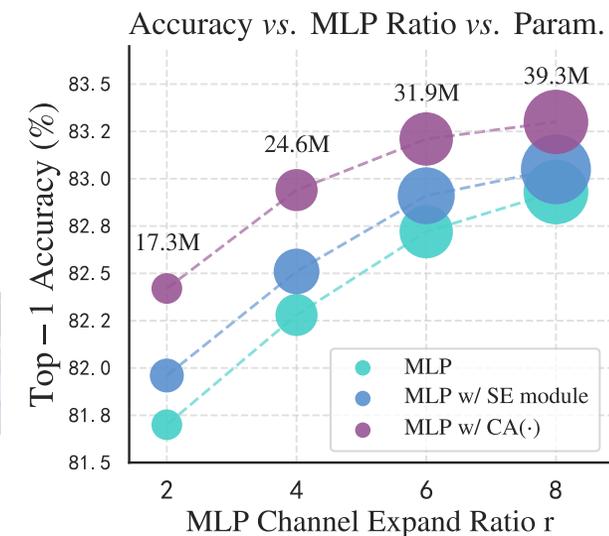
$$Y = \text{GELU}\left(\text{DW}_{3\times 3}\left(\text{Conv}_{1\times 1}\left(\text{Norm}(X)\right)\right)\right),$$

$$Z = \text{Conv}_{1\times 1}(\text{CA}(Y)) + X.$$

$$\text{CA}(X) = X + \gamma_c \odot (X - \text{GELU}(XW_r))$$



Channel energy ranks and channel saliency maps (CSM)<sup>[1]</sup>



Modules	Top-1 Acc (%)	Params. (M)	FLOPs (G)
Baseline	76.6	4.75	1.01
+Gating branch	77.3	5.09	1.07
+DW <sub>7×7</sub>	77.5	5.14	1.09
+Multi-order DW(·)	78.0	5.17	1.10
+FD(·)	78.3	5.18	1.10
+SE module	78.6	5.29	1.14
+CA(·)	<b>79.0</b>	5.20	1.10

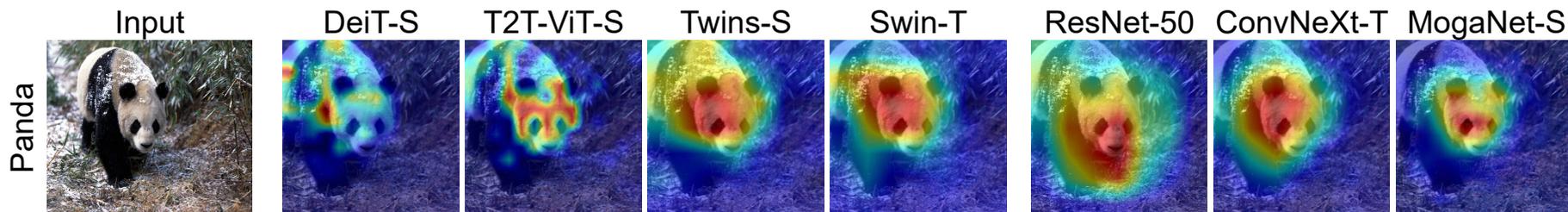
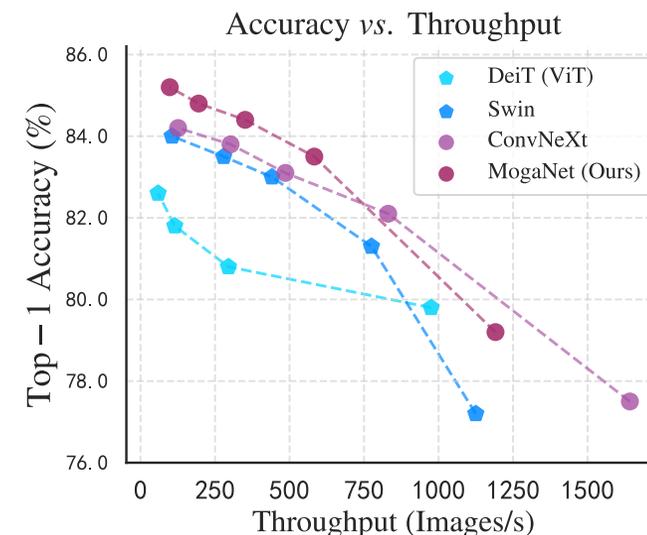
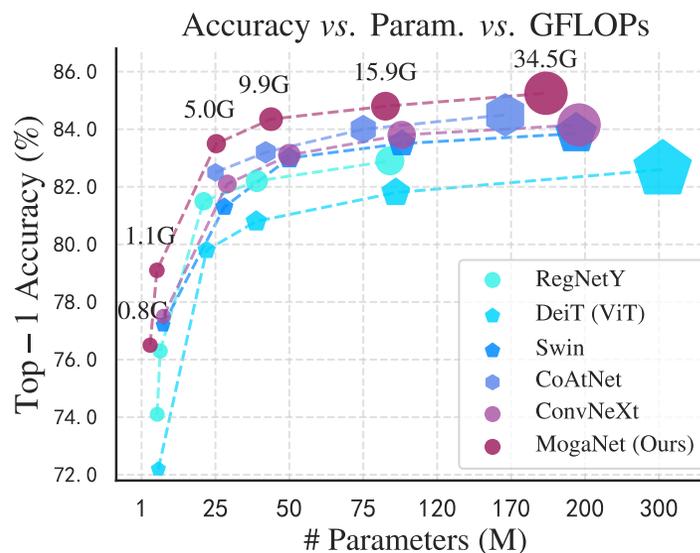
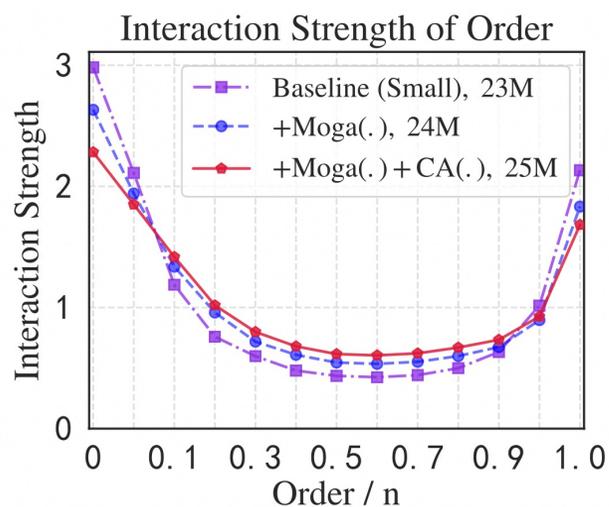
Ablation of MogaNet-S on ImageNet

[1] Reflash dropout in image supe-resolution. CVPR, 2022.

# Multi-order Interaction: MogaNet

- Great scalability and efficiency of parameters.
- Relieving representation bottleneck.

Modules	Top-1 Acc (%)
ConvNeXt-T	82.1
Baseline	82.2
<b>Moga Block</b>	<b>83.4</b>
–FD(·)	83.2
–Multi-DW(·)	83.1
–Moga(·)	82.7
–CA(·)	82.9



# MogaNet: ImageNet Classification

## Light weight (3-10M)

Architecture	Date	Type	Image Param.		FLOPs (G)	Top-1 Acc (%)
			Size	(M)		
ResNet-18	CVPR'2016	C	224 <sup>2</sup>	11.7	1.80	71.5
ShuffleNetV2 2×	ECCV'2018	C	224 <sup>2</sup>	5.5	0.60	75.4
EfficientNet-B0	ICML'2019	C	224 <sup>2</sup>	5.3	0.39	77.1
RegNetY-800MF	CVPR'2020	C	224 <sup>2</sup>	6.3	0.80	76.3
DeiT-T <sup>†</sup>	ICML'2021	T	224 <sup>2</sup>	5.7	1.08	74.1
PVT-T	ICCV'2021	T	224 <sup>2</sup>	13.2	1.60	75.1
T2T-ViT-7	ICCV'2021	T	224 <sup>2</sup>	4.3	1.20	71.7
ViT-C	NIPS'2021	T	224 <sup>2</sup>	4.6	1.10	75.3
SReT-T <sub>Distill</sub>	ECCV'2022	T	224 <sup>2</sup>	4.8	1.10	77.6
PiT-Ti	ICCV'2021	H	224 <sup>2</sup>	4.9	0.70	74.6
LeViT-S	ICCV'2021	H	224 <sup>2</sup>	7.8	0.31	76.6
CoaT-Lite-T	ICCV'2021	H	224 <sup>2</sup>	5.7	1.60	77.5
Swin-1G	ICCV'2021	H	224 <sup>2</sup>	7.3	1.00	77.3
MobileViT-S	ICLR'2022	H	256 <sup>2</sup>	5.6	4.02	78.4
MobileFormer-294M	CVPR'2022	H	224 <sup>2</sup>	11.4	0.59	77.9
ConvNext-XT	CVPR'2022	C	224 <sup>2</sup>	7.4	0.60	77.5
VAN-B0	CVMJ'2023	C	224 <sup>2</sup>	4.1	0.88	75.4
ParC-Net-S	ECCV'2022	C	256 <sup>2</sup>	5.0	3.48	78.6
<b>MogaNet-XT</b>	Ours	C	256 <sup>2</sup>	3.0	1.04	77.2
<b>MogaNet-T</b>	Ours	C	224 <sup>2</sup>	5.2	1.10	79.0
<b>MogaNet-T<sup>§</sup></b>	Ours	C	256 <sup>2</sup>	5.2	1.44	<b>80.0</b>

Architecture	Input size	Learning rate	Warmup epochs	Rand Augment	3-Augment	EMA	Top-1 Acc (%)
MogaNet-XT	224 <sup>2</sup>	1 × 10 <sup>-3</sup>	5	7/0.5	×	×	76.5
MogaNet-XT	224 <sup>2</sup>	2 × 10 <sup>-3</sup>	20	×	✓	×	77.1
MogaNet-XT	256 <sup>2</sup>	1 × 10 <sup>-3</sup>	5	7/0.5	×	×	77.2
MogaNet-XT	256 <sup>2</sup>	2 × 10 <sup>-3</sup>	20	×	✓	×	77.6
MogaNet-T	224 <sup>2</sup>	1 × 10 <sup>-3</sup>	5	7/0.5	×	×	79.0
MogaNet-T	224 <sup>2</sup>	2 × 10 <sup>-3</sup>	20	×	✓	×	79.4
MogaNet-T	256 <sup>2</sup>	1 × 10 <sup>-3</sup>	5	7/0.5	×	×	79.6
MogaNet-T	256 <sup>2</sup>	2 × 10 <sup>-3</sup>	20	×	✓	×	<b>80.0</b>

## Normal size (25-50M)

Architecture	Date	Type	Image Param.		FLOPs (G)	Top-1 Acc (%)
			Size	(M)		
DeiT-S	ICML'2021	T	224 <sup>2</sup>	22	4.6	79.8
Swin-T	ICCV'2021	T	224 <sup>2</sup>	28	4.5	81.3
CSWin-T	CVPR'2022	T	224 <sup>2</sup>	23	4.3	82.8
LITV2-S	NIPS'2022	T	224 <sup>2</sup>	28	3.7	82.0
CoaT-S	ICCV'2021	H	224 <sup>2</sup>	22	12.6	82.1
CoAtNet-0	NIPS'2021	H	224 <sup>2</sup>	25	4.2	82.7
UniFormer-S	ICLR'2022	H	224 <sup>2</sup>	22	3.6	82.9
RegNetY-4GF <sup>†</sup>	CVPR'2020	C	224 <sup>2</sup>	21	4.0	81.5
ConvNeXt-T	CVPR'2022	C	224 <sup>2</sup>	29	4.5	82.1
SLaK-T	ICLR'2023	C	224 <sup>2</sup>	30	5.0	82.5
HorNet-T <sub>7×7</sub>	NIPS'2022	C	224 <sup>2</sup>	22	4.0	82.8
<b>MogaNet-S</b>	Ours	C	224 <sup>2</sup>	25	5.0	<b>83.4</b>
Swin-S	ICCV'2021	T	224 <sup>2</sup>	50	8.7	83.0
Focal-S	NIPS'2021	T	224 <sup>2</sup>	51	9.1	83.6
CSWin-S	CVPR'2022	T	224 <sup>2</sup>	35	6.9	83.6
LITV2-M	NIPS'2022	T	224 <sup>2</sup>	49	7.5	83.3
CoaT-M	ICCV'2021	H	224 <sup>2</sup>	45	9.8	83.6
CoAtNet-1	NIPS'2021	H	224 <sup>2</sup>	42	8.4	83.3
UniFormer-B	ICLR'2022	H	224 <sup>2</sup>	50	8.3	83.9
FAN-B-Hybrid	ICML'2022	H	224 <sup>2</sup>	50	11.3	83.9
EfficientNet-B6	ICML'2019	C	528 <sup>2</sup>	43	19.0	84.0
RegNetY-8GF <sup>†</sup>	CVPR'2020	C	224 <sup>2</sup>	39	8.1	82.2
ConvNeXt-S	CVPR'2022	C	224 <sup>2</sup>	50	8.7	83.1
FocalNet-S (LRF)	NIPS'2022	C	224 <sup>2</sup>	50	8.7	83.5
HorNet-S <sub>7×7</sub>	NIPS'2022	C	224 <sup>2</sup>	50	8.8	84.0
SLaK-S	ICLR'2023	C	224 <sup>2</sup>	55	9.8	83.8
<b>MogaNet-B</b>	Ours	C	224 <sup>2</sup>	44	9.9	<b>84.3</b>

Training and inference at the resolution of 224<sup>2</sup> or 256<sup>2</sup>.

## Large size (80-200M)

DeiT-B	ICML'2021	T	224 <sup>2</sup>	86	17.5	81.8
Swin-B	ICCV'2021	T	224 <sup>2</sup>	89	15.4	83.5
Focal-B	NIPS'2021	T	224 <sup>2</sup>	90	16.4	84.0
CSWin-B	CVPR'2022	T	224 <sup>2</sup>	78	15.0	84.2
DeiT III-B	ECCV'2022	T	224 <sup>2</sup>	87	18.0	83.8
BoTNet-T7	CVPR'2021	H	256 <sup>2</sup>	79	19.3	84.2
CoAtNet-2	NIPS'2021	H	224 <sup>2</sup>	75	15.7	84.1
FAN-B-Hybrid	ICML'2022	H	224 <sup>2</sup>	77	16.9	84.3
RegNetY-16GF	CVPR'2020	C	224 <sup>2</sup>	84	16.0	82.9
ConvNeXt-B	CVPR'2022	C	224 <sup>2</sup>	89	15.4	83.8
RepLkNet-31B	CVPR'2022	C	224 <sup>2</sup>	79	15.3	83.5
FocalNet-B (LRF)	NIPS'2022	C	224 <sup>2</sup>	89	15.4	83.9
HorNet-B <sub>7×7</sub>	NIPS'2022	C	224 <sup>2</sup>	87	15.6	84.3
SLaK-B	ICLR'2023	C	224 <sup>2</sup>	95	17.1	84.0
<b>MogaNet-L</b>	Ours	C	224 <sup>2</sup>	83	15.9	<b>84.7</b>
Swin-L <sup>‡</sup>	ICCV'2021	T	384 <sup>2</sup>	197	104	87.3
DeiT III-L <sup>‡</sup>	ECCV'2022	T	384 <sup>2</sup>	304	191	87.7
CoAtNet-3 <sup>‡</sup>	NIPS'2021	H	384 <sup>2</sup>	168	107	87.6
RepLkNet-31L <sup>‡</sup>	CVPR'2022	C	384 <sup>2</sup>	172	96	86.6
ConvNeXt-L	CVPR'2022	C	224 <sup>2</sup>	198	34.4	84.3
ConvNeXt-L <sup>‡</sup>	CVPR'2022	C	384 <sup>2</sup>	198	101	87.5
ConvNeXt-XL <sup>‡</sup>	CVPR'2022	C	384 <sup>2</sup>	350	179	87.8
HorNet-L <sup>‡</sup>	NIPS'2022	C	384 <sup>2</sup>	202	102	87.7
<b>MogaNet-XL</b>	Ours	C	224 <sup>2</sup>	181	34.5	85.1
<b>MogaNet-XL<sup>‡</sup></b>	Ours	C	384 <sup>2</sup>	181	102	<b>87.8</b>

Architecture	Date	Type	Param. (M)	100-epoch		300-epoch			
				Train	Test	Acc (%)	Train	Test	Acc (%)
ConvNeXt-T (Liu et al., 2022b)	CVPR'2022	C	29	160 <sup>2</sup>	224 <sup>2</sup>	78.8	224 <sup>2</sup>	224 <sup>2</sup>	82.1
ConvNeXt-S (Liu et al., 2022b)	CVPR'2022	C	50	160 <sup>2</sup>	224 <sup>2</sup>	81.7	224 <sup>2</sup>	224 <sup>2</sup>	83.1
ConvNeXt-B (Liu et al., 2022b)	CVPR'2022	C	89	160 <sup>2</sup>	224 <sup>2</sup>	82.1	224 <sup>2</sup>	224 <sup>2</sup>	83.8
ConvNeXt-L (Liu et al., 2022b)	CVPR'2022	C	189	160 <sup>2</sup>	224 <sup>2</sup>	82.8	224 <sup>2</sup>	224 <sup>2</sup>	84.3
ConvNeXt-XL (Liu et al., 2022b)	CVPR'2022	C	350	160 <sup>2</sup>	224 <sup>2</sup>	82.9	224 <sup>2</sup>	224 <sup>2</sup>	84.5
HorNet-T <sub>7×7</sub> (Rao et al., 2022)	NIPS'2022	C	22	160 <sup>2</sup>	224 <sup>2</sup>	80.1	224 <sup>2</sup>	224 <sup>2</sup>	82.8
HorNet-S <sub>7×7</sub> (Rao et al., 2022)	NIPS'2022	C	50	160 <sup>2</sup>	224 <sup>2</sup>	81.2	224 <sup>2</sup>	224 <sup>2</sup>	84.0
VAN-B0 (Guo et al., 2023)	CVMJ'2023	C	4	160 <sup>2</sup>	224 <sup>2</sup>	72.6	224 <sup>2</sup>	224 <sup>2</sup>	75.8
VAN-B2 (Guo et al., 2023)	CVMJ'2023	C	27	160 <sup>2</sup>	224 <sup>2</sup>	81.0	224 <sup>2</sup>	224 <sup>2</sup>	82.8
VAN-B3 (Guo et al., 2023)	CVMJ'2023	C	45	160 <sup>2</sup>	224 <sup>2</sup>	81.9	224 <sup>2</sup>	224 <sup>2</sup>	83.9
<b>MogaNet-XT</b>	Ours	C	3	160 <sup>2</sup>	224 <sup>2</sup>	72.8	224 <sup>2</sup>	224 <sup>2</sup>	76.5
<b>MogaNet-T</b>	Ours	C	5	160 <sup>2</sup>	224 <sup>2</sup>	75.4	224 <sup>2</sup>	224 <sup>2</sup>	79.0
<b>MogaNet-S</b>	Ours	C	25	160 <sup>2</sup>	224 <sup>2</sup>	81.1	224 <sup>2</sup>	224 <sup>2</sup>	83.4
<b>MogaNet-B</b>	Ours	C	44	160 <sup>2</sup>	224 <sup>2</sup>	82.2	224 <sup>2</sup>	224 <sup>2</sup>	84.3
<b>MogaNet-L</b>	Ours	C	83	160 <sup>2</sup>	224 <sup>2</sup>	83.2	224 <sup>2</sup>	224 <sup>2</sup>	84.7

# MogaNet: COCO Object Det. and Seg.

## RetinaNet (1×)

Architecture	Type	#P. (M)	FLOPs		RetinaNet 1×				
			(G)	AP	AP <sub>50</sub>	AP <sub>75</sub>	AP <sup>S</sup>	AP <sub>M</sub>	AP <sub>L</sub>
RegNet-800M	C	17	168	35.6	54.7	37.7	19.7	390	47.8
PVTv2-B0	T	13	160	37.1	57.2	39.2	23.4	40.4	49.2
<b>MogaNet-XT</b>	C	12	167	<b>39.7</b>	<b>60.0</b>	<b>42.4</b>	<b>23.8</b>	<b>43.6</b>	<b>51.7</b>
ResNet-18	C	21	189	31.8	49.6	33.6	16.3	34.3	43.2
RegNet-1.6G	C	20	185	37.4	56.8	39.8	22.4	41.1	49.2
RegNet-3.2G	C	26	218	39.0	58.4	41.9	22.6	43.5	50.8
PVT-T	T	23	183	36.7	56.9	38.9	22.6	38.8	50.0
PoolFormer-S12	T	22	207	36.2	56.2	38.2	20.8	39.1	48.0
PVTv2-B1	T	24	187	41.1	61.4	43.8	26.0	44.6	54.6
<b>MogaNet-T</b>	C	14	173	<b>41.4</b>	<b>61.5</b>	<b>44.4</b>	<b>25.1</b>	<b>45.7</b>	<b>53.6</b>
ResNet-50	C	37	239	36.3	55.3	38.6	19.3	40.0	48.8
Swin-T	T	38	245	41.8	62.6	44.7	25.2	45.8	54.7
PVT-S	T	34	226	40.4	61.3	43.0	25.0	42.9	55.7
Twins-SVT-S	T	34	209	42.3	63.4	45.2	26.0	45.5	56.5
Focal-T	T	39	265	43.7	-	-	-	-	-
PoolFormer-S36	T	41	272	39.5	60.5	41.8	22.5	42.9	52.4
PVTv2-B2	T	35	281	44.6	65.7	47.6	28.6	48.5	59.2
CMT-S	H	45	231	44.3	65.5	47.5	27.1	48.3	59.1
<b>MogaNet-S</b>	C	35	253	<b>45.8</b>	<b>66.6</b>	<b>49.0</b>	<b>29.1</b>	<b>50.1</b>	<b>59.8</b>
ResNet-101	C	57	315	38.5	57.8	41.2	21.4	42.6	51.1
PVT-M	T	54	258	41.9	63.1	44.3	25.0	44.9	57.6
Focal-S	T	62	367	45.6	-	-	-	-	-
PVTv2-B3	T	55	263	46.0	67.0	49.5	28.2	50.0	61.3
PVTv2-B4	T	73	315	46.3	67.0	49.6	29.0	50.1	62.7
<b>MogaNet-B</b>	C	54	355	<b>47.7</b>	<b>68.9</b>	<b>51.0</b>	<b>30.5</b>	<b>52.2</b>	<b>61.7</b>
ResNeXt-101-64	C	95	473	41.0	60.9	44.0	23.9	45.2	54.0
PVTv2-B5	T	92	335	46.1	66.6	49.5	27.8	50.2	62.0
<b>MogaNet-L</b>	C	92	477	<b>48.7</b>	<b>69.5</b>	<b>52.6</b>	<b>31.5</b>	<b>53.4</b>	<b>62.7</b>

Inference input size 800×1280

## Mask R-CNN (1×)

Architecture	Type	#P. (M)	FLOPs		Mask R-CNN 1×					
			(G)	AP <sup>b</sup>	AP <sub>50</sub> <sup>b</sup>	AP <sub>75</sub> <sup>b</sup>	AP <sup>m</sup>	AP <sub>50</sub> <sup>m</sup>	AP <sub>75</sub> <sup>m</sup>	
RegNet-800M	C	27	187	37.5	57.9	41.1	34.3	56.0	36.8	
<b>MogaNet-XT</b>	C	23	185	<b>40.7</b>	<b>62.3</b>	<b>44.4</b>	<b>37.6</b>	<b>59.6</b>	<b>40.2</b>	
ResNet-18	C	31	207	34.0	54.0	36.7	31.2	51.0	32.7	
RegNet-1.6G	C	29	204	38.9	60.5	43.1	35.7	57.4	38.9	
PVT-T	T	33	208	36.7	59.2	39.3	35.1	56.7	37.3	
PoolFormer-S12	T	32	207	37.3	59.0	40.1	34.6	55.8	36.9	
<b>MogaNet-T</b>	C	25	192	<b>42.6</b>	<b>64.0</b>	<b>46.4</b>	<b>39.1</b>	<b>61.3</b>	<b>42.0</b>	
ResNet-50	C	44	260	38.0	58.6	41.4	34.4	55.1	36.7	
RegNet-6.4G	C	45	307	41.1	62.3	45.2	37.1	59.2	39.6	
PVT-S	T	44	245	40.4	62.9	43.8	37.8	60.1	40.3	
Swin-T	T	48	264	42.2	64.6	46.2	39.1	61.6	42.0	
MViT-T	T	46	326	45.9	<b>68.7</b>	50.5	42.1	<b>66.0</b>	45.4	
PoolFormer-S36	T	32	207	41.0	63.1	44.8	37.7	60.1	40.0	
Focal-T	T	49	291	44.8	67.7	49.2	41.0	64.7	44.2	
PVTv2-B2	T	45	309	45.3	67.1	49.6	41.2	64.2	44.4	
LITv2-S	T	47	261	44.9	67.0	49.5	40.8	63.8	44.2	
CMT-S	H	45	249	44.6	66.8	48.9	40.7	63.9	43.4	
Conformer-S/16	H	58	341	43.6	65.6	47.7	39.7	62.6	42.5	
UniFormer-S	H	41	269	45.6	68.1	49.7	41.6	64.8	45.0	
ConvNeXt-T	C	48	262	44.2	66.6	48.3	40.1	63.3	42.8	
FocalNet-T (SRF)	C	49	267	45.9	68.3	50.1	41.3	65.0	44.3	
FocalNet-T (LRF)	C	49	268	46.1	68.2	50.6	41.5	65.1	44.5	
<b>MogaNet-S</b>	C	45	272	<b>46.7</b>	<b>68.0</b>	<b>51.3</b>	<b>42.2</b>	<b>65.4</b>	<b>45.5</b>	
ResNet-101	C	63	336	40.4	61.1	44.2	36.4	57.7	38.8	
RegNet-12G	C	64	423	42.2	63.7	46.1	38.0	60.5	40.5	
PVT-M	T	64	302	42.0	64.4	45.6	39.0	61.6	42.1	
Swin-S	T	69	354	44.8	66.6	48.9	40.9	63.4	44.2	
Focal-S	T	71	401	47.4	69.8	51.9	42.8	66.6	46.1	
PVTv2-B3	T	65	397	47.0	68.1	51.7	42.5	65.7	45.7	
LITv2-M	T	68	315	46.5	68.0	50.9	42.0	65.1	45.0	
UniFormer-B	H	69	399	47.4	69.7	52.1	43.1	66.0	46.5	
ConvNeXt-S	C	70	348	45.4	67.9	50.0	41.8	65.2	45.1	
<b>MogaNet-B</b>	C	63	373	<b>47.9</b>	<b>70.0</b>	<b>52.7</b>	<b>43.2</b>	<b>67.0</b>	<b>46.6</b>	
Swin-B	T	107	496	46.9	69.6	51.2	42.3	65.9	45.6	
PVTv2-B5	T	102	557	47.4	68.6	51.9	42.5	65.7	46.0	
ConvNeXt-B	C	108	486	47.0	69.4	51.7	42.7	66.3	46.0	
FocalNet-B (SRF)	C	109	496	48.8	70.7	53.5	43.3	67.5	46.5	
<b>MogaNet-L</b>	C	102	495	<b>49.4</b>	<b>70.7</b>	<b>54.1</b>	<b>44.1</b>	<b>68.1</b>	<b>47.6</b>	

## Cascade Mask R-CNN (3×)

Architecture	Type	#P. (M)	FLOPs		Cascade Mask R-CNN +MS 3×					
			(G)	AP <sup>bb</sup>	AP <sub>50</sub> <sup>b</sup>	AP <sub>75</sub> <sup>b</sup>	AP <sup>m</sup>	AP <sub>50</sub> <sup>m</sup>	AP <sub>75</sub> <sup>m</sup>	
ResNet-50	C	77	739	46.3	64.3	50.5	40.1	61.7	43.4	
Swin-T	T	86	745	50.4	69.2	54.7	43.7	66.6	47.3	
Focal-T	T	87	770	51.5	70.6	55.9	-	-	-	
ConvNeXt-T	C	86	741	50.4	69.1	54.8	43.7	66.5	47.3	
FocalNet-T (SRF)	C	86	746	51.5	70.1	55.8	44.6	67.7	48.4	
<b>MogaNet-S</b>	C	78	750	<b>51.6</b>	<b>70.8</b>	<b>56.3</b>	<b>45.1</b>	<b>68.7</b>	<b>48.8</b>	
ResNet-101-32	C	96	819	48.1	66.5	52.4	41.6	63.9	45.2	
Swin-S	T	107	838	51.9	70.7	56.3	45.0	68.2	48.8	
ConvNeXt-S	C	108	827	51.9	70.8	56.5	45.0	68.4	49.1	
<b>MogaNet-B</b>	C	101	851	<b>52.6</b>	<b>72.0</b>	<b>57.3</b>	<b>46.0</b>	<b>69.6</b>	<b>49.7</b>	
Swin-B	T	145	982	51.9	70.5	56.4	45.0	68.1	48.9	
ConvNeXt-B	C	146	964	52.7	71.3	57.2	45.6	68.9	49.5	
<b>MogaNet-L</b>	C	140	974	<b>53.3</b>	<b>71.8</b>	<b>57.8</b>	<b>46.1</b>	<b>69.2</b>	<b>49.8</b>	
Swin-L <sup>‡</sup>	T	253	1382	53.9	72.4	58.8	46.7	70.1	50.8	
ConvNeXt-L <sup>‡</sup>	C	255	1354	54.8	73.8	59.8	47.6	71.3	51.7	
ConvNeXt-XL <sup>‡</sup>	C	407	1898	55.2	74.2	59.9	47.7	71.6	52.2	
RepLkNet-31L <sup>‡</sup>	C	229	1321	53.9	72.5	58.6	46.5	70.0	50.6	
HorNet-L <sup>‡</sup>	C	259	1399	56.0	-	-	48.6	-	-	
<b>MogaNet-XL<sup>‡</sup></b>	C	238	1355	<b>56.2</b>	<b>75.0</b>	<b>61.2</b>	<b>48.8</b>	<b>72.6</b>	<b>53.3</b>	

- Object Detection: RetinaNet.
- Instance Segmentation: (Cascade) Mask R-CNN.
- Multi-scale fine-tuning with IN-21K pre-trained models.

# MogaNet: ADE20K Semantic Segmentation

## Semantic FPN (80K)

Method	Architecture	Date	Crop size	Param. (M)	FLOPs (G)	mIoU <sup>ss</sup> (%)
Semantic FPN (80K)	PVT-S	ICCV'2021	512 <sup>2</sup>	28	161	39.8
	Twins-S	NIPS'2021	512 <sup>2</sup>	28	162	44.3
	Swin-T	ICCV'2021	512 <sup>2</sup>	32	182	41.5
	Uniformer-S	ICLR'2022	512 <sup>2</sup>	25	247	46.6
	LITV2-S	NIPS'2022	512 <sup>2</sup>	31	179	44.3
	VAN-B2	CVMJ'2023	512 <sup>2</sup>	30	164	46.7
<b>MogaNet-S</b>	Ours		512 <sup>2</sup>	29	189	<b>47.7</b>

## MogaNet + Semantic FPN

Method	Backbone	Pretrain	Params	FLOPs	Iters	mIoU	mAcc
Semantic FPN	MogaNet-XT	ImageNet-1K	6.9M	101.4G	80K	40.3	52.4
Semantic FPN	MogaNet-T	ImageNet-1K	9.1M	107.8G	80K	43.1	55.4
Semantic FPN	MogaNet-S	ImageNet-1K	29.1M	189.7G	80K	47.7	59.8
Semantic FPN	MogaNet-B	ImageNet-1K	47.5M	293.6G	80K	49.3	61.6
Semantic FPN	MogaNet-L	ImageNet-1K	86.2M	418.7G	80K	50.2	63.0

- Semantic FPN (80K) with 512×2048 inference resolutions.
- UperNet (160K) with 512×2048 or 640×2560 inference resolutions using IN-1K or IN-21K models.

## ADE20K UperNet (160K)

Architecture	Date	Type	Crop size	Param. (M)	FLOPs (G)	mIoU <sup>ss</sup> (%)
ResNet-18	CVPR'2016	C	512 <sup>2</sup>	41	885	39.2
<b>MogaNet-XT</b>	Ours	C	512 <sup>2</sup>	30	856	<b>42.2</b>
ResNet-50	CVPR'2016	C	512 <sup>2</sup>	67	952	42.1
<b>MogaNet-T</b>	Ours	C	512 <sup>2</sup>	33	862	<b>43.7</b>
DeiT-S	ICML'2021	T	512 <sup>2</sup>	52	1099	44.0
Swin-T	ICCV'2021	T	512 <sup>2</sup>	60	945	46.1
TwinsP-S	NIPS'2021	T	512 <sup>2</sup>	55	919	46.2
Twins-S	NIPS'2021	T	512 <sup>2</sup>	54	901	46.2
Focal-T	NIPS'2021	T	512 <sup>2</sup>	62	998	45.8
Uniformer-S <sub>h32</sub>	ICLR'2022	H	512 <sup>2</sup>	52	955	47.0
UniFormer-S	ICLR'2022	H	512 <sup>2</sup>	52	1008	47.6
ConvNeXt-T	CVPR'2022	C	512 <sup>2</sup>	60	939	46.7
FocalNet-T (SRF)	NIPS'2022	C	512 <sup>2</sup>	61	944	46.5
HorNet-T <sub>7×7</sub>	NIPS'2022	C	512 <sup>2</sup>	52	926	48.1
<b>MogaNet-S</b>	Ours	C	512 <sup>2</sup>	55	946	<b>49.2</b>
Swin-S	ICCV'2021	T	512 <sup>2</sup>	81	1038	48.1
Twins-B	NIPS'2021	T	512 <sup>2</sup>	89	1020	47.7
Focal-S	NIPS'2021	T	512 <sup>2</sup>	85	1130	48.0
Uniformer-B <sub>h32</sub>	ICLR'2022	H	512 <sup>2</sup>	80	1106	49.5
ConvNeXt-S	CVPR'2022	C	512 <sup>2</sup>	82	1027	48.7
FocalNet-S (SRF)	NIPS'2022	C	512 <sup>2</sup>	83	1035	49.3
SLaK-S	ICLR'2023	C	512 <sup>2</sup>	91	1028	49.4
<b>MogaNet-B</b>	Ours	C	512 <sup>2</sup>	74	1050	<b>50.1</b>
Swin-B	ICCV'2021	T	512 <sup>2</sup>	121	1188	49.7
Focal-B	NIPS'2021	T	512 <sup>2</sup>	126	1354	49.0
ConvNeXt-B	CVPR'2022	C	512 <sup>2</sup>	122	1170	49.1
RepLKNet-31B	CVPR'2022	C	512 <sup>2</sup>	112	1170	49.9
FocalNet-B (SRF)	NIPS'2022	C	512 <sup>2</sup>	124	1180	50.2
SLaK-B	ICLR'2023	C	512 <sup>2</sup>	135	1185	50.2
<b>MogaNet-L</b>	Ours	C	512 <sup>2</sup>	113	1176	<b>50.9</b>
Swin-L <sup>‡</sup>	ICCV'2021	T	640 <sup>2</sup>	234	2468	52.1
ConvNeXt-L <sup>‡</sup>	CVPR'2022	C	640 <sup>2</sup>	245	2458	53.7
RepLKNet-31L <sup>‡</sup>	CVPR'2022	C	640 <sup>2</sup>	207	2404	52.4
<b>MogaNet-XL<sup>‡</sup></b>	Ours	C	640 <sup>2</sup>	214	2451	<b>54.0</b>

# MogaNet: 2D/3D Pose Estimation

## COCO 2D Human Pose with TopDown baseline (256×192)

Architecture	Type	Crop size	#P. (M)	FLOPs (G)	AP (%)	AP <sup>50</sup> (%)	AP <sup>75</sup> (%)	AR (%)
MobileNetV2	C	256 × 192	10	1.6	64.6	87.4	72.3	70.7
ShuffleNetV2 2×	C	256 × 192	8	1.4	59.9	85.4	66.3	66.4
<b>MogaNet-XT</b>	C	256 × 192	6	1.8	<b>72.1</b>	<b>89.7</b>	<b>80.1</b>	<b>77.7</b>
RSN-18	C	256 × 192	9	2.3	70.4	88.7	77.9	77.1
<b>MogaNet-T</b>	C	256 × 192	8	2.2	<b>73.2</b>	<b>90.1</b>	<b>81.0</b>	<b>78.8</b>
ResNet-50	C	256 × 192	34	5.5	72.1	89.9	80.2	77.6
HRNet-W32	C	256 × 192	29	7.1	74.4	90.5	81.9	78.9
Swin-T	T	256 × 192	33	6.1	72.4	90.1	80.6	78.2
PVT-S	T	256 × 192	28	4.1	71.4	89.6	79.4	77.3
PVTV2-B2	T	256 × 192	29	4.3	73.7	90.5	81.2	79.1
Uniformer-S	H	256 × 192	25	4.7	74.0	90.3	82.2	79.5
ConvNeXt-T	C	256 × 192	33	5.5	73.2	90.0	80.9	78.8
<b>MogaNet-S</b>	C	256 × 192	29	6.0	<b>74.9</b>	<b>90.7</b>	<b>82.8</b>	<b>80.1</b>
ResNet-101	C	256 × 192	53	12.4	71.4	89.3	79.3	77.1
ResNet-152	C	256 × 192	69	15.7	72.0	89.3	79.8	77.8
HRNet-W48	C	256 × 192	64	14.6	75.1	90.6	82.2	80.4
Swin-B	T	256 × 192	93	18.6	72.9	89.9	80.8	78.6
Swin-L	T	256 × 192	203	40.3	74.3	90.6	82.1	79.8
Uniformer-B	H	256 × 192	54	9.2	75.0	90.6	83.0	80.4
ConvNeXt-S	C	256 × 192	55	9.7	73.7	90.3	81.9	79.3
ConvNeXt-B	C	256 × 192	94	16.4	74.0	90.7	82.1	79.5
<b>MogaNet-B</b>	C	256 × 192	47	10.9	<b>75.3</b>	<b>90.9</b>	<b>83.3</b>	<b>80.7</b>

Architecture	Type	Crop size	#P. (M)	FLOPs (G)	AP (%)	AP <sup>50</sup> (%)	AP <sup>75</sup> (%)	AR (%)
MobileNetV2	C	384 × 288	10	3.6	67.3	87.9	74.3	72.9
ShuffleNetV2 2×	C	384 × 288	8	3.1	63.6	86.5	70.5	69.7
<b>MogaNet-XT</b>	C	384 × 288	6	4.2	<b>74.7</b>	<b>90.1</b>	<b>81.3</b>	<b>79.9</b>
RSN-18	C	384 × 288	9	5.1	72.1	89.5	79.8	78.6
<b>MogaNet-T</b>	C	384 × 288	8	4.9	<b>75.7</b>	<b>90.6</b>	<b>82.6</b>	<b>80.9</b>
HRNet-W32	C	384 × 288	29	16.0	75.8	90.6	82.7	81.0
Uniformer-S	H	384 × 288	25	11.1	75.9	90.6	83.4	81.4
ConvNeXt-T	C	384 × 288	33	33.1	75.3	90.4	82.1	80.5
<b>MogaNet-S</b>	C	384 × 288	29	13.5	<b>76.4</b>	<b>91.0</b>	<b>83.3</b>	<b>81.4</b>
ResNet-152	C	384 × 288	69	35.6	74.3	89.6	81.1	79.7
HRNet-W48	C	384 × 288	64	32.9	76.3	90.8	82.0	81.2
Swin-B	T	384 × 288	93	39.2	74.9	90.5	81.8	80.3
Swin-L	T	384 × 288	203	86.9	76.3	91.2	83.0	81.4
HRFormer-B	T	384 × 288	54	30.7	77.2	91.0	83.6	82.0
ConvNeXt-S	C	384 × 288	55	21.8	75.8	90.7	83.1	81.0
ConvNeXt-B	C	384 × 288	94	36.6	75.9	90.6	83.1	81.1
Uniformer-B	C	384 × 288	54	14.8	76.7	90.8	84.0	81.4
<b>MogaNet-B</b>	C	384 × 288	47	24.4	<b>77.3</b>	<b>91.4</b>	<b>84.0</b>	<b>82.2</b>

## COCO 2D Human Pose with TopDown baseline (384×288)

Architecture	Type	Hand			Face		
		#P. (M)	FLOPs (G)	PA-MPJPE (mm)↓	#P. (M)	FLOPs (G)	3DRMSE ↓
MobileNetV2	C	4.8	0.3	8.33	4.9	0.4	2.64
ResNet-18	C	13.0	1.8	7.51	13.1	2.4	2.40
<b>MogaNet-T</b>	C	<b>6.5</b>	<b>1.1</b>	<b>6.82</b>	<b>6.6</b>	<b>1.5</b>	<b>2.36</b>
ResNet-50	C	26.9	4.1	6.85	27.0	5.4	2.48
ResNet-101	C	45.9	7.9	6.44	46.0	10.3	2.47
DeiT-S	T	23.4	4.3	7.86	23.5	5.5	2.52
Swin-T	T	30.2	4.6	6.97	30.3	6.1	2.45
Swin-S	T	51.0	13.8	6.50	50.9	8.5	2.48
ConvNeXt-T	C	29.9	4.5	6.18	30.0	5.8	2.34
ConvNeXt-S	C	51.5	8.7	6.04	51.6	11.4	2.27
HorNet-T	C	23.7	4.3	6.46	23.8	5.6	2.39
<b>MogaNet-S</b>	C	<b>26.6</b>	<b>5.0</b>	<b>6.08</b>	<b>26.7</b>	<b>6.5</b>	<b>2.24</b>

## 3D Human Pose with Expose

- 3D Face: FFHQ (256<sup>2</sup>)
- 3D Hand: FreiHand (224<sup>2</sup>)

# MogaNet: Video Prediction

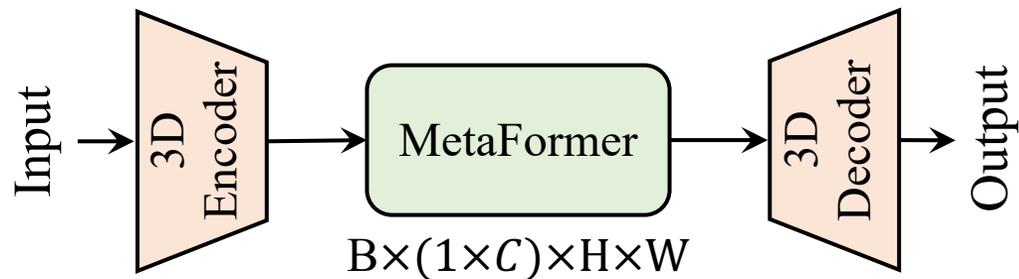
## Moving MNIST ( $10 \times 1 \times 64 \times 64$ )

Architecture	#P. (M)	FLOPs (G)	FPS (s)	200 epochs			2000 epochs		
				MSE↓	MAE↓	SSIM↑	MSE↓	MAE↓	SSIM↑
ViT	46.1	16.9	290	35.15	95.87	0.9139	19.74	61.65	0.9539
Swin	46.1	16.4	294	29.70	84.05	0.9331	19.11	59.84	0.9584
Uniformer	44.8	16.5	296	30.38	85.87	0.9308	18.01	57.52	0.9609
MLP-Mixer	38.2	14.7	334	29.52	83.36	0.9338	18.85	59.86	0.9589
ConvMixer	3.9	5.5	658	32.09	88.93	0.9259	22.30	67.37	0.9507
Poolformer	37.1	14.1	341	31.79	88.48	0.9271	20.96	64.31	0.9539
SimVP	58.0	19.4	209	32.15	89.05	0.9268	21.15	64.15	0.9536
ConvNeXt	37.3	14.1	344	26.94	77.23	0.9397	17.58	55.76	0.9617
VAN	44.5	16.0	288	26.10	76.11	0.9417	16.21	53.57	0.9646
HorNet	45.7	16.3	287	29.64	83.26	0.9331	17.40	55.70	0.9624
<b>MogaNet</b>	46.8	16.5	255	<b>25.57</b>	<b>75.19</b>	<b>0.9429</b>	<b>15.67</b>	<b>51.84</b>	<b>0.9661</b>

## MMNIST-CIFAR ( $10 \times 3 \times 64 \times 64$ )

	Method	Params (M)	FLOPs (G)	FPS	MSE ↓	MAE ↓	SSIM ↑	PSNR ↑
Recurrent-based	ConvLSTM	15.0	56.8	113	73.31	338.56	0.9204	23.09
	PredNet	12.5	8.4	659	286.70	514.14	0.8139	17.49
	PredRNN	23.8	116.0	54	50.09	225.04	0.9499	24.90
	PredRNN++	38.6	171.7	38	<b>44.19</b>	198.27	<b>0.9567</b>	<b>25.60</b>
	MIM	38.0	179.2	37	48.63	213.44	0.9521	25.08
	E3D-LSTM	51.0	298.9	18	80.79	214.86	0.9314	22.89
	PhyDNet	3.1	15.3	182	142.54	700.37	0.8276	19.92
	MAU	4.5	17.8	201	58.84	255.76	0.9408	24.19
	PredRNNv2	23.9	116.6	52	57.27	252.29	0.9419	24.24
	DMVFN	3.5	0.2	1145	298.73	606.92	0.7765	17.07
Recurrent-free	SimVP	58.0	19.4	209	59.83	214.54	0.9414	24.15
	TAU	44.7	16.0	283	48.17	<b>177.35</b>	0.9539	25.21
	SimVPv2	46.8	16.5	282	51.13	185.13	0.9512	24.93
	ViT	46.1	16.9	290	64.94	234.01	0.9354	23.90
	Swin Transformer	46.1	16.4	294	57.11	207.45	0.9443	24.34
	Uniformer	44.8	16.5	296	56.96	207.51	0.9442	24.38
	MLP-Mixer	38.2	14.7	334	57.03	206.46	0.9446	24.34
	ConvMixer	3.9	5.5	658	59.29	219.76	0.9403	24.17
	Poolformer	37.1	14.1	341	60.98	219.50	0.9399	24.16
	ConvNext	37.3	14.1	344	51.39	187.17	0.9503	24.89
VAN	44.5	16.0	288	59.59	221.32	0.9398	25.20	
HorNet	45.7	16.3	287	55.79	202.73	0.9456	24.49	
<b>MogaNet</b>	46.8	16.5	255	49.48	184.11	0.9521	25.07	

- Replacing the MetaFormer blocks in SimVP.
- Comparison with MMNIST and MMNIST-CIFAR.





# State-Space Models: Mamba

Structured state space  $h'(t) = Ah(t) + Bx(t)$  (1a)

$h_t = \bar{A}h_{t-1} + \bar{B}x_t$  (2a)

$\bar{K} = (C\bar{B}, C\bar{A}\bar{B}, \dots, C\bar{A}^k\bar{B}, \dots)$  (3a)

sequence models (S4)  $y(t) = Ch(t)$  (1b)

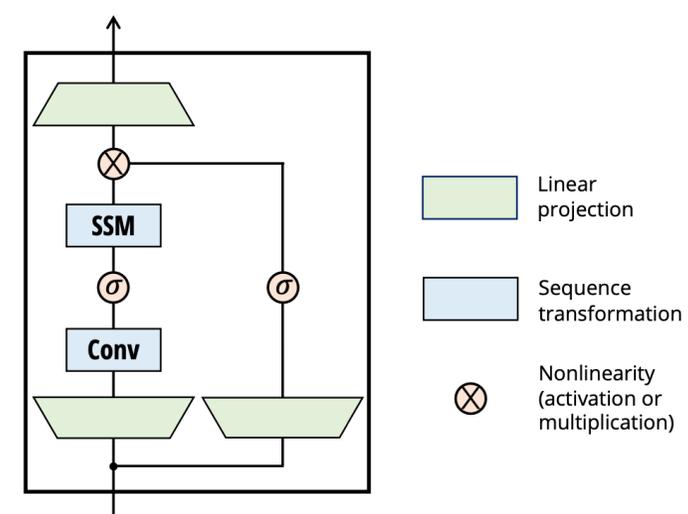
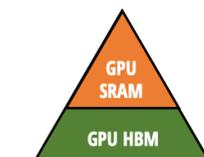
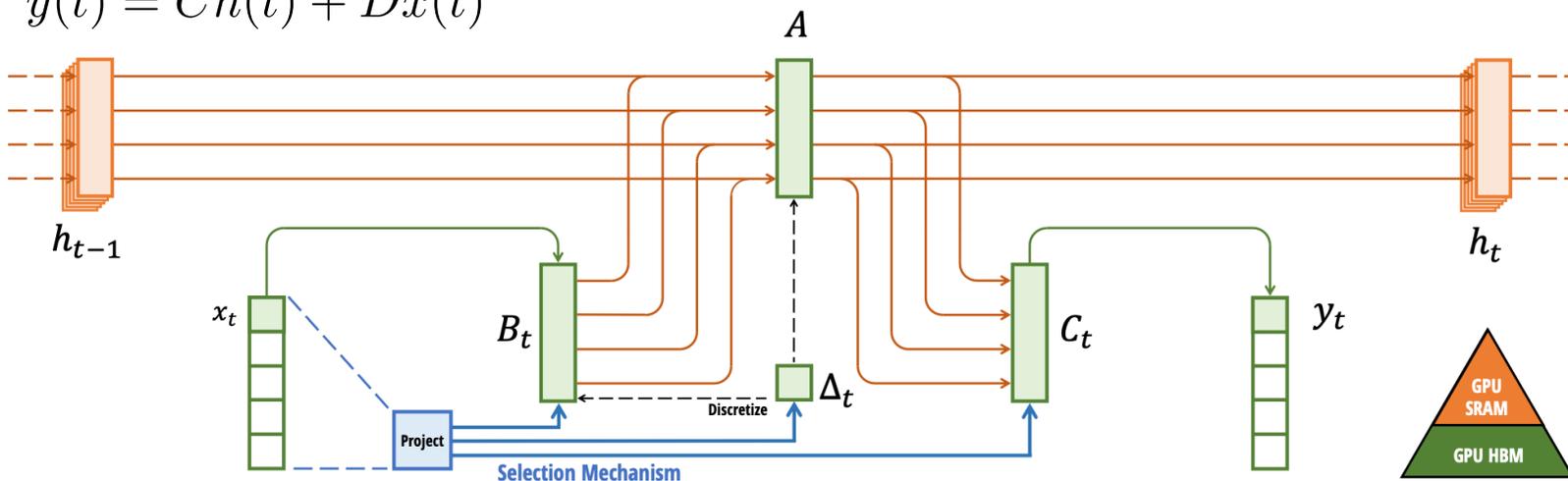
$y_t = Ch_t$  (2b)

$y = x * \bar{K}$  (3b)

$x(t) \in \mathbb{R}^L \rightarrow y(t) \in \mathbb{R}^L, A \in \mathbb{C}^{N \times N}, B, C \in \mathbb{C}^N, D \in \mathbb{C}^1$

$h'(t) = Ah(t) + Bx(t)$

$y(t) = Ch(t) + Dx(t)$



Mamba

$g_t = \sigma(\text{Linear}(x_t))$

$h_t = (1 - g_t)h_{t-1} + g_t x_t$

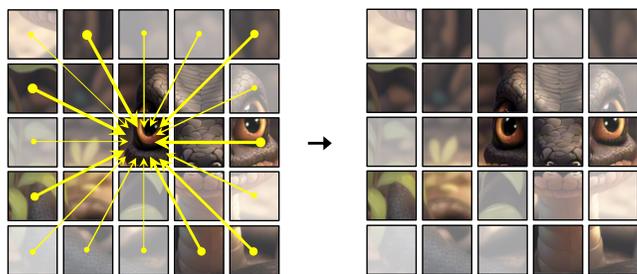
Model	Params	Accuracy (%) at Sequence Length					
		2 <sup>10</sup>	2 <sup>12</sup>	2 <sup>14</sup>	2 <sup>16</sup>	2 <sup>18</sup>	2 <sup>20</sup>
HyenaDNA	1.4M	28.04	28.43	41.17	42.22	31.10	54.87
Mamba	1.4M	31.47	27.50	27.66	40.72	42.41	<b>71.67</b>
Mamba	7M	30.00	29.01	31.48	43.73	56.60	<b>81.31</b>

Great Apes DNA Classification

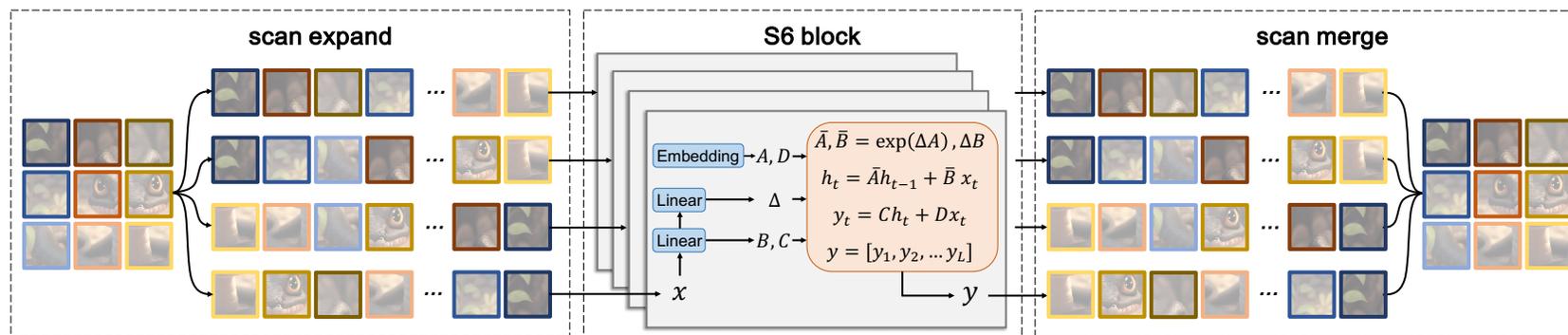
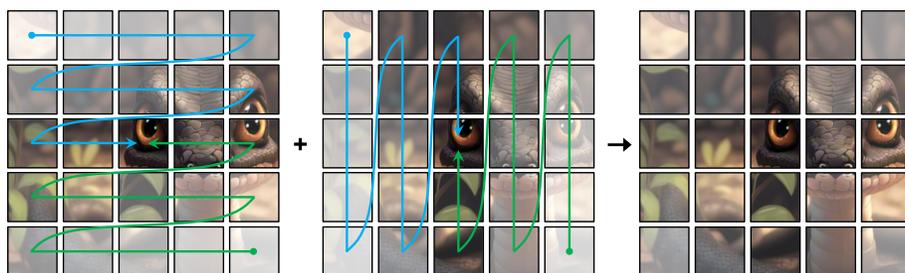
[1] Linear-Time Sequence Modeling with Selective State Spaces. arXiv, 2023.

# State-Space Models: VMamba

(a) Attention  
 $O(N^2)$  complexity



(b) Cross-Scan  
 $O(N)$  complexity



## ADE20K Segmentation

method	crop size	mIoU (SS)	mIoU (MS)	#param.	FLOPs
ResNet-50	512 <sup>2</sup>	42.1	42.8	67M	953G
DeiT-S + MLN	512 <sup>2</sup>	43.8	45.1	58M	1217G
Swin-T	512 <sup>2</sup>	44.4	45.8	60M	945G
ConvNeXt-T	512 <sup>2</sup>	46.0	46.7	60M	939G
VMamba-T	512 <sup>2</sup>	47.3	48.3	55M	939G

## ImageNet-1K Classification

method	image size	#param.	FLOPs	ImageNet top-1 acc.
RegNetY-4G [36]	224 <sup>2</sup>	21M	4.0G	80.0
RegNetY-8G [36]	224 <sup>2</sup>	39M	8.0G	81.7
RegNetY-16G [36]	224 <sup>2</sup>	84M	16.0G	82.9
EffNet-B3 [42]	300 <sup>2</sup>	12M	1.8G	81.6
EffNet-B4 [42]	380 <sup>2</sup>	19M	4.2G	82.9
EffNet-B5 [42]	456 <sup>2</sup>	30M	9.9G	83.6
EffNet-B6 [42]	528 <sup>2</sup>	43M	19.0G	84.0
ViT-B/16 [10]	384 <sup>2</sup>	86M	55.4G	77.9
ViT-L/16 [10]	384 <sup>2</sup>	307M	190.7G	76.5
DeiT-S [45]	224 <sup>2</sup>	22M	4.6G	79.8
DeiT-B [45]	224 <sup>2</sup>	86M	17.5G	81.8
DeiT-B [45]	384 <sup>2</sup>	86M	55.4G	83.1
Swin-T [28]	224 <sup>2</sup>	29M	4.5G	81.3
Swin-S [28]	224 <sup>2</sup>	50M	8.7G	83.0
Swin-B [28]	224 <sup>2</sup>	88M	15.4G	83.5
S4ND-ViT-B [35]	224 <sup>2</sup>	89M	-	80.4
VMamba-T	224 <sup>2</sup>	22M	4.5G	82.2
VMamba-S	224 <sup>2</sup>	44M	9.1G	83.5

# Thank you!

---



Paper: MogaNet



Code: MogaNet



Homepage



lisiyuan@westlake.edu.cn