

OverLoCK: An Overview-first-Look-Closely-next ConvNet with Context-Mixing Dynamic Kernels

Meng Lou

Yizhou Yu

School of Computing and Data Science, The University of Hong Kong

louloum@connect.hku.hk, yizhouy@acm.org

Abstract

Top-down attention plays a crucial role in the human vision system, wherein the brain initially obtains a rough overview of a scene to discover salient cues (i.e., overview first), followed by a more careful finer-grained examination (i.e., look closely next). However, modern ConvNets remain confined to a pyramid structure that successively downsamples the feature map for receptive field expansion, neglecting this crucial biomimetic principle. We present OverLoCK, the first pure ConvNet backbone architecture that explicitly incorporates a top-down attention mechanism. Unlike pyramid backbone networks, our design features a branched architecture with three synergistic sub-networks: 1) a Base-Net that encodes low/mid-level features; 2) a lightweight Overview-Net that generates dynamic top-down attention through coarse global context modeling (i.e., overview first); and 3) a robust Focus-Net that performs finer-grained perception guided by top-down attention (i.e., look closely next). To fully unleash the power of top-down attention, we further propose a novel context-mixing dynamic convolution (ContMix) that effectively models long-range dependencies while preserving inherent local inductive biases even when the input resolution increases, addressing critical limitations in existing convolutions. Our OverLoCK exhibits a notable performance improvement over existing methods. For instance, OverLoCK-T achieves a Top-1 accuracy of 84.2%, significantly surpassing ConvNeXt-B while using only around one-third of the FLOPs/parameters. On object detection, our OverLoCK-S clearly surpasses MogaNet-B by 1% in AP^b. On semantic segmentation, our OverLoCK-T remarkably improves UniRepLNet-T by 1.7% in mIoU. Code is publicly available at <https://bit.ly/OverLoCK>.

1. Introduction

Top-down neural attention [17, 40, 58] is a crucial perception mechanism in the human vision system, which suggests that the brain initially processes a visual scene to quickly

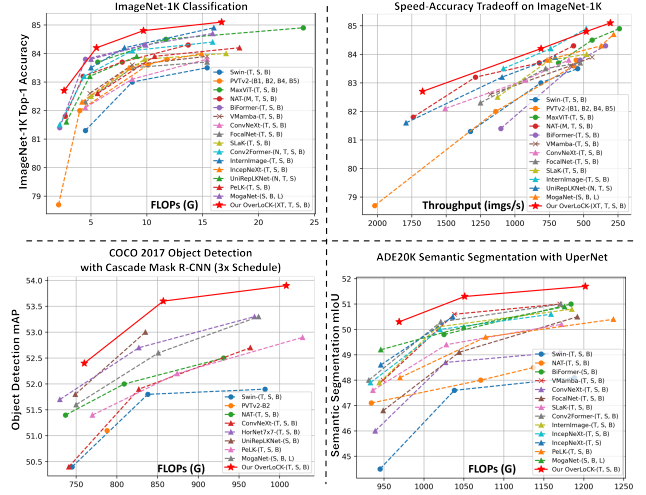


Figure 1. Performance comparisons between our OverLoCK and other representative backbone networks on vision tasks.

form an overall high-level perception, which goes back to fuse with the sensory input, enabling the brain to make more accurate judgments, such as object locations, shapes, and categories. Many previous works have incorporated such top-down attention into vision models, but some of them are unsuitable for building modern vision backbones due to incompatible model designs [1, 6, 28, 53, 72] while the remaining methods primarily focus on recurrent architectures [3, 4, 54, 60, 77], which introduce additional computational overhead due to recurrent operations, resulting in a suboptimal trade-off between performance and computational complexity.

A key property of the top-down attention mechanism is the use of feedback signals as explicit guidance to locate meaningful regions in a scene [58]. However, the classic hierarchical architecture employed in most existing vision backbones [23, 43–45, 66, 67] contrasts with this biological mechanism, as it progressively encodes features from lower to higher levels so that the input features of a layer rely solely on features from previous layers. Hence, there is a lack of explicit top-down semantic guidance in

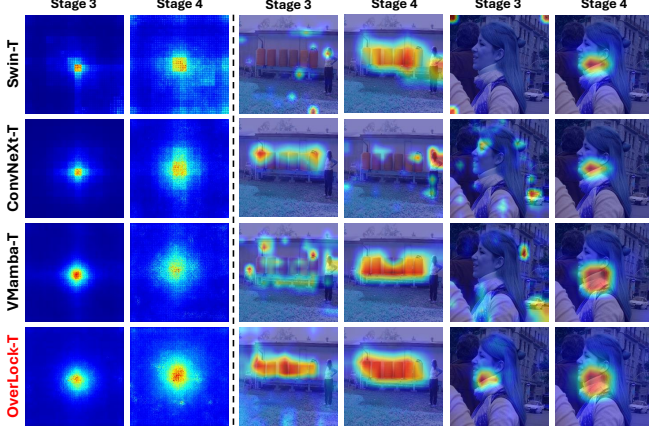


Figure 2. (a) Comparison of Effective Receptive Fields (ERF) [49] at the last layer of deep stages (i.e., Stages 3 and 4) among backbone networks. The results are obtained by averaging over 300 images from ImageNet-1K validation set. As shown, despite being a pure ConvNet, OverLoCK-T has a larger ERF than VMamba-T that emphasizes global modeling, in both Stages 3 and 4. (b) Visualizations of class activation maps computed using Grad-CAM [59] for the output of deep stages (i.e., Stages 3 and 4). The category labels of these two images are “Barrel” and “Neck Brace”. The results demonstrate that although classic hierarchical models can capture long-range dependencies to varying degrees, they struggle to localize objects with the correct category label, especially in Stage 3, which is farther from the classifier. In contrast, our proposed new network architecture can produce more accurate class activation maps in both Stages 3 and 4.

the operations at intermediate layers. To investigate this, we visualize the class activation maps [59] and the effective receptive fields (ERFs) [49] of three representative hierarchical vision models. Swin-T [44], ConvNeXt-T [45], and VMamba-T [43]. As shown in Figure 2, these image classification models struggle to accurately localize objects with the correct category label in the feature maps, especially in Stage 3, which is farther from the classifier layer, despite capturing long-range dependencies in varying degrees. Therefore, *how to develop a modern ConvNet that leverages the top-down attention mechanism while achieving an excellent performance-complexity trade-off remains an open problem.*

On the basis of the above discussion, we propose a biomimetic Deep-stage Decomposition Strategy (DDS) inspired by the top-down attention mechanism in the human vision system. Unlike previous works, our goal is to enhance both feature maps and kernel weights in ConvNets with guidance from dynamic top-down semantic contexts. As illustrated in Figure 3, DDS decomposes a network into three sub-networks: Base-Net, Overview-Net, and Focus-Net. Specifically, a Base-Net encodes low-level and mid-level information, the output is fed into a lightweight Overview-Net to rapidly gather a semantically meaningful but low-quality context representation, analogous to the “overview” process in visual perception. Subsequently, we

designate the output of the Overview-Net as a context prior, which, along with the output of the Base-Net, is fed into a deeper and more powerful Focus-Net to obtain more accurate and informative high-level representations, akin to the “look closely” process in visual perception.

As a top-down context contains information across the entire input image, to fully unleash its power and absorb its information into convolution kernels, the Focus-Net should utilize a powerful dynamic convolution as the token mixer, capable of adaptively modeling long-range dependencies to produce large receptive fields while preserving local inductive biases to capture nuanced local details. Nonetheless, we find that existing convolutions cannot meet these requirements simultaneously. Unlike self-attention mechanisms [14, 64, 66, 67, 81] and State Space Models [16, 18, 43, 47, 82] that can adaptively model long-range dependencies at various input resolutions, large kernel convolutions [12, 13, 42, 45, 75] and dynamic convolutions [8, 36, 37] are still confined to finite regions due to fixed kernel sizes even when input images have an increasingly large resolution, suggesting a weak long-range modeling ability. Although deformable convolutions [10, 69] can alleviate these issues to a certain extent, a deformable kernel shape sacrifices the inherent inductive bias of convolutions, giving rise to a relatively weak local perception ability. Hence, *enabling pure convolutions to possess dynamic global modeling capabilities comparable to those of Transformer- and Mamba-based models while preserving strong inductive biases remains a challenge.*

To tackle this problem, we introduce a novel **Context-Mixing** Dynamic Convolution (ContMix) that dynamically models long-range dependencies while maintaining strong inductive biases. Specifically, for every token in the input feature map, we compute its affinity with respect to a set of region centers across the top-down context feature map, yielding an affinity map. Subsequently, we utilize a learnable linear layer to transform every row of the affinity map, generating spatially varying dynamic convolution kernels. In this regard, every kernel weight carries global information from the top-down semantic context. Consequently, during convolution operations using our dynamic convolution kernels, each token interacts with the global information encoded in the kernels, thereby capturing long-range dependencies despite the fixed size of convolution kernels.

Equipped with the proposed DDS and ContMix, we propose a novel **Overview-first-Look-Closely-next** ConvNet with context-mixing dynamic Kernels (OverLoCK). As shown in Figure 1, our OverLoCK demonstrates superior performance in comparison to representative ConvNet-, Transformer-, and Mamba-based models while striking an excellent balance between speed and accuracy. For example, on the ImageNet-1K dataset, OverLoCK-T achieves a Top-1 accuracy of 84.2%, outperforming UniRepLKNet-

T [13] by 1% and surpassing VMamba-T [43] by 1.6%. On downstream tasks, OverLoCK also demonstrates leading performance. For example, OverLoCK-S outperforms MogaNet-B [39] by 1.2% in mIoU on semantic segmentation and surpasses PeLK-S [7] by 1.4% in AP^b on object detection. Additionally, our method is capable of generating a larger ERF with a strong local inductive bias and more reasonable feature responses in comparison to other competitors, as shown in Figure 2.

2. Related Work

Evolution of ConvNets. Since the debut of AlexNet [33], ConvNets gradually became the dominant architecture in computer vision. VGGNet [61] introduced the concept of stacking small kernels to build deep networks. ResNet [23] and DenseNet [30] further proposed skip-connections to address the gradient vanishing/exploding issues in deep networks. However, with the rise of Vision Transformers [14, 44, 48, 63, 66, 81], ConvNets’ dominance in vision tasks has been challenged. Hence, recent methods have proposed increasingly larger kernel sizes to mimic the self-attention mechanism [14] and establish long-range dependencies [7, 12, 13, 42, 45, 70, 75]. ConvNeXt [45] pioneered the use of 7×7 kernels to build vision backbones, surpassing the performance of Swin Transformer [44]. RepLKNet [12] further explored the promising performance of very large kernels by using 31×31 kernels. On the other hand, gated mechanisms have been extensively explored in ConvNets [39, 50, 51, 56, 74]. For instance, MogaNet [39] introduced a multi-order gated aggregation module to enhance the capacity for refining multi-scale feature representations. StarNet [50] unveiled the underlying reason for the superior performance of element-wise multiplications in the gated mechanism. More recently, RDNet [32] rethought the design of DenseNet and proposed an efficient densely-connected ConvNet. Unlike previous work, this paper focuses on improving the performance of ConvNets from both the architectural and mixer perspectives.

Dynamic Convolutions. Dynamic convolution has been demonstrated to be effective in improving the performance of ConvNets [8, 22, 36, 73] by enhancing feature representation with input-dependent filters. Beyond regular channel-varying modeling, some methods [20, 37, 55, 76] have also proposed spatially varying modeling, which can generate distinct convolution weights for individual pixels in a feature map. Moreover, to enable both the weights and shape of convolution kernels to change dynamically, InternImage [69] re-designed deformable convolutions [10], achieving notable performance gains. However, previous works have failed to simultaneously model long-range dependencies while preserving a strong local inductive bias, a limitation that our new dynamic convolution effectively addresses.

Biomimetic Vision Models. The human vision system has inspired the design of many excellent vision backbone networks. For instance, several advanced vision backbones [7, 52, 74] have been inspired by the peripheral perception mechanism [35], achieving notable performance. Likewise, the top-down attention mechanism [17, 40, 58] has promoted developments in computer vision and machine learning, such as enhancing performance on specific tasks [1, 6, 28, 53, 72], exploring new learning algorithms [78], and designing generic architectures with a recurrent style [3, 4, 54, 77]. Recently, AbsViT [60] introduced a feedback-based Vision Transformer backbone that reuses network outputs to recalibrate early features. In contrast to the above works, we propose a novel modern ConvNet-based vision backbone network that can efficiently generate and utilize top-down guidance, achieving significant performance gains across diverse vision tasks.

3. Methodology

3.1. Deep-stage Decomposition

Overview. Driven by the “Overview-first-Look-Closely-next” mechanism of the human vision system [17, 40], we propose a deep-stage decomposition strategy (DDS) which, in contrast to classic hierarchical architectures, decomposes the network into three distinct sub-networks: Base-Net, Overview-Net, and Focus-Net. As shown in Figure 3, the Base-Net produces a mid-level feature map by progressively downsampling an input image to $\frac{H}{16} \times \frac{W}{16}$ via three embedding layers. This mid-level feature map is fed into both the lightweight Overview-Net as well as the deeper and more powerful Focus-Net. The Overview-Net quickly produces a semantically meaningful but low-quality overview feature map, that serves as an overall understanding of the input image, by immediately downsampling the mid-level feature map to $\frac{H}{32} \times \frac{W}{32}$. This overview feature map is in fact used as a feedback signal fused into all building blocks of the Focus-Net to provide overall contextual information. Thus it is called the *context prior*. Finally, guided by the *context prior*, the Focus-Net goes back to progressively refine the mid-level feature map while enlarging the receptive field to obtain more accurate and informative high-level representations. Note that two backbone networks actually “live” in the above design, one created by cascading Base-Net and Overview-Net and the other by cascading Base-Net and Focus-Net. Each backbone consists of four stages defined by four embedding layers and their following network building blocks. Our DDS design minimizes the overhead by having one Base-Net “serving” mid-level features for both backbones.

During pre-training on ImageNet-1K, to achieve representation learning in both Focus-Net and Overview-Net, each of them is connected to its own classifier head and

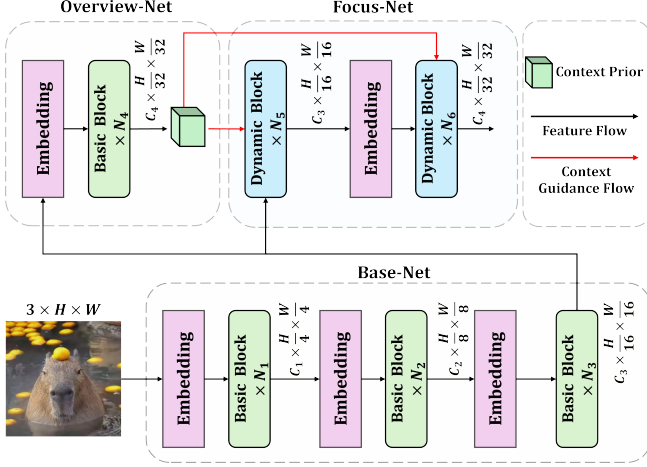


Figure 3. The architecture of our OverLoCK network.

the same classification loss is imposed on both classifiers. When the pre-trained network is transferred to downstream tasks, we no longer apply auxiliary supervision signals to Overview-Net as it has already learned high-level representations during the pre-training stage. Besides, applying auxiliary supervision in dense prediction tasks makes the training process time-consuming. Focus-Net is always used to make predictions in classification tasks. In dense prediction tasks, we use features from Base-Net at $\frac{H}{4} \times \frac{W}{4}$ and $\frac{H}{8} \times \frac{W}{8}$ resolutions as well as features from Focus-Net at $\frac{H}{16} \times \frac{W}{16}$ and $\frac{H}{32} \times \frac{W}{32}$ resolutions to construct a feature pyramid. These four groups of features also correspond to Stages 1 to 4 of our proposed ConvNet backbone network.

Base-Net and Overview-Net. As shown in Figure 4 (a), we adopt the Basic Block as the building blocks of Base-Net and Overview-Net. The input feature is first fed into a residual 3×3 DWConv to perform local perception. The output is then forwarded to a block consisting of a Layer Normalization [34] layer, a Dilated RepConv layer [13], a SE Layer [31], and a ConvFFN [67].

Focus-Net. As illustrated in Figure 4 (b), Focus-Net employs a more complex building block termed Dynamic Block mainly consisting of a residual 3×3 DWConv, a Gated Dynamic Spatial Aggregator (GDSA), and a ConvFFN. The pipeline of GDSA is given in Figure 4 (c), where it uses the proposed ContMix (Section 3.2) as the core token mixer and additionally introduces a gated mechanism to eliminate contextual noise [18, 39, 50]. Note that the Dynamic Blocks before the embedding layer in Focus-Net belong to Stage 3 of the Base-Net+Focus-Net backbone.

Context Flow. There is a dynamic context flow within Focus-Net. The *context prior* from Overview-Net not only provides guidance at both feature and kernel weight levels within Focus-Net, but also is updated within every block along the forward pass. Let us denote the *context prior* and feature map at the entrance of the i -th block as $\mathbf{P}_i \in \mathbb{R}^{C_p \times H \times W}$ and $\mathbf{Z}_i \in \mathbb{R}^{C_z \times H \times W}$, respectively. \mathbf{P}_i and

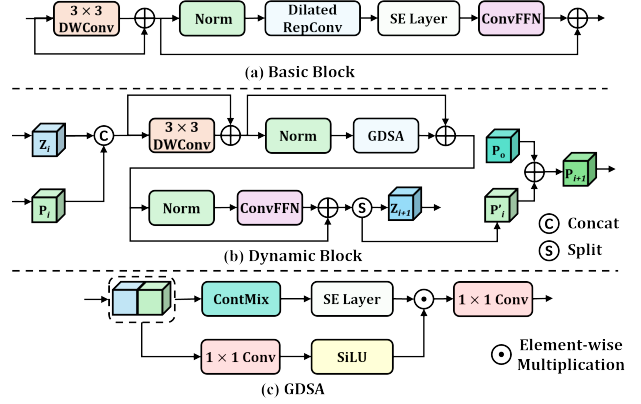


Figure 4. Structures of network building blocks.

\mathbf{Z}_i are fused together through concatenation before being fed into the block (Figure 4 (b)). Inside the block, feature-level guidance is achieved within GDSA by computing a dynamic gate to modulate the feature map using GDSA’s input feature (Figure 4 (c)), which is the result of applying a 1×1 convolution followed by SiLU activation [15] to the aforementioned concatenated feature map. Subsequently, the dynamic gate is element-wise multiplied with the output of its parallel branch. On the other hand, to achieve weight-level guidance, the *context prior* is injected into dynamic convolutions by utilizing \mathbf{P}_i to compute dynamic kernel weights in ContMix, which will be elaborated in the next subsection. Before exiting the block, the fused feature map is split into $\mathbf{P}'_i \in \mathbb{R}^{C_p \times H \times W}$ and $\mathbf{Z}_{i+1} \in \mathbb{R}^{C_z \times H \times W}$, which can be regarded as the disentangled and updated *context prior* and feature map. To prevent the *context prior* from dilution, we add the initial *context prior* \mathbf{P}_o to \mathbf{P}'_i , i.e., $\mathbf{P}_{i+1} = \alpha \mathbf{P}'_i + \beta \mathbf{P}_o$, where α and β are learnable scalars, both initialized to 1 before training.

We perform channel reduction and spatial upsampling on the original *context prior* to save computation and match the input resolution of Focus-Net, respectively. This results in the initial *context prior* \mathbf{P}_o of the context flow.

3.2. Dynamic Convolution with Context-Mixing

In this section, we explore a solution that equips convolutions with the capability of long-range dependency modeling so that they can better handle varying input resolutions. Meanwhile, we still wish them to preserve strong inductive biases. To achieve these goals while fully leveraging the power of the context prior from Overview-Net, we propose a novel dynamic convolution that has a **Context-Mixing** ability, namely ContMix. Our key idea is to represent the relation between a token and its context using the set of affinity values between this individual token and all the tokens at a set of region centers in a feature map. These affinity values can then be aggregated to define token-wise dynamic convolution kernels in a learnable manner, thereby injecting contextual knowledge into every weight of the convolu-

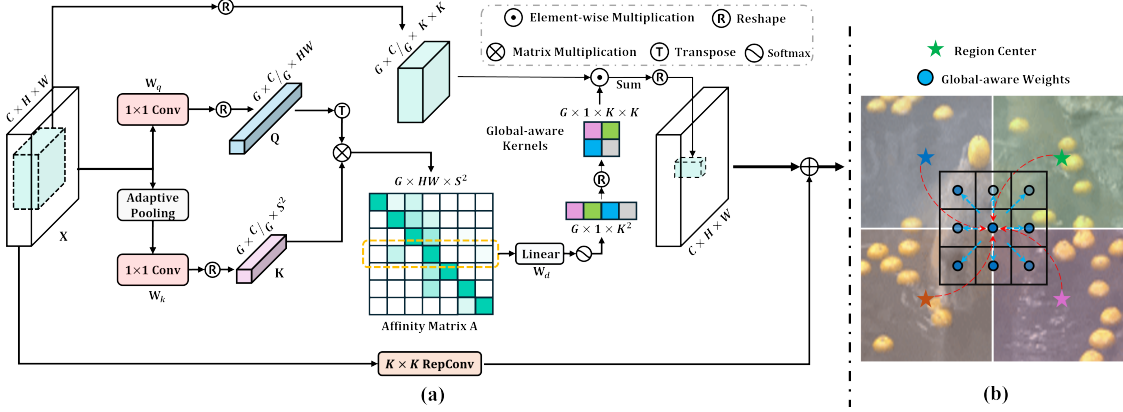


Figure 5. (a) A schematic diagram of our proposed dynamic convolution (ContMix). (b) An illustration of ContMix’s ability in capturing long-range dependencies and preserving inductive biases.

tion kernels. Once such dynamic kernels have been applied to the feature map via sliding windows, every token in the feature map becomes modulated by the approximate global information gathered through the region centers. Thus long-range dependencies can be effectively modeled.

Token-wise Global Context Representation. As shown in Figure 5, given an input feature map $X \in \mathbb{R}^{C \times H \times W}$, we first transform it into two parts, namely, $Q \in \mathbb{R}^{C \times HW} = \text{Re}(W_q X)$ and $K \in \mathbb{R}^{C \times S^2} = \text{Re}(W_k \text{Pool}(X))$, where W_q and W_k denote 1×1 convolutional layers, $\text{Re}(\cdot)$ refers to the reshape operation, K represents the aggregation of X into $S \times S$ region centers via adaptive average pooling. Next, we evenly divide the channels of Q and K into G groups to obtain $\{Q^g\}_{g=1}^G$ and $\{K^g\}_{g=1}^G$ such that $Q^g \in \mathbb{R}^{\frac{C}{G} \times HW}$ and $K^g \in \mathbb{R}^{\frac{C}{G} \times S^2}$. Groups here are analogous to heads in multi-head attention [14]. As every pair of Q^g and K^g have been flattened into 2D matrices, simple matrix multiplications between them computes G affinity matrices $\{A^g\}_{g=1}^G = \{Q^{gT} K^g\}_{g=1}^G$ where $A^g \in \mathbb{R}^{HW \times S^2}$. The i -th row of affinity matrix A^g , A^g_i , holds the affinity values between the i -th token in Q^g and all tokens in K^g .

Token-wise Global Context Mixing. To generate more robust feature representations, we define G spatially varying $K \times K$ dynamic kernels. First, we use another learnable linear layer $W_d \in \mathbb{R}^{S^2 \times K^2}$ to aggregate the token-wise affinity values stored as matrix rows in every affinity matrix A^g by performing a matrix multiplication between A^g and W_d . Note that all G affinity matrices share

the same W_d for saving computational efficiency. Then, a softmax function is employed to normalize the aggregated affinities. These two operations can be formulated as $D^g = \text{softmax}(A^g W_d) \in \mathbb{R}^{HW \times K^2}$. Finally, every row of D^g can be reshaped into the target kernel shape to produce an input-dependent kernel at every token position. During the convolution operation, the channels of feature map X are also evenly divided into G groups, and channels within the same group share the same dynamic kernel.

Implementation. Our ContMix is a general plug-and-play module. In the Dynamic Block of our OverLoCK network, ContMix is customized as follows. The aforementioned Q and K matrices are computed using the channels of X corresponding to Z_i and P_i (the latest *context prior*), respectively. This setting gives rise to better performance in comparison to computing both Q and K using the current fused feature X . In addition, we empirically set S to 7, ensuring that our ContMix enjoys linear-time complexity. Meanwhile, many previous works [12, 13, 42] suggest that combining large and small kernels can lead to better extraction of multi-scale features. Therefore, we allocate half of the groups in ContMix to large kernels and the remaining groups to small kernels, whose size is set to 5×5 following previous works, enabling the modeling of long-range dependencies and local details using different kernels. We also employ a Dilated RepConv layer with $K \times K$ kernels to increase channel diversity.

3.3. Network Architecture

Our OverLoCK network has four architectural variants, including Extreme-Tiny (XT), Tiny (T), Small (S), and Base (B). As listed in Table 1, we control the model size using four variables: *Channels*, *Blocks*, *Kernel Sizes*, and *Groups*. For instance, in OverLoCK-XT, *Channels* = $\{[56, 112, 256], [256], [256, 336]\}$, indicating that the channel counts in the three stages of Base-Net are [56, 112, 256],

Table 1. The configurations of OverLoCK variants.

OverLoCK	Channels	Blocks	Kernel Sizes	Groups
XT	$\{[56, 112, 256], [256], [256, 336]\}$	$\{[2, 2, 3], [2], [6, 2]\}$	$\{[17, 15, 13], [7], [13, 7]\}$	[4, 6]
T	$\{[64, 128, 256], [512], [256, 512]\}$	$\{[4, 4, 6], [2], [12, 2]\}$	$\{[17, 15, 13], [7], [13, 7]\}$	[4, 8]
S	$\{[64, 128, 320], [512], [320, 512]\}$	$\{[6, 6, 8], [3], [16, 3]\}$	$\{[17, 15, 13], [7], [13, 7]\}$	[5, 8]
B	$\{[80, 160, 384], [576], [384, 576]\}$	$\{[8, 8, 10], [4], [20, 4]\}$	$\{[17, 15, 13], [7], [13, 7]\}$	[6, 9]

Table 2. A comparison of image classification performance on ImageNet-1K with Table 3. A comparison of object detection and instance segmentation performance on the COCO dataset using Mask R-CNN. FLOPs are calculated for the 800×1280 resolution.

Method	# T	# F (G)	# P (M)	Acc. (%)
PVTv2-B1[67]	T	2.1	14	78.7
QuadTree-B-b1[62]	T	2.3	14	80.0
RegionViT-T[5]	T	2.4	14	80.4
UniFormer-XS[38]	H	2.0	17	82.0
CrossFormer-T[68]	T	2.9	28	81.5
BiFormer-T[81]	T	2.2	13	81.4
NAT-M[20]	T	2.7	20	81.8
GCViT-Xt[21]	T	2.6	20	82.0
ConvNeXt-N[45]	C	2.7	16	80.9
VAN-B1[19]	C	2.5	14	81.1
Conv2Former-N[27]	C	2.2	15	81.5
UniRepLkNet-N[13]	C	2.8	18	81.6
OverLoCK-Xt	C	2.6	16	82.7

Method	# T	# F (G)	# P (M)	Acc. (%)
Swin-S[44]	T	8.7	50	83.0
PVTv2-B4[67]	T	10.1	63	83.6
UniFormer-B[38]	H	8.3	50	83.9
MaxViT-S[64]	T	11.7	69	84.5
NAT-S[20]	T	7.8	51	83.7
BiFormer-B[81]	T	9.8	57	84.3
VMamba-S[43]	M	8.7	50	83.6
ConvNeXt-S[45]	C	8.7	50	83.1
FocalNet-S[74]	C	8.7	50	83.5
SLaK-S[42]	C	9.8	55	83.8
RDNet-S[32]	C	8.7	50	83.7
InternImage-S[69]	C	8.0	50	84.2
InceptionNeXt-S[75]	C	8.4	49	83.5
PeLk-S[7]	C	10.7	50	83.9
UniRepLkNet-S[13]	C	9.1	56	83.9
MogaNet-B[39]	C	9.9	44	84.3
OverLoCK-S	C	9.7	56	84.8

Method	# T	# F (G)	# P (M)	Acc. (%)
Swin-B[44]	T	15.4	88	83.5
PVTv2-B5[67]	T	11.8	82	83.8
NAT-B[20]	T	13.7	90	84.3
MaxViT-B[64]	T	24.0	120	84.9
VMamba-B[43]	M	15.4	89	83.9
ConvNeXt-B[45]	C	15.4	89	83.8
FocalNet-B[74]	C	15.4	89	83.7
SLaK-B[42]	C	17.1	95	84.0
RDNet-B[32]	C	15.4	87	84.4
InternImage-B[69]	C	16.0	97	84.9
InceptionNeXt-B[75]	C	14.9	87	84.0
PeLk-B[7]	C	18.3	89	84.2
MogaNet-L[39]	C	15.9	83	84.7
OverLoCK-B	C	16.7	95	85.1

Backbone	# F (G)	# P (M)	1 × Schedule	3 × Schedule
			AP^b AP^m	AP^b AP^m
Swin-T[44]	267	48	42.7 39.3	46.0 41.6
PVTv2-B2[67]	309	45	45.3 41.2	47.8 43.1
UniFormer-S[38]	269	41	45.6 41.6	48.2 43.4
NAT-T[20]	258	48	- -	47.8 42.6
BiFormer-S[81]	295	46	47.8 43.2	- -
VMamba-T[43]	271	50	47.3 42.7	48.8 43.7
ConvNeXt-T[45]	262	48	44.2 40.1	46.2 41.7
FocalNet-T[74]	268	49	46.1 41.5	48.0 42.9
InternImage-T[69]	270	49	47.2 42.5	49.1 43.7
RDNet-T[32]	278	43	- -	47.3 42.2
MogaNet-S[39]	272	45	46.7 42.2	48.5 43.1
OverLoCK-T	281	52	48.3 43.3	49.6 43.9
Swin-S[44]	354	69	44.8 40.9	48.2 43.2
PVTv2-B3[67]	397	65	47.0 42.5	48.4 43.2
UniFormer-B[38]	399	69	47.4 43.1	50.3 44.8
NAT-S[20]	330	70	- -	48.4 43.2
BiFormer-B[81]	426	76	48.6 43.7	- -
VMamba-S[43]	384	70	48.7 43.7	49.9 44.2
ConvNeXt-S[45]	348	70	45.4 41.8	47.9 42.9
FocalNet-S[74]	365	72	48.3 43.1	49.3 43.8
InternImage-S[69]	340	69	47.8 43.3	49.7 44.5
MogaNet-B[39]	373	63	47.9 43.2	50.3 44.4
OverLoCK-S	366	75	49.4 44.0	51.0 45.0
Swin-B[44]	496	107	46.9 42.3	48.6 43.3
PVTv2-B5[67]	557	102	47.4 42.5	48.4 42.9
VMamba-B[43]	485	108	49.2 43.9	- -
ConvNeXt-B[45]	486	108	47.0 42.7	48.5 43.5
FocalNet-B[74]	507	111	49.0 43.5	49.8 44.1
InternImage-B[69]	501	115	48.8 44.0	50.3 44.8
MogaNet-L[39]	495	102	49.4 44.1	50.5 44.5
OverLoCK-B	511	114	49.9 44.4	51.4 45.3

the channel count in Overview-Net is 256, and the channel counts in the two stages of Focus-Net are [256, 336]. *Blocks* and *Kernel Sizes* are similarly defined. Additionally, *Groups*=[4, 6] indicates that the number of groups in the dynamic kernels of ContMix in the two stages of Focus-Net is 4 and 6, respectively.

4. Experiments

In this section, we present comprehensive experimental evaluations on various vision tasks, commencing with image classification. Then, we transfer the pre-trained models to downstream tasks, including object detection and semantic segmentation. Due to space constraints, we only report a subset of the results in this section, with additional experimental results provided in the [Appendix](#).

4.1. Image Classification

Setup. We conduct experiments on the ImageNet-1K dataset [11] and adhere to the same experimental setting described in DeiT [63] to ensure a fair comparison. Specifically, all models are trained for 300 epochs using the AdamW optimizer [46]. The stochastic depth rate [29] is set to 0.1, 0.15, 0.4, and 0.5 for OverLoCK-Xt, -T, -S, and -B models, respectively. All experiments are conducted on 8 NVIDIA H800 GPUs.

Results. As shown in Table 2, our pure ConvNet model achieves notable performance improvements over other competitors. For instance, OverLoCK-Xt surpasses a

strong Transformer-based model (BiFormer-T [81]) and a recent large kernel ConvNet (UniRepLkNet-N [13]) by significant 1.3% and 1.1% in Top-1 accuracy, respectively. For Tiny models, our OverLoCK-T also attains the best performance compared with other methods, achieving 84.2% Top-1 accuracy, which improves upon MogaNet-S [39] and PeLk-T [7] by 0.8% and 1.6% in Top-1 accuracy, respectively. When scaling up to larger models, our OverLoCK still maintains a significant advantage. Specifically, OverLoCK-S improves upon BiFormer-B and UniRepLkNet-S by notable 0.5% and 0.9% in Top-1 accuracy, respectively, with comparable computational complexity. Regarding the largest model, OverLoCK-B achieves an impressive 85.1% Top-1 accuracy, outperforming MaxViT-B by 0.2% in Top-1 accuracy with significantly lower computational complexity. Meanwhile, we evaluate the throughput of different models using a batch size of 128 on a single NVIDIA L40S GPU. Figure 1 demonstrates that our OverLoCK achieves an excellent tradeoff between speed and accuracy. For instance, OverLoCK-S surpasses MogaNet-B with over 100 imgs/s in throughput, while significantly increasing Top-1 accuracy from 84.3% to 84.8%. Similarly, OverLoCK-Xt exceeds BiFormer-T by over 600 imgs/s in throughput, while remarkably improving Top-1 accuracy by 1.3%. Overall, to the best of our knowledge, OverLoCK is the first pure ConvNet model to achieve such substantial performance gains over strong baselines on ImageNet-1K.

Table 4. A comparison of object detection and instance segmentation performance on the COCO dataset using Cascade Mask R-CNN. FLOPs are calculated for the 800×1280 resolution.

Backbone	# F (G)	# P (M)	3× Schedule	
			AP^b	AP^m
Swin-T[44]	745	86	50.4	43.7
PVTv2-B2[67]	788	83	51.1	-
NAT-T[20]	737	85	51.4	44.5
ConvNeXt-T[45]	741	86	50.4	43.7
HorNet-T[56]	730	80	51.7	44.8
RDNet-T[32]	757	81	51.6	44.6
PeLK-T[7]	770	86	51.4	44.6
UniRepLkNet-T[13]	749	89	51.8	44.9
MogaNet-S[39]	750	78	51.6	45.1
OverLoCK-T	760	90	52.4	45.4
Swin-S[44]	838	107	51.8	44.7
NAT-S[20]	809	108	52.0	44.9
ConvNeXt-S[45]	827	108	51.9	45.0
HorNet-S[56]	827	108	52.7	45.6
RDNet-S[32]	832	108	52.3	45.3
PeLK-S[7]	874	108	52.2	45.3
UniRepLkNet-S[13]	835	113	53.0	45.9
MogaNet-B[39]	851	101	52.6	46.0
OverLoCK-S	857	114	53.6	46.4
Swin-B[44]	982	145	51.9	45.0
NAT-B[20]	931	147	52.5	45.2
ConvNeXt-B[45]	964	146	52.7	45.6
HorNet-B[56]	969	144	53.3	46.1
RDNet-B[32]	971	144	52.3	45.3
PeLK-B[7]	1028	147	52.9	45.9
MogaNet-L[39]	974	149	53.3	46.1
OverLoCK-B	1008	154	53.9	46.8

4.2. Object Detection and Instance Segmentation

Setup. We evaluate our network architecture on object detection and instance segmentation tasks using the COCO 2017 dataset [41]. We employ both Mask R-CNN [24] and Cascade Mask R-CNN [2] frameworks, adopting the same experimental setting in Swin [44]. The backbone networks are initially pre-trained on ImageNet-1K and subsequently fine-tuned for 12 epochs (1× schedule) and 36 epochs (3× schedule with multi-scale training).

Results. As shown in Tables 3 and 4, OverLoCK demonstrates notable superiority over other methods. For instance, using the Mask R-CNN 1× schedule, OverLoCK-S surpasses BiFormer-B and MogaNet-B by 0.8% and 1.5% in AP^b , respectively. When using Cascade Mask R-CNN, OverLoCK-S improves upon PeLK-S and UniRepLkNet-S by 1.4% and 0.6% in AP^b , respectively. Notably, we observe an interesting phenomenon: *although ConvNet-based methods achieve comparable performance with Transformer-based methods on image classification tasks, there is a significant performance gap on detection tasks.* For example, both MogaNet-B and BiFormer-B achieve Top-1 accuracy of 84.3% on ImageNet-1K, but the former lags behind the latter on detection tasks. This validates our previous argument that ConvNet’s fixed kernel size leads to limited receptive fields, resulting in performance degrada-

Table 5. Comparison of semantic segmentation performance on the ADE20K dataset. FLOPs are calculated for the 512×2048 resolution.

Backbone	UperNet 160K		
	F (G)	P (M)	mIoU
Swin-T[44]	945	60	44.5/45.8
UniFormer-S[38]	1008	52	47.6/48.5
NAT-T[20]	934	58	47.1/48.4
BiFormer-S[81]	1025	55	49.8/50.8
VMamba-T[43]	949	62	48.0/48.8
ConvNeXt-T[45]	939	60	46.0/46.7
FocalNet-T[74]	949	61	46.8/47.8
SLaK-T[42]	936	65	47.6/—
RDNet-T[32]	961	58	47.6/48.6
InternImage-T[69]	944	59	47.9/48.1
InceptionNeXt-T[75]	933	56	47.9/—
PeLK-T[7]	970	62	48.1/—
UniRepLkNet-T[13]	946	61	48.6/49.1
MogaNet-S[39]	946	55	49.2/—
OverLoCK-T	969	63	50.3/50.8
Swin-S[44]	1038	81	47.6/49.5
UniFormer-B[38]	1227	80	50.0/50.8
NAT-S[20]	1071	82	48.0/49.5
BiFormer-B[81]	1184	88	51.0/51.7
VMamba-S[43]	1038	82	50.6/51.2
ConvNeXt-S[45]	1027	82	48.7/49.6
FocalNet-S[74]	1044	84	49.1/50.1
SLaK-S[42]	1028	91	49.4/—
RDNet-S[32]	1040	86	48.7/49.8
InternImage-S[69]	1017	80	50.1/50.9
InceptionNeXt-S[75]	1020	78	50.0/—
PeLK-S[7]	1077	84	49.7/—
UniRepLkNet-S[13]	1036	86	50.5/51.0
MogaNet-B[39]	1050	74	50.1/—
OverLoCK-S	1051	85	51.3/51.9
Swin-B[44]	1188	121	48.1/49.7
NAT-B[20]	1137	123	48.5/49.7
FocalNet-B[74]	1192	126	50.5/51.4
SLaK-B[42]	1172	135	50.2/—
RDNet-B[32]	1187	127	49.6/50.5
InternImage-B[75]	1185	128	50.8/51.3
InceptionNeXt-B[75]	1159	115	50.6/—
PeLK-B[7]	1237	126	50.4/—
MogaNet-L[39]	1176	113	50.9/—
OverLoCK-B	1202	124	51.7/52.3

Table 6. A comprehensive roadmap that incrementally evolves a simple baseline to our OverLoCK-XT model. †: The auxiliary loss for Overview-Net is only used in the classification task.

Method	# F (G)	# P (M)	Top-1 (%)	mIoU (%)
PlainNet	2.5	15.7	76.3	38.8
w/ Dilated RepConv	2.5	15.7	76.6	39.3
w/ SE	2.5	16.0	77.1	39.6
w/ Local Conv	2.5	16.1	78.0	40.2
FFN→ConvFFN	2.6	16.3	78.5	41.1
Recurrent Model	4.4	16.4	76.8	39.5
DDS Model	2.6	15.6	79.0	41.6
w/o feature feed	2.7	16.3	78.2	40.0
Static→Dynamic	2.6	15.8	80.0	42.9
w/ Aux Loss†	2.6	15.8	80.2	43.1
w/ Initial Prior	2.6	15.8	80.4	43.4
w/ Gate (OverLoCK-XT)	2.6	16.4	80.8	43.8

Table 7. A comparison of dynamic token mixers

Method	# F (G)	# P (M)	Top-1 (%)	mIoU (%)
DyConv[8]	1.7	12.5	77.5	36.3
ODConv[36]	1.7	13.8	77.9	36.7
Involution[37]	1.9	12.2	78.3	37.5
VOLO[76]	2.3	14.1	75.5	35.6
DCNv3[69]	2.6	16.3	78.7	37.5
Shifted Window[44]	2.2	13.6	78.4	37.4
Natten[20]	2.2	13.6	78.6	38.1
ContMix (Ours)	2.2	13.1	78.8	39.0

tion when using large input resolutions. Conversely, our OverLoCK effectively captures long-range dependencies even at large resolutions, resulting in excellent performance.

4.3. Semantic Segmentation

Setup. We conduct experiments on semantic segmentation using the ADE20K dataset [80] with the UperNet framework [71]. For a fair comparison, we initialize all backbone networks with weights pre-trained on ImageNet-1K, following the same training settings as outlined in Swin [44].

Results. Table 5 demonstrates that our OverLoCK achieves leading performance on semantic segmentation. For instance, OverLoCK-T outperforms MogaNet-S and UniRepLkNet-T by 1.1% and 1.7% in terms of mIoU, respectively, and surpasses VMamba-T, which emphasizes global modeling, by 2.3% in mIoU. This advantage is consistently observed in both Small and Base models. Moreover, *we find that the issue of limited receptive fields also negatively impacts the performance of ConvNets on segmentation tasks*, e.g., MogaNet-B lags behind BiFormer-B by 0.9% despite having the same accuracy on classification. In contrast, our OverLoCK effectively alleviates this issue.

4.4. Ablation Studies

Setup. We conduct comprehensive ablation studies on image classification and semantic segmentation tasks to eval-

uate the effectiveness of individual components in OverLoCK. Specifically, we train each model variant on the ImageNet-1K dataset for 120 epochs following [7, 42] while keeping the remaining training settings consistent with those described in Section 4.1. Subsequently, we fine-tune the pre-trained models on the ADE20K dataset for 80K iteration steps, with a batch size of 32 for faster training, while keeping the remaining settings identical to those outlined in Section 4.3. Due to the page limit, more ablation studies are presented in the [Appendix](#).

A detailed roadmap to our OverLoCK model. First, we aim to develop a powerful baseline model using static large kernel convolutions. To this end, we first evaluate the performance of different components in the Basic Block (Figure 4 (a)). Specifically, we construct a hierarchical model using a vanilla convolutional layer followed by a vanilla FFN [14] as the building block. The model consists of four stages with the number of blocks in every stage set to [2, 2, 9, 4] and the number of channels in every stage set to [56, 112, 304, 400]. The kernel sizes of the four stages are consistent with the XT model. This model is denoted as “PlainNet” and achieves a Top-1/mIoU of 76.3%/38.8%, as listed in Table 6. Then, we convert the vanilla convolutional layer to the Dilated RepConv layer [13], referring to “w/ Dilated RepConv” (Top-1/mIoU: 76.6%/39.3%). Next, we incrementally add an SE Layer (Top-1/mIoU: 77.1%/39.6%), a 3×3 DWConv (Top-1/mIoU: 78.0%/40.2%), and replace the vanilla FFN with ConvFFN (Top-1/mIoU: 78.5%/41.1%). The resulting network is termed “Baseline”.

Subsequently, we explore three strategies to inject top-down attention into this Baseline network. (1) Inspired by AbsViT [60], we construct a recurrent model by upsampling the output of Stage 4 and concatenating it with the input of Stage 3, termed “Recurrent Model”. However, this model’s performance drops to 76.8%/39.5% with higher complexity, indicating that recurrent designs are unsuitable for modern ConvNet-based backbones. (2) We employ our proposed DDS to decompose the Baseline network into three inter-connected sub-networks. The outputs from both BaseNet and Overview-Net are concatenated and fed into FocusNet. To ensure similar computational costs with the Baseline model, the numbers of blocks and channels in the three sub-networks are set to [56, 112, 304], [400], [304, 400] and [2, 2, 3], [2], [6, 2], respectively. This model, denoted as “DDS Model”, achieves a Top-1 accuracy of 79.0%/41.6%. (3) In the “DDS Model”, we only feed the projected output of Overview-Net into FocusNet, without concatenating it with the output of BaseNet, the resulting model is termed “w/o feature feed”. This model decreases the performance, demonstrating the importance of feeding the output from BaseNet into FocusNet.

Finally, we evaluate the impact of weight-level context

guidance by replacing every existing block in the FocusNet with the proposed Dynamic Block, excluding the gate module, and ensuring that the context prior update flow does not use the initial *context prior*. This modification maintains the same numbers of blocks and channels as our XT model, ensuring comparable computational complexity for a fair comparison. This model, termed “Static \rightarrow Dynamic”, notably improves the Top-1/mIoU to 80.0%/42.9%. Next, to make Overview-Net produce semantically meaningful context features, we use an auxiliary classification loss to supervise its output. The resulting model is termed “w/ Aux Loss”, which further improves the Top-1/mIoU by 0.2%/0.2%. Subsequently, we incorporate the initial *context prior* into each dynamic block, as described in Section 5, to prevent the dilution of meaningful information within the *context prior* during the updating process. This variant, labeled as “w/ Initial Prior” enhances the Top-1/mIoU by 0.2%/0.3%. Lastly, we evaluate the impact of context-guided feature modulation by adding the gate module. This results in our XT model, which further boosts the Top-1/mIoU to 80.8%/43.8%. Summarily, our proposed method plays a vital role in significant performance improvements.

A comparison of token mixers. To conduct a fair comparison of dynamic token mixers, we construct a Swin-like architecture [44] by setting the numbers of blocks and channels in the four stages to [2, 2, 6, 2] and [64, 128, 256, 512], respectively, and employing non-overlapping patch embedding and a standard feed-forward network (FFN). We implement DyConv and ODConv in a separable convolution style [9] to ensure comparable computational complexity with other methods. Additionally, we set the kernel/window size to 7×7 for all methods except for VOLO, where larger kernels incur significantly more parameters. From Table 7, our **Context-Mixing** Dynamic Kernel (ContMix) achieves the best result on both image classification and semantic segmentation tasks. Notably, although ContMix exhibits similar performance as Natten and DCNv3 on classification tasks with low-resolution inputs, it demonstrates a clear advantage on semantic segmentation tasks with higher-resolution inputs. This is because ContMix captures long-range dependencies while preserving local inductive biases.

5. Conclusion

This paper proposes a biomimetic Deep-stage Decomposition (DDS) mechanism that injects semantically meaningful contexts into the intermediate stages of the network and also presents a novel dynamic convolution with context-mixing capacity, dubbed ContMix, which captures long-range dependencies while preserving strong inductive biases. By integrating these components, we propose a powerful, pure ConvNet-based vision backbone network, termed OverLoCK, achieving clearly superior performance compared to strong baselines.

OverLoCK: An Overview-first-Look-Closely-next ConvNet with Context-Mixing Dynamic Kernels

Supplementary Material

A. More Ablation Studies

On the basis of the training settings outlined in Section 4.4, we additionally conduct a series of in-depth ablation experiments to meticulously examine the impact of every component in our proposed method.

Impact of Kernel Sizes. We compared the performance under various settings of kernel sizes, as outlined in Table 6 (the definition of the kernel size in our proposed method is given in Section 3.3). The results indicate that the configuration $\{[17, 15, 13], [7], [13, 7]\}$ yields the optimal performance on both image classification and semantic segmentation tasks. Further enlarging the kernels does not lead to additional improvements.

Table A. Ablation study of the kernel size setting.

Kernel Sizes	# F (G)	# P (M)	Top-1 (%)	mIoU (%)
$\{[19, 17, 15], [7], [15, 7]\}$	2.8	16.5	80.7	43.8
$\{[17, 15, 13], [7], [13, 7]\}$	2.6	16.4	80.8	43.8
$\{[13, 11, 9], [7], [9, 7]\}$	2.6	16.3	80.5	43.5
$\{[9, 9, 7], [7], [7, 7]\}$	2.6	16.1	80.6	43.3
$\{[7, 7, 7], [7], [7, 7]\}$	2.5	16.1	80.4	43.1

Impact of Stage Ratio. The *Stage Ratio* means the ratio between the number of blocks in the last stage of Base-Net and the number of blocks in the first stage of Focus-Net. In the default setting of the OverLoCK model, the stage ratio is 1:2 with the intention of allocating more network blocks to Focus-Net for extracting robust contextual information. In this section, we investigate the impact of *Stage Ratio*. Apart from the default setting of 1:2, we further set *Stage Ratio* to 1:1 and 1:3 while maintaining the total number of network blocks constant. The results presented in Table B demonstrate that a *Stage Ratio* of 1:2 yields the best outcomes. We posit that this is because a too small *Stage Ratio* results in insufficient number of blocks in Focus-Net, thereby hindering the extraction of discriminative deep features. Conversely, an excessively large *Stage Ratio* leads to a shortage of blocks in Base-Net, thereby providing insufficient contextual guidance.

Table B. Ablation study of different stage ratio settings.

Stage Ratio	# F (G)	# P (M)	Top-1 (%)	mIoU (%)
1:1	2.7	16.1	80.4	42.9
1:2	2.6	16.4	80.8	43.8
1:3	2.7	15.9	80.6	43.6

Impact of Channel Reduction Factor. In the default configuration of the OverLoCK model, we employ a 1×1 convolution to reduce the number of output channels of Overview-Net by a factor of 4 and concatenate this result with the output of Base-Net before forwarding it to Focus-Net. We term this reduction as the *Channel Reduction Factor (CRF)*. Therefore, the value of *CRF* determines the number of channels in the *context prior*, thereby influencing the guidance capability. In this regard, we investigate the effects of different *CRF* settings. It is important to note that during the adjustment of *CRF*, we also modify the number of channels of Focus-Net to maintain similar complexities across different model variants. The results in Table C demonstrate that *CRF*=4 yields the optimal performance.

Table C. Ablation study of channel reduction factor settings.

CRF	# F (G)	# P (M)	Top-1 (%)	mIoU (%)
2	2.6	16.1	80.5	42.9
4	2.6	16.4	80.8	43.8
6	2.7	16.6	80.7	43.4
8	2.7	16.7	80.6	43.0

Impact of Auxiliary Loss. To explore the effects of applying the auxiliary loss to Overview-Net, we adjust the weight of the auxiliary loss, drawing inspiration from prior research [79]. Given that the architectures of the models in this comparison are consistent, we opt not conduct further experiments on segmentation tasks for the sake of simplicity. The results presented in Table D indicate that the utilization of an auxiliary loss improves accuracy, while varying the weight of the auxiliary loss does not lead to a notable impact on performance. This observation aligns with findings in previous study [79].

Table D. Ablation study of auxiliary loss.

Aux Loss Ratio	0	0.2	0.4	0.8	1.0
Top-1 (%)	80.4	80.7	80.8	80.7	80.7

Effectiveness of our DDS-based Top-down Network.

To evaluate the effectiveness of the proposed DDS, we reconstruct our OverLoCK-XT model as a standard hierarchical network. To be specific, we eliminate the top-down attention mechanism by removing the Overview-Net while keeping the same types of layers in the Base-Net and Focus-Net. To maintain comparable complexity with other mod-

els, the number of channels and layers in the four stages are set to [64, 112, 256, 360] and [2, 2, 9, 4], respectively. This model is denoted as the “Hierarchical Model”. Additionally, we compare it with the “Baseline” model in Table 6 which is a fully static ConvNet. As shown in Table E, the “Hierarchical Model” results in a noticeable performance drop, demonstrating the effectiveness of our DDS-based top-down context guidance. However, when compared with the “Baseline” model, it still exhibits significant advantages, clearly indicating the superiority of our proposed dynamic convolution module.

Table E. Effectiveness of the proposed DDS-based top-down network.

Method	# F (G)	# P (M)	Top-1 (%)	mIoU (%)
Baseline Model	2.6	16.3	78.5	41.1
Hierarchical Model	2.7	16.2	79.2	41.9
OverLoCK-XT	2.6	16.4	80.7	43.8

Ablation Study of ContMix. We conduct a comprehensive comparison of various components within our proposed ContMix framework, as presented in Section 3.2. As listed in Table F, we initially compute \mathbf{Q} and \mathbf{K} using the fused feature map instead of utilizing the channels of \mathbf{X} corresponding to \mathbf{Z}_i and \mathbf{P}_i (the latest *context prior*). This model variant, referred to as “Fusion Affinity”, results in a marginal performance decline. Subsequently, we interchange the features used to generate the \mathbf{Q} and \mathbf{K} matrices. This model, denoted as “Reverse QK”, also exhibits a decrease in performance. Furthermore, we individually eliminate the Softmax function (referred to as “w/o Softmax”), remove the RepConv (referred to as “w/o RepConv”), and substitute small kernels with large kernels (referred to as “w/o Small Kernel”). These alterations decrease performance on both classification and segmentation tasks.

Table F. Ablation study of ContMix.

Method	# F (G)	# P (M)	Top-1 (%)	mIoU (%)
Baseline	2.6	16.4	80.8	43.8
Fusion Affinity	2.7	16.6	80.7	43.5
Reverse QK	2.7	16.4	80.6	42.9
w/o Softmax	2.6	16.4	80.5	43.5
w/o RepConv	2.5	16.1	80.6	43.4
w/o Small Kernel	2.8	16.6	80.7	43.3

B. Additional Experiments on Image Classification

B.1. Large Resolution Evaluation

Following previous works [32, 45, 75], we further investigate the image classification performance on the ImageNet-1K dataset at a higher resolution (i.e., 384×384). Specifi-

Table G. A comparison of image classification with 384×384 inputs.

Method	Type	# F (G)	# P (M)	Acc. (%)
Swin-B	T	47.1	88	84.5
MaxViT-B	T	74.2	120	85.7
ConvNeXt-B	C	45.2	88	85.1
InceptionNeXt-B	C	43.6	87	85.2
RDNet-L	C	101.9	186	85.8
PeLK-B-101	C	68.3	90	85.8
OverLoCK-B	C	50.4	95	86.2

Table H. Robustness comparisons of different models.

Models	# F (G)	# P (G)	1K	V2	A	R	Sketch
Swin-T	4.5	28	81.3	69.7	21.1	41.5	29.3
VMamba-T	4.9	29	82.6	72.0	27.0	45.4	32.9
ConvNeXt-T	4.5	29	82.1	72.5	24.2	47.2	33.8
HorNet-T	4.0	22	82.8	72.3	26.6	46.6	34.1
SLaK-T	5.0	30	82.5	72.0	30.0	45.3	32.4
NAT-T	4.3	28	83.2	72.2	33.0	44.9	31.9
RDNet-T	5.0	24	82.8	72.9	27.7	49.0	37.0
UniRepLKNet-T	4.9	25	83.2	72.8	34.8	49.4	36.9
MogaNet-S	5.0	33	83.4	72.6	33.4	49.7	37.8
OverLoCK-T	5.5	33	84.2	74.0	39.4	53.3	40.6
Swin-S	8.7	50	83.0	72.0	32.5	45.2	32.3
VMamba-S	8.7	50	83.6	73.2	33.2	49.4	37.0
ConvNeXt-S	8.7	50	83.1	72.5	31.3	49.6	37.1
HorNet-S	8.8	50	84.0	73.6	36.2	49.7	36.9
SLaK-S	9.8	55	83.8	73.6	39.3	50.9	37.5
NAT-S	7.8	51	83.7	73.2	37.4	47.3	34.3
RDNet-S	8.7	50	83.7	73.8	33.5	52.8	39.8
UniRepLKNet-S	9.1	56	83.9	73.7	38.3	50.6	36.9
MogaNet-B	9.9	44	84.3	74.3	40.4	50.1	38.6
OverLoCK-S	9.7	56	84.8	74.9	45.0	57.2	45.8
Swin-B	15.4	88	83.5	72.4	35.4	46.5	32.7
VMamba-B	15.4	89	83.9	73.5	37.2	49.5	38.5
ConvNeXt-B	15.4	89	83.8	73.7	36.7	51.2	38.2
HorNet-B	15.6	87	84.3	73.9	39.9	51.2	38.1
SLaK-B	17.1	95	84.0	74.0	41.6	50.8	38.5
NAT-B	13.7	90	84.3	74.1	41.4	49.7	36.6
RDNet-B	15.4	87	84.4	74.2	38.1	52.7	40.1
MogaNet-L	15.9	83	84.7	74.0	41.0	52.2	39.0
OverLoCK-B	16.7	95	85.1	75.4	47.7	58.5	46.0

cally, we pre-train the base model on 224×224 inputs and then fine-tune it on 384×384 inputs for 30 epochs. As shown in Table G, our OverLoCK-B model achieves superior performance under high-resolution input conditions. Notably, OverLoCK-B surpasses MaxViT-B by 0.5% in Top-1 accuracy while reducing the parameter count by over one-third. Compared to PeLK-B, a large kernel ConvNet, our method also demonstrates significant improvements. These results further validate the robustness of our proposed method in handling large-resolution inputs.

B.2. Robustness Evaluation

We further assess the robustness of our models using the ImageNet out-of-distribution (OOD) benchmarks, including ImageNet-V2 [57], ImageNet-A [26], ImageNet-R [25], and ImageNet-Sketch [65]. As shown in Table H, our method demonstrates excellent robustness on differ-

Table I. Speed comparison among various models. Throughput (Thr.) is tested on a single NVIDIA L40S GPU with a batch size of 128 and an image size of $3 \times 224 \times 224$.

Method	# F (G)	# P (M)	Thr. (imgs/s)	Acc. (%)
Swin-T	4.5	28	1324	81.3
Swin-S	8.7	50	812	83.0
Swin-B	15.4	88	544	83.5
MaxViT-T	5.6	31	683	83.7
MaxViT-S	11.7	69	439	84.5
MaxViT-B	24.0	120	241	84.9
NAT-M	2.7	20	1740	81.8
NAT-T	4.3	28	1287	83.2
NAT-S	7.8	51	823	83.7
NAT-B	13.7	90	574	84.3
BiFormer-T	2.2	13	1103	81.4
BiFormer-S	4.5	26	527	83.8
BiFormer-B	9.8	57	341	84.3
VMamba-T	4.9	29	1179	82.6
VMamba-S	8.7	50	596	83.6
VMamba-B	15.4	89	439	83.9
ConvNeXt-T	4.5	29	1507	82.1
ConvNeXt-S	8.7	50	926	83.1
ConvNeXt-B	15.4	89	608	83.8

Method	# F (G)	# P (M)	Thr. (imgs/s)	Acc. (%)
FocalNet-T	4.5	29	1251	82.3
FocalNet-S	8.7	50	777	83.5
FocalNet-B	15.4	89	481	83.7
SLaK-T	5.0	30	1126	82.5
SLaK-S	9.8	55	747	83.8
SLaK-B	17.1	95	478	83.7
InternImage-T	5.0	30	1084	83.5
InternImage-S	8.0	50	740	84.2
InternImage-B	16.0	97	481	84.9
UniRepLKNet-N	2.8	18	1792	81.6
UniRepLKNet-T	4.9	31	1094	83.2
UniRepLKNet-S	9.1	56	707	83.9
MogaNet-S	5.0	25	766	83.4
MogaNet-B	9.9	44	373	84.3
MogaNet-L	15.9	83	282	84.7
OverLoCK-XT	2.6	16	1672	82.7
OverLoCK-T	5.5	33	810	84.2
OverLoCK-S	9.7	56	480	84.8
OverLoCK-B	16.7	95	306	85.1

ent datasets, outperforming representative ConvNets, Vision Transformers, and Vision Mamba. Notably, although OverLoCK-B improves over MogaNet-L by 0.4% in Top-1 accuracy on ImageNet-1K, it achieves significant gains on OOD datasets, with improvements of 1.4% on ImageNet-V2, 6.7% on ImageNet-A, 6.3% on ImageNet-R, and 6.8% on ImageNet-Sketch. These results showcase the strong robustness of our pure ConvNet.

C. Speed Analysis

We provide a comparison of speed-accuracy trade-off in Figure 1. More details are listed in Table 1, where an OverLoCK variant often achieves faster speed and higher accuracy simultaneously than a larger variant of another network, demonstrating an excellent trade-off between speed and accuracy. For instance, OverLoCK-XT achieves 1672 imgs/s in throughput, improving upon Swin-T by over 300 imgs/s, while significantly enhancing Top-1 accuracy by 1.4%. Also, OverLoCK-T achieves about 200 imgs/s improvement in throughput compared to ConvNeXt-B while achieving better performance at the cost of only around one-third of the FLOPS. When compared to more advanced models, OverLoCK still exhibits significant advantages. For example, OverLoCK-S surpasses MogaNet-B by over 100 imgs/s in throughput while increasing Top-1 accuracy from 84.3% to 84.8%. Likewise, OverLoCK-XT surpasses BiFormer-T by over 600 imgs/s in throughput while remarkably improving Top-1 accuracy by 1.3%.

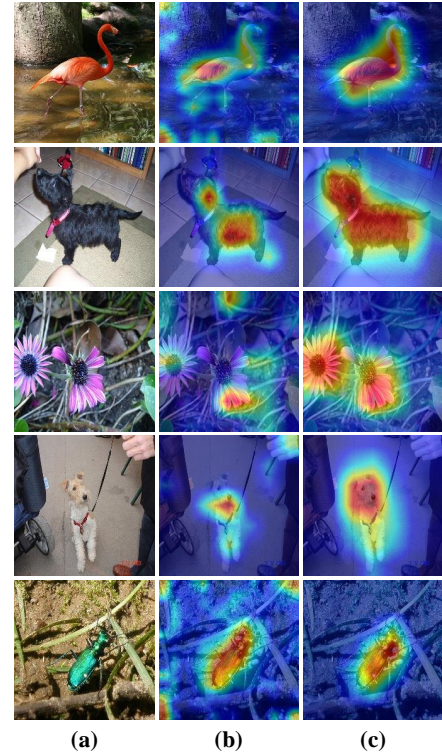


Figure A. Class activation maps of the proposed OverLoCK network. (a), (b), and (c) show the input images, class activation maps of Overview-Net, and class activation maps of Focus-Net, respectively.

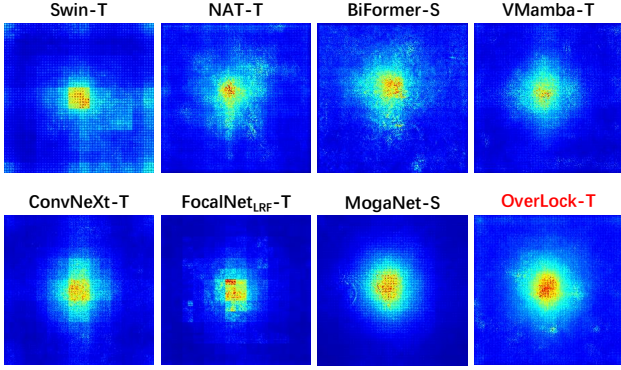


Figure B. Comparison of ERF among various models.

D. Visualization Analysis

D.1. Effect of Context Guidance

To visually understand the effect of context guidance, we separately visualize the class activation maps generated by Overview-Net and Focus-Net in OverLoCK-T using Grad-CAM [59] for the ImageNet-1K validation set. As shown in Figure A, Overview-Net first produces a coarse localization of an object, and when this signal is used as the top-down guidance for Focus-Net, the object’s location and shape becomes more accurate.

D.2. Effective Receptive Field Analysis

To visually demonstrate the representation capacity of OverLoCK, we compare the Effective Receptive Field (ERF) [49] of our OverLoCK-T with that of other representative models with comparable complexity. The visualizations are generated using over 300 randomly sampled images with a resolution of 224×224 from the ImageNet-1K validation set. As shown in Figure B, our model not only produces global responses but also exhibits significant local sensitivity, indicating that OverLoCK can effectively model both global and local contexts simultaneously.

References

- [1] Peter Anderson, Xiaodong He, Chris Buehler, Damien Teney, Mark Johnson, Stephen Gould, and Lei Zhang. Bottom-up and top-down attention for image captioning and visual question answering. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pages 6077–6086, 2018. 1, 3
- [2] Zhaowei Cai and Nuno Vasconcelos. Cascade r-cnn: High quality object detection and instance segmentation. *IEEE Transactions on Pattern Analysis and Machine Intelligence*, 43(5):1483–1498, 2019. 7
- [3] Chunshui Cao, Xianming Liu, Yi Yang, Yinan Yu, Jiang Wang, Zilei Wang, Yongzhen Huang, Liang Wang, Chang Huang, Wei Xu, et al. Look and think twice: Capturing top-down visual attention with feedback convolutional neural networks. In *Proceedings of the IEEE international conference on computer vision*, pages 2956–2964, 2015. 1, 3
- [4] Chunshui Cao, Yongzhen Huang, Yi Yang, Liang Wang, Zilei Wang, and Tieniu Tan. Feedback convolutional neural network for visual localization and segmentation. *IEEE Transactions on Pattern Analysis and Machine Intelligence*, 41(7):1627–1640, 2018. 1, 3
- [5] Chun-Fu Chen, Rameswar Panda, and Quanfu Fan. Regionvit: Regional-to-local attention for vision transformers. In *International Conference on Learning Representations*, 2022. 6
- [6] Hao Chen, Kunyang Sun, Zhi Tian, Chunhua Shen, Yongming Huang, and Youliang Yan. Blendmask: Top-down meets bottom-up for instance segmentation. In *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*, pages 8573–8581, 2020. 1, 3
- [7] Honghao Chen, Xiangxiang Chu, Yongjian Ren, Xin Zhao, and Kaiqi Huang. Pelk: Parameter-efficient large kernel convnets with peripheral convolution. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, 2024. 3, 6, 7, 8
- [8] Yinpeng Chen, Xiyang Dai, Mengchen Liu, Dongdong Chen, Lu Yuan, and Zicheng Liu. Dynamic convolution: Attention over convolution kernels. In *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*, pages 11030–11039, 2020. 2, 3, 7
- [9] François Chollet. Xception: Deep learning with depthwise separable convolutions. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pages 1251–1258, 2017. 8
- [10] Jifeng Dai, Haozhi Qi, Yuwen Xiong, Yi Li, Guodong Zhang, Han Hu, and Yichen Wei. Deformable convolutional networks. In *Proceedings of the IEEE international conference on computer vision*, pages 764–773, 2017. 2, 3
- [11] Jia Deng, Wei Dong, Richard Socher, Li-Jia Li, Kai Li, and Li Fei-Fei. Imagenet: A large-scale hierarchical image database. In *IEEE conference on computer vision and pattern recognition*, pages 248–255, 2009. 6
- [12] Xiaohan Ding, Xiangyu Zhang, Jungong Han, and Guiguang Ding. Scaling up your kernels to 31×31 : Revisiting large kernel design in cnns. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 11963–11975, 2022. 2, 3, 5
- [13] Xiaohan Ding, Yiyuan Zhang, Yixiao Ge, Sijie Zhao, Lin Song, Xiangyu Yue, and Ying Shan. Unireplknet: A universal perception large-kernel convnet for audio, video, point cloud, time-series and image recognition. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, 2024. 2, 3, 4, 5, 6, 7, 8
- [14] Alexey Dosovitskiy, Lucas Beyer, Alexander Kolesnikov, Dirk Weissenborn, Xiaohua Zhai, Thomas Unterthiner, Mostafa Dehghani, Matthias Minderer, Georg Heigold, Sylvain Gelly, et al. An image is worth 16×16 words: Transformers for image recognition at scale. In *International Conference on Learning Representations*, 2021. 2, 3, 5, 8
- [15] Stefan Elfving, Eiji Uchibe, and Kenji Doya. Sigmoid-weighted linear units for neural network function approxima-

- tion in reinforcement learning. *Neural networks*, 107:3–11, 2018. 4
- [16] Yunxiang Fu, Meng Lou, and Yizhou Yu. Segman: Omni-scale context modeling with state space models and local attention for semantic segmentation. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, 2025. 2
- [17] Charles D Gilbert and Mariano Sigman. Brain states: top-down influences in sensory processing. *Neuron*, 54(5):677–696, 2007. 1, 3
- [18] Albert Gu and Tri Dao. Mamba: Linear-time sequence modeling with selective state spaces. In *Conference on Language Modeling*, 2024. 2, 4
- [19] Meng-Hao Guo, Cheng-Ze Lu, Zheng-Ning Liu, Ming-Ming Cheng, and Shi-Min Hu. Visual attention network. *Computational Visual Media*, 9(4):733–752, 2023. 6
- [20] Ali Hassani, Steven Walton, Jiachen Li, Shen Li, and Humphrey Shi. Neighborhood attention transformer. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 6185–6194, 2023. 3, 6, 7
- [21] Ali Hatamizadeh, Hongxu Yin, Greg Heinrich, Jan Kautz, and Pavlo Molchanov. Global context vision transformers. In *International Conference on Machine Learning*, pages 12633–12646. PMLR, 2023. 6
- [22] Junjun He, Zhongying Deng, and Yu Qiao. Dynamic multi-scale filters for semantic segmentation. In *Proceedings of the IEEE/CVF International Conference on Computer Vision*, pages 3562–3572, 2019. 3
- [23] Kaiming He, Xiangyu Zhang, Shaoqing Ren, and Jian Sun. Deep residual learning for image recognition. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, pages 770–778, 2016. 1, 3
- [24] Kaiming He, Georgia Gkioxari, Piotr Dollár, and Ross Girshick. Mask r-cnn. In *Proceedings of the IEEE International Conference on Computer Vision*, pages 2961–2969, 2017. 7
- [25] Dan Hendrycks, Steven Basart, Norman Mu, Saurav Kadavath, Frank Wang, Evan Dorundo, Rahul Desai, Tyler Zhu, Samyak Parajuli, Mike Guo, et al. The many faces of robustness: A critical analysis of out-of-distribution generalization. In *Proceedings of the IEEE/CVF international conference on computer vision*, pages 8340–8349, 2021. 10
- [26] Dan Hendrycks, Kevin Zhao, Steven Basart, Jacob Steinhardt, and Dawn Song. Natural adversarial examples. In *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*, pages 15262–15271, 2021. 10
- [27] Qibin Hou, Cheng-Ze Lu, Ming-Ming Cheng, and Jiashi Feng. Conv2former: A simple transformer-style convnet for visual recognition. *IEEE Transactions on Pattern Analysis and Machine Intelligence*, 2024. 6
- [28] Peiyun Hu and Deva Ramanan. Bottom-up and top-down reasoning with hierarchical rectified gaussians. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, pages 5600–5609, 2016. 1, 3
- [29] Gao Huang, Yu Sun, Zhuang Liu, Daniel Sedra, and Kilian Q Weinberger. Deep networks with stochastic depth. In *European Conference on Computer Vision*, pages 646–661. Springer, 2016. 6
- [30] Gao Huang, Zhuang Liu, Laurens Van Der Maaten, and Kilian Q Weinberger. Densely connected convolutional networks. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, pages 4700–4708, 2017. 3
- [31] Hu Jie, Shen Li, Sun Gang, and Samuel Albanie. Squeeze-and-excitation networks. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, 2018. 4
- [32] Donghyun Kim, Byeongho Heo, and Dongyoon Han. Densenets reloaded: Paradigm shift beyond resnets and vits. In *European Conference on Computer Vision*, 2024. 3, 6, 7, 10
- [33] Alex Krizhevsky, Ilya Sutskever, and Geoffrey E Hinton. Imagenet classification with deep convolutional neural networks. *Advances in neural information processing systems*, 25, 2012. 3
- [34] Jimmy Lei Ba, Jamie Ryan Kiros, and Geoffrey E Hinton. Layer normalization. *ArXiv e-prints*, pages arXiv-1607, 2016. 4
- [35] Jerome Y Lettvin et al. On seeing sidelong. *The Sciences*, 16(4):10–20, 1976. 3
- [36] Chao Li, Aojun Zhou, and Anbang Yao. Omni-dimensional dynamic convolution. In *International Conference on Learning Representations*, 2022. 2, 3, 7
- [37] Duo Li, Jie Hu, Changhu Wang, Xiangtai Li, Qi She, Lei Zhu, Tong Zhang, and Qifeng Chen. Involution: Inverting the inherence of convolution for visual recognition. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 12321–12330, 2021. 2, 3, 7
- [38] Kunchang Li, Yali Wang, Junhao Zhang, Peng Gao, Guanglu Song, Yu Liu, Hongsheng Li, and Yu Qiao. Uniformer: Unifying convolution and self-attention for visual recognition. *IEEE Transactions on Pattern Analysis and Machine Intelligence*, 2023. 6, 7
- [39] Siyuan Li, Zedong Wang, Zicheng Liu, Cheng Tan, Haitao Lin, Di Wu, Zhiyuan Chen, Jiangbin Zheng, and Stan Z Li. Moganet: Multi-order gated aggregation network. In *International Conference on Learning Representations*, 2023. 3, 4, 6, 7
- [40] Zhaoping Li. *Understanding vision: theory, models, and data*. Oxford University Press, USA, 2014. 1, 3
- [41] Tsung-Yi Lin, Michael Maire, Serge Belongie, James Hays, Pietro Perona, Deva Ramanan, Piotr Dollár, and C Lawrence Zitnick. Microsoft coco: Common objects in context. In *European Conference on Computer Vision*, pages 740–755. Springer, 2014. 7
- [42] Shiwei Liu, Tianlong Chen, Xiaohan Chen, Xuxi Chen, Qiao Xiao, Boqian Wu, Mykola Pechenizkiy, Decebal Mocanu, and Zhangyang Wang. More convnets in the 2020s: Scaling up kernels beyond 51x51 using sparsity. In *International Conference on Learning Representations*, 2023. 2, 3, 5, 6, 7, 8
- [43] Yue Liu, Yunjie Tian, Yuzhong Zhao, Hongtian Yu, Lingxi Xie, Yaowei Wang, Qixiang Ye, and Yunfan Liu. Vmamba: Visual state space model. *Advances in Neural Information Processing Systems*, 2024. 1, 2, 3, 6, 7
- [44] Ze Liu, Yutong Lin, Yue Cao, Han Hu, Yixuan Wei, Zheng Zhang, Stephen Lin, and Baining Guo. Swin transformer:

- Hierarchical vision transformer using shifted windows. In *Proceedings of the IEEE/CVF International Conference on Computer Vision*, pages 10012–10022, 2021. 2, 3, 6, 7, 8
- [45] Zhuang Liu, Hanzi Mao, Chao-Yuan Wu, Christoph Feichtenhofer, Trevor Darrell, and Saining Xie. A convnet for the 2020s. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 11976–11986, 2022. 1, 2, 3, 6, 7, 10
- [46] Ilya Loshchilov and Frank Hutter. Decoupled weight decay regularization. *arXiv preprint arXiv:1711.05101*, 2017. 6
- [47] Meng Lou, Yunxiang Fu, and Yizhou Yu. Sparx: A sparse cross-layer connection mechanism for hierarchical vision mamba and transformer networks. In *Proceedings of the AAAI Conference on Artificial Intelligence*, 2025. 2
- [48] Meng Lou, Shu Zhang, Hong-Yu Zhou, Sibe Yang, Chuan Wu, and Yizhou Yu. Transxnet: Learning both global and local dynamics with a dual dynamic token mixer for visual recognition. *IEEE Transactions on Neural Networks and Learning Systems*, 2025. 3
- [49] Wenjie Luo, Yujia Li, Raquel Urtasun, and Richard Zemel. Understanding the effective receptive field in deep convolutional neural networks. *Advances in Neural Information Processing Systems*, 29, 2016. 2, 12
- [50] Xu Ma, Xiyang Dai, Yue Bai, Yizhou Wang, and Yun Fu. Rewrite the stars. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 5694–5703, 2024. 3, 4
- [51] Xu Ma, Xiyang Dai, Jianwei Yang, Bin Xiao, Yinpeng Chen, Yun Fu, and Lu Yuan. Efficient modulation for vision networks. In *International Conference on Learning Representations*, 2024. 3
- [52] Juhong Min, Yucheng Zhao, Chong Luo, and Minsu Cho. Peripheral vision transformer. *Advances in Neural Information Processing Systems*, 35:32097–32111, 2022. 3
- [53] Sarthak Mittal, Alex Lamb, Anirudh Goyal, Vikram Voleti, Murray Shanahan, Guillaume Lajoie, Michael Mozer, and Yoshua Bengio. Learning to combine top-down and bottom-up signals in recurrent neural networks with attention over modules. In *International Conference on Machine Learning*, pages 6972–6986. PMLR, 2020. 1, 3
- [54] Bo Pang, Yizhuo Li, Jiefeng Li, Muchen Li, Hanwen Cao, and Cewu Lu. Tdaf: Top-down attention framework for vision tasks. In *Proceedings of the AAAI Conference on Artificial Intelligence*, pages 2384–2392, 2021. 1, 3
- [55] Prajit Ramachandran, Niki Parmar, Ashish Vaswani, Irwan Bello, Anselm Levskaya, and Jon Shlens. Stand-alone self-attention in vision models. *Advances in neural information processing systems*, 32, 2019. 3
- [56] Yongming Rao, Wenliang Zhao, Yansong Tang, Jie Zhou, Ser Nam Lim, and Jiwen Lu. Hornet: Efficient high-order spatial interactions with recursive gated convolutions. *Advances in Neural Information Processing Systems*, 35: 10353–10366, 2022. 3, 7
- [57] Benjamin Recht, Rebecca Roelofs, Ludwig Schmidt, and Vaishaal Shankar. Do imagenet classifiers generalize to imagenet? In *International Conference on Machine Learning*, pages 5389–5400. PMLR, 2019. 10
- [58] Yuri B Saalman, Ivan N Pigarev, and Trichur R Vidyasagar. Neural mechanisms of visual attention: how top-down feedback highlights relevant locations. *Science*, 316(5831): 1612–1615, 2007. 1, 3
- [59] Ramprasaath R Selvaraju, Michael Cogswell, Abhishek Das, Ramakrishna Vedantam, Devi Parikh, and Dhruv Batra. Grad-cam: Visual explanations from deep networks via gradient-based localization. *International journal of computer vision*, 128:336–359, 2020. 2, 12
- [60] Baifeng Shi, Trevor Darrell, and Xin Wang. Top-down visual attention from analysis by synthesis. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 2102–2112, 2023. 1, 3, 8
- [61] Karen Simonyan and Andrew Zisserman. Very deep convolutional networks for large-scale image recognition. In *International Conference on Learning Representations*, 2015. 3
- [62] Shitao Tang, Jiahui Zhang, Siyu Zhu, and Ping Tan. Quadtree attention for vision transformers. In *International Conference on Learning Representations*, 2022. 6
- [63] Hugo Touvron, Matthieu Cord, Matthijs Douze, Francisco Massa, Alexandre Sablayrolles, and Hervé Jégou. Training data-efficient image transformers & distillation through attention. In *International Conference on Machine Learning*, pages 10347–10357. PMLR, 2021. 3, 6
- [64] Zhengzhong Tu, Hossein Talebi, Han Zhang, Feng Yang, Peyman Milanfar, Alan Bovik, and Yinxiao Li. Maxvit: Multi-axis vision transformer. In *European Conference on Computer Vision*. Springer, 2022. 2, 6
- [65] Haohan Wang, Songwei Ge, Zachary Lipton, and Eric P Xing. Learning robust global representations by penalizing local predictive power. *Advances in Neural Information Processing Systems*, 32, 2019. 10
- [66] Wenhai Wang, Enze Xie, Xiang Li, Deng-Ping Fan, Kaitao Song, Ding Liang, Tong Lu, Ping Luo, and Ling Shao. Pyramid vision transformer: A versatile backbone for dense prediction without convolutions. In *Proceedings of the IEEE/CVF Conference on Computer Vision*, pages 568–578, 2021. 1, 2, 3
- [67] Wenhai Wang, Enze Xie, Xiang Li, Deng-Ping Fan, Kaitao Song, Ding Liang, Tong Lu, Ping Luo, and Ling Shao. Pvt v2: Improved baselines with pyramid vision transformer. *Computational Visual Media*, 8(3):415–424, 2022. 1, 2, 4, 6, 7
- [68] Wenxiao Wang, Wei Chen, Qibo Qiu, Long Chen, Boxi Wu, Binbin Lin, Xiaofei He, and Wei Liu. Crossformer++: A versatile vision transformer hinging on cross-scale attention. *IEEE Transactions on Pattern Analysis and Machine Intelligence*, 2023. 6
- [69] Wenhai Wang, Jifeng Dai, Zhe Chen, Zhenhang Huang, Zhiqi Li, Xizhou Zhu, Xiaowei Hu, Tong Lu, Lewei Lu, Hongsheng Li, et al. Internimage: Exploring large-scale vision foundation models with deformable convolutions. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, 2023. 2, 3, 6, 7
- [70] Sanghyun Woo, Shoubhik Debnath, Ronghang Hu, Xinlei Chen, Zhuang Liu, In So Kweon, and Saining Xie. Convnext v2: Co-designing and scaling convnets with masked

- autoencoders. In *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*, pages 16133–16142, 2023. 3
- [71] Tete Xiao, Yingcheng Liu, Bolei Zhou, Yuning Jiang, and Jian Sun. Unified perceptual parsing for scene understanding. In *European Conference on Computer Vision*, pages 418–434, 2018. 7
- [72] Huijuan Xu and Kate Saenko. Ask, attend and answer: Exploring question-guided spatial attention for visual question answering. In *European Conference on Computer Vision*, pages 451–466. Springer, 2016. 1, 3
- [73] Brandon Yang, Gabriel Bender, Quoc V Le, and Jiquan Ngiam. Condconv: Conditionally parameterized convolutions for efficient inference. *Advances in Neural Information Processing Systems*, 32, 2019. 3
- [74] Jianwei Yang, Chunyuan Li, Xiyang Dai, and Jianfeng Gao. Focal modulation networks. *Advances in Neural Information Processing Systems*, 35:4203–4217, 2022. 3, 6, 7
- [75] Weihao Yu, Pan Zhou, Shuicheng Yan, and Xinchao Wang. Inceptionnext: When inception meets convnext. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, 2024. 2, 3, 6, 7, 10
- [76] Li Yuan, Qibin Hou, Zihang Jiang, Jiashi Feng, and Shuicheng Yan. Volo: Vision outlooker for visual recognition. *IEEE transactions on pattern analysis and machine intelligence*, 45(5):6575–6586, 2022. 3, 7
- [77] Amir R Zamir, Te-Lin Wu, Lin Sun, William B Shen, Bertram E Shi, Jitendra Malik, and Silvio Savarese. Feedback networks. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pages 1308–1317, 2017. 1, 3
- [78] Jianming Zhang, Sarah Adel Bargal, Zhe Lin, Jonathan Brandt, Xiaohui Shen, and Stan Sclaroff. Top-down neural attention by excitation backprop. *International Journal of Computer Vision*, 126(10):1084–1102, 2018. 3
- [79] Hengshuang Zhao, Jianping Shi, Xiaojuan Qi, Xiaogang Wang, and Jiaya Jia. Pyramid scene parsing network. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pages 2881–2890, 2017. 9
- [80] Bolei Zhou, Hang Zhao, Xavier Puig, Sanja Fidler, Adela Barriuso, and Antonio Torralba. Scene parsing through ade20k dataset. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pages 633–641, 2017. 7
- [81] Lei Zhu, Xinjiang Wang, Zhanghan Ke, Wayne Zhang, and Rynson WH Lau. Biformer: Vision transformer with bi-level routing attention. In *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*, pages 10323–10333, 2023. 2, 3, 6, 7
- [82] Lianghui Zhu, Bencheng Liao, Qian Zhang, Xinlong Wang, Wenyu Liu, and Xinggang Wang. Vision mamba: Efficient visual representation learning with bidirectional state space model. In *International Conference on Machine Learning*, 2024. 2