

Frequency Dynamic Convolution for Dense Image Prediction

Linwei Chen¹ Lin Gu^{2,3} Liang Li⁴ Chenggang Yan^{5,6} Ying Fu^{1*}

¹Beijing Institute of Technology ²RIKEN ³The University of Tokyo

⁴Chinese Academy of Sciences ⁵Hangzhou Dianzi University ⁶Tsinghua University

chenlinwei@bit.edu.cn; lin.gu@riken.jp; liang.li@ict.ac.cn; cgyan@hdu.edu.cn; fuying@bit.edu.cn

Abstract

While Dynamic Convolution (DY-Conv) has shown promising performance by enabling adaptive weight selection through multiple parallel weights combined with an attention mechanism, the frequency response of these weights tends to exhibit high similarity, resulting in high parameter costs but limited adaptability. In this work, we introduce Frequency Dynamic Convolution (FDConv), a novel approach that mitigates these limitations by learning a fixed parameter budget in the Fourier domain. FDConv divides this budget into frequency-based groups with disjoint Fourier indices, enabling the construction of frequency-diverse weights without increasing the parameter cost. To further enhance adaptability, we propose Kernel Spatial Modulation (KSM) and Frequency Band Modulation (FBM). KSM dynamically adjusts the frequency response of each filter at the spatial level, while FBM decomposes weights into distinct frequency bands in the frequency domain and modulates them dynamically based on local content. Extensive experiments on object detection, segmentation, and classification validate the effectiveness of FDConv. We demonstrate that when applied to ResNet-50, FDConv achieves superior performance with a modest increase of +3.6M parameters, outperforming previous methods that require substantial increases in parameter budgets (e.g., CondConv +90M, KW +76.5M). Moreover, FDConv seamlessly integrates into a variety of architectures, including ConvNeXt, Swin-Transformer, offering a flexible and efficient solution for modern vision tasks. The code is made publicly available at <https://github.com/Linwei-Chen/FDConv>.

1. Introduction

Convolution, the core operation in ConvNets, has driven decades of advancement in computer vision [1, 3–6, 14, 18, 23, 29, 35, 38, 44, 47, 61, 77, 80]. Essential for capturing lo-

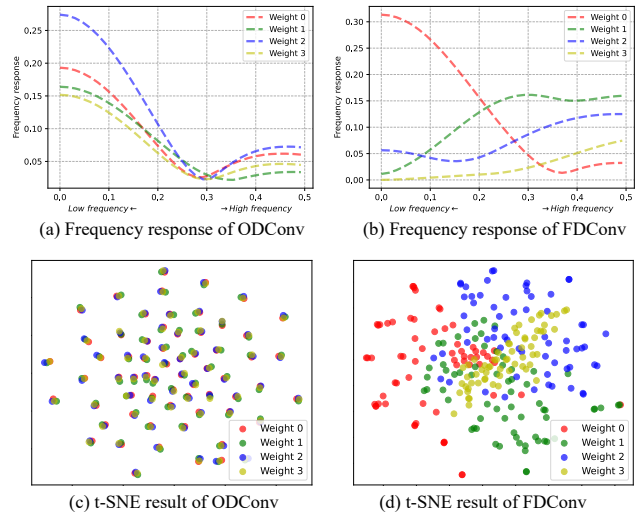


Figure 1. Weight frequency responses and t-SNE analyses. We set the number of weights to 4 to align with ODConv [32]. (a) The frequency responses of the four parallel weights in ODConv are highly similar, indicating limited diversity. (b) In contrast, FDConv shows distinct frequency responses for each weight, spanning different parts of the frequency spectrum. (c) The t-SNE plot for ODConv reveals that the filters in the four weights are closely clustered, suggesting a lack of diversity. (d) The t-SNE plot for FDConv shows that the filters in the four weights have different distributions, indicating greater diversity.

cal patterns and building hierarchical representations, it remains fundamental in modern architectures [11, 47, 62, 71].

Building upon the success of standard convolution, Dynamic Convolution (DY-Conv) [63, 73] offers a more adaptive and efficient approach. Unlike standard convolution with fixed weights, DY-Conv uses multiple parallel weights combined by an attention module, allowing sample-specific weight adaptation with minimal extra computation.

However, our analysis in Figure 1 reveals that traditional dynamic convolution [31, 32, 63, 73] lack of frequency responses diversity in their parallel weights. As shown in Figure 1(a), these weights exhibit highly similar frequency responses, while the t-SNE visualization in Figure 1(c) indicates that filters within ODConv [32] are clus-

*Corresponding Author

tered closely. Despite a significant increase in parameters (*e.g.*, $4\times$ in [31, 32, 63, 73]), this limited frequency diversity reduces the model’s ability to adaptively capture frequency information. For example, extracting low-frequency components helps suppress noise [3], while high-frequency components capture details and boundaries [2, 36, 51, 78], which are vital for foreground-background differentiation.

To address these limitations, we propose Frequency Dynamic Convolution (FDConv), as shown in Figure 2. It is designed to enhance frequency adaptability without incurring excessive parameter overhead. Our approach is based on three core modules: Fourier Disjoint Weight, Kernel Spatial Modulation, and Frequency Band Modulation.

Unlike traditional methods [32, 63, 73], which learn weights in the spatial domain, the Fourier Disjoint Weight (FDW) constructs kernel weights by learning spectral coefficients in the Fourier domain. These coefficients are divided into frequency-based groups, each with a disjoint set of Fourier indices. An inverse Discrete Fourier Transform (iDFT) is then applied to these groups, converting them into spatial weights. This disjoint grouping enables each weight to exhibit distinct frequency responses (as shown in Figure 1(b)), ensuring high diversity among the learned weights (also shown in Figure 1(d)).

Kernel Spatial Modulation (KSM) enhances flexibility by precisely adjusting the frequency response of each filter at the spatial level within the kernel. By combining local and global channel information, KSM generates a dense matrix of modulation values that finely tunes each individual weight element. This fine-grained control enables FDConv to dynamically adapt each filter element, allowing for frequency response adjustment across the entire kernel.

Frequency Band Modulation (FBM) decomposes weights into distinct frequency bands in the frequency domain, enabling spatially variant frequency modulation. It allows each frequency band of the weight to be adjusted independently across spatial locations. Unlike traditional dynamic convolutions, which apply fixed frequency responses across spatial dimensions, FBM decomposes weights into distinct frequency bands and dynamically modulates them based on local content. This design enables the model to selectively emphasize or suppress frequency bands across different regions, adaptively capturing diverse frequency information in a spatially variant manner.

Moreover, unlike previous works [32, 63, 73], which increase parameter costs by a factor of n (where n is the number of weights, typically $n < 10$ [32, 63, 73]), our FDConv maintains a fixed parameter budget while generating a large number of frequency-diverse weight kernels ($n > 10$) by dividing parameters in the Fourier domain into disjoint frequency-based groups. This design allows the model to efficiently learn weights with distinct frequency responses without burdening parameter cost.

Extensive experiments on object detection, instance segmentation, semantic segmentation, and image classification validate the effectiveness of FDConv. For example, when applied to ResNet-50, FDConv achieves superior performance with a modest increase of +3.6M parameters, outperforming previous methods that require substantial increases in parameter budgets (*e.g.*, CondConv +90M, DY-Conv +75.3M, ODConv +65.1M, KW +76.5M) [31, 32, 63, 73]. FDConv can be seamlessly integrated into various architectures, including ConvNeXt and Swin Transformer, where it replaces the linear layer (as a 1×1 convolution), offering a versatile and efficient solution.

- We conduct a comprehensive exploration of dynamic convolution using frequency analysis. Our findings reveal that the parameters of traditional dynamic convolution methods exhibit high homogeneity in frequency response across learned parallel weights, resulting in high parameter redundancy and limited adaptability.
- We introduce the Fourier Disjoint Weight (FDW), Kernel Spatial Modulation (KSM), and Frequency Band Modulation (FBM) strategies. FDW constructs multiple weights with diversified frequency responses without increasing the parameter cost, KSM enhances the representation power by adjusting weights element-wise, and FBM improves convolution by precisely extracting frequency bands in a spatially variant manner.
- We demonstrate that our approach can be easily integrated into existing ConvNets and vision transformers. Comprehensive experiments on segmentation tasks show that it surpasses previous state-of-the-art dynamic convolution methods, requiring only a minor increase in parameters, consistently demonstrating its effectiveness.

2. Related Work

Feature Recalibration. Feature recalibration through attention mechanisms has proven highly effective in deep learning models. Methods such as RAN [64], SE [26], CBAM [69], GE [25], SRM [30], ECA [67], and SimAtt [74] focus on adaptively emphasizing informative features or suppressing irrelevant ones across channels and spatial dimensions of the feature map, *i.e.*, channel and spatial attention. In contrast, our approach introduces frequency-specific recalibration for convolution weights.

Dynamic Weight Networks. Recently, dynamic networks have shown to be effective in various computer vision tasks. Dynamic Filter Networks [28] and Kernel Prediction Networks [53] generate sample-adaptive filters conditioned on the input. In contrast, Hypernetworks [20] generate weights for a larger recurrent network instead of ConvNets. Building upon similar idea, CARAFE [66] and Involution [33] have developed efficient modules that predict spatially variant convolution weights.

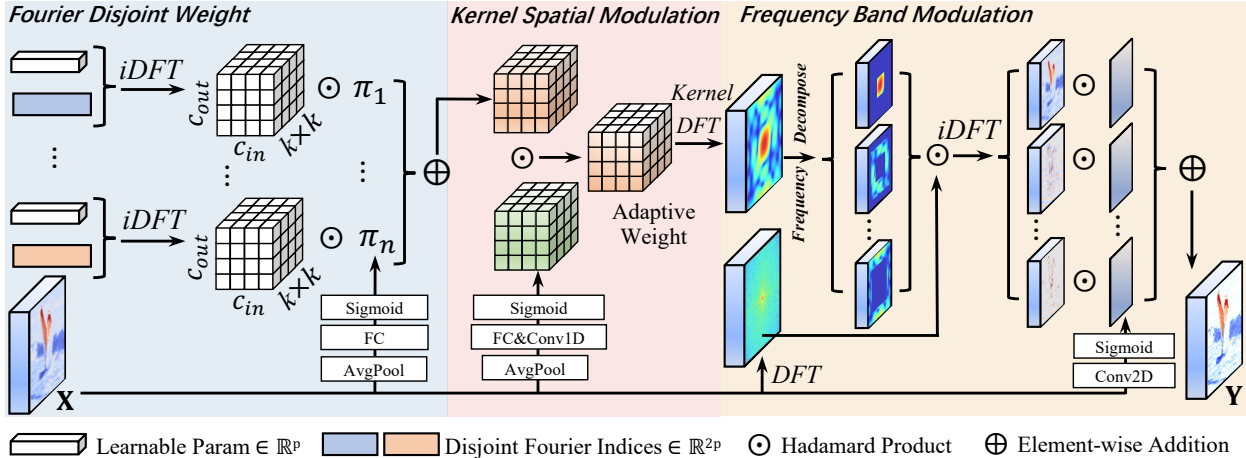


Figure 2. Illustration of the proposed Frequency Dynamic Convolution, which consists of the Fourier Disjoint Weight (FDW), Kernel Spatial Modulation (KSM), and Frequency Band Modulation (FBM) modules. FC indicates fully connected layer.

Dynamic convolution methods [63, 73] learn multiple parallel weights and adaptively mix them linearly using attention modules. CondConv [73] uses a sigmoid function for weight fusion, while DY-Conv improves upon this by using a softmax function [63]. Inspired by SE [26], WeightNet [50], CGC [43], and WE [56], these methods design various attention modules to adjust convolutional weights in ConvNets. ODCConv [32] further enhances the attention module by predicting channel-wise, filter-wise, and spatial-wise attention values to adjust the weights.

To mitigate the increased parameter overhead of multiple weights, methods like DCD [37] and PEDConv [24] use matrix decomposition techniques to construct low-rank weight matrices, reducing computational complexity. More recently, KW [31] introduced a decomposition approach where kernel weights are divided into smaller, shareable units across different stages and layers, enabling dynamic kernel reconstruction with fewer parameters.

In contrast, our FDCConv addresses the heavy parameter cost and limited diversity of weights from the frequency aspect, offering a new solution.

Frequency Domain Learning. Frequency-domain analysis has long been a cornerstone of signal processing [15, 54]. Recently, these techniques have been leveraged in deep learning, influencing model optimization strategies [75], robustness [49], and generalization abilities [65] in Deep Neural Networks (DNNs). Moreover, the integration of frequency-domain methods into DNNs has proven effective for learning global features [9, 19, 27, 39, 58] and enhancing domain-generalizable representations [40]. FcaNet [55] demonstrates that frequency information benefits feature recalibration, while FreqFusion [2] shows its advantages in feature fusion. Some studies [6, 16, 17] have improved downsampling operations by addressing high-frequency components that lead to aliasing. FADC [7] enhances di-

lated convolution by adjusting dilation based on the frequency characteristics of the features. Our method incorporates a frequency-based perspective into dynamic convolution, improving its ability to learn diversified weights for capturing a wider range of frequency information.

3. Method

An overview of the proposed Frequency Dynamic Convolution (FDCConv) framework is shown in Figure 2. This section first introduces the concept of Fourier Disjoint Weights, followed by a detailed exploration of two key strategies: Kernel Spatial Modulation and Frequency Band Modulation, which are designed to fully leverage the frequency adaptability of FDCConv in the kernel spatial and frequency domains, respectively.

3.1. Fourier Disjoint Weight

Dynamic Convolution. For a standard convolutional layer, it can be formulated as $\mathbf{Y} = \mathbf{W} * \mathbf{X}$, where $\mathbf{X} \in \mathbb{R}^{h \times w \times C_{in}}$ and $\mathbf{Y} \in \mathbb{R}^{h \times w \times C_{out}}$ are the input and output features, respectively. Here, C_{in} and C_{out} represent the number of input and output feature channels, and $h \times w$ denotes the spatial size. The weight $\mathbf{W} \in \mathbb{R}^{k \times k \times C_{in} \times C_{out}}$ consists of C_{out} convolutional filters, each with a spatial size of $k \times k$.

Dynamic convolution [63, 73] enhances the adaptability of convolutional layers by replacing the static weight \mathbf{W} in standard convolution with a combination of n distinct weights $\{\mathbf{W}_1, \dots, \mathbf{W}_n\}$, each of the same dimension. The contribution of each kernel is modulated by a set of attention-based coefficients $\{\pi_1, \dots, \pi_n\}$, which are dynamically generated. Typically, these coefficients are derived by applying global average pooling on the input, followed by a fully connected (FC) layer. This dynamic convolution operation can be expressed as

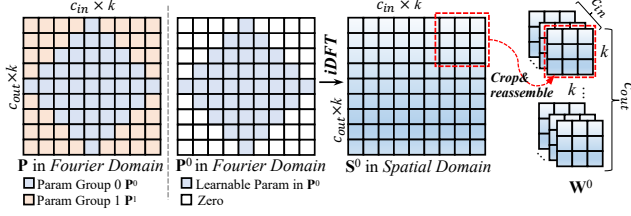


Figure 3. Illustration of Fourier Disjoint Weight (FDW). The left figure illustrates the division of parameters into disjoint groups, ranging from low frequencies (center) to high frequencies (border). In this example, $n = 2$ groups are shown. The right figure demonstrates how to obtain the convolution weights from the learnable parameter group 0. It first transforms the learnable parameters with specific Fourier indices (with all other Fourier indices set to zero) using the inverse Discrete Fourier Transform (iDFT). The resulting spatial weights are then obtained by cropping the iDFT result into $k \times k$ patches and reshaping them into a weight tensor of size $k \times k \times C_{in} \times C_{out}$.

$$\mathbf{W} = \pi_1 \mathbf{W}_1 + \dots + \pi_n \mathbf{W}_n. \quad (1)$$

Despite the increased parameter cost by a factor of n , we expect dynamic convolution to learn diverse weights. However, our analysis reveals, as shown in Figure 1, that the frequency responses of parallel weights are highly similar. This lack of frequency diversity limits the model’s ability to adaptively capture features across different frequency bands, reducing the flexibility of dynamic convolution.

Overview of Fourier Disjoint Weight. To construct multiple parallel weights with high frequency response diversity without increasing parameter costs, we propose Fourier Disjoint Weight (FDW). Unlike previous methods [32, 63, 73], which are limited to a small number of kernels ($n < 10$) due to high parameter costs, FDW can generate $n > 10$ diversified weights.

The core concept of FDW is learning spectral coefficients in the Fourier domain with disjoint sets of Fourier indices, rather than in the traditional spatial domain. FDW involves three steps to construct n weights: 1) *Fourier disjoint grouping*. Divide a fixed number of parameters into n groups with disjoint Fourier indices. 2) *Fourier to spatial transformation*. Convert each group of parameters from the Fourier domain to the spatial domain using the Inverse Discrete Fourier Transform (iDFT). 3) *Reassembling*. Crop the transformed results in the spatial domain into $k \times k$ patches and reassemble them into the standard weight shape of $k \times k \times C_{in} \times C_{out}$.

Fourier Disjoint Grouping. Given a parameter budget of $k \times k \times C_{in} \times C_{out}$, FDW first treats these parameters as learnable spectral coefficients in the Fourier domain, reshaping them into $\mathbf{P} \in \mathbb{R}^{k \times k \times C_{in} \times C_{out}}$. Each parameter is associated with a Fourier index (u, v) , i.e., coordinates in the Fourier domain that indicate frequency. FDW then sorts these parameters from low to high frequency based on the L_2 norm

of the Fourier index, $\|(u, v)\|_2$, and divides them uniformly into disjoint n set, $\{\mathbf{P}^0, \dots, \mathbf{P}^{n-1}\}$.

As shown on the left side of Figure 3, we divide the learnable parameters into $n = 2$ groups for simplicity in the demonstration, where the center represents low frequencies and the border represents high frequencies. Moreover, the number of groups, n , can be set to a large value ($n > 10$), allowing for the generation of a large number of diversified weights without increasing the parameter cost.

Fourier to Spatial Transformation. To obtain the weights, FDW transforms each group of parameters into the spatial domain using the inverse Discrete Fourier Transform (iDFT). This can be formulated as:

$$\mathbf{S}_{p,q}^i = \sum_{u=0}^{kC_{in}-1} \sum_{v=0}^{kC_{out}-1} \mathbf{P}_{u,v}^i e^{i2\pi \left(\frac{p}{kC_{in}} u + \frac{q}{kC_{out}} v \right)} \quad (2)$$

where $\mathbf{P}_{u,v}^i$ is the parameter with Fourier index (u, v) in the i -th group \mathbf{P}^i . As shown on the right side of Figure 3, if the Fourier index (u, v) is assigned to the i -th group, then $\mathbf{P}_{u,v}^i = \mathbf{P}_{u,v}$, otherwise, $\mathbf{P}_{u,v}^i = 0$. $\mathbf{S}_{p,q}^i \in \mathbb{R}^{k \times k \times C_{in} \times C_{out}}$ represents the element at position (p, q) in the converted results in the spatial domain.

Reassembling. As shown on the right side of Figure 3, the final i -th weight $\mathbf{W}^i \in \mathbb{R}^{k \times k \times C_{in} \times C_{out}}$ can be obtained by cropping $\mathbf{S}^i \in \mathbb{R}^{k \times k \times C_{in} \times C_{out}}$ into $C_{in} \times C_{out}$ patches of size $k \times k$ and reassembling them to form \mathbf{W}^i .

Since the parameters are divided according to frequency, \mathbf{S}^i contains only the frequency components of a specific band. Therefore, each weight \mathbf{W}^i derived from \mathbf{S}^i exhibits a distinct frequency response compared to \mathbf{W}^j when $i \neq j$. This ensures that the frequency responses of the constructed weights are diversified. After a linear mixture, as described in Equation (1), FDW can adaptively adjust the frequency response of the combined weight based on the input sample.

Note that FDW can also be applied to linear layers in modern vision architectures, such as Transformers [12, 46], which are equivalent to convolutions with a kernel size of 1.

3.2. Kernel Spatial Modulation

By ensuring the diversity of frequency responses of parallel weights in the Fourier domain, FDW can adaptively adjust the frequency response of the combined weight based on the input sample after a linear mixture, as described in Equation (1). However, this weight-wise mixture is too coarse and cannot independently adjust the frequency response of each $k \times k$ filter in the weight.

To address this issue, we propose Kernel Spatial Modulation (KSM), which predicts a dense modulation matrix $\alpha \in \mathbb{R}^{k \times k \times C_{in} \times C_{out}}$, instead of a sparse vector. As shown in Figure 4, KSM consists of a local channel branch and a global channel branch.

Local Channel Branch. While global fully connected layers are commonly used for modulation value predic-

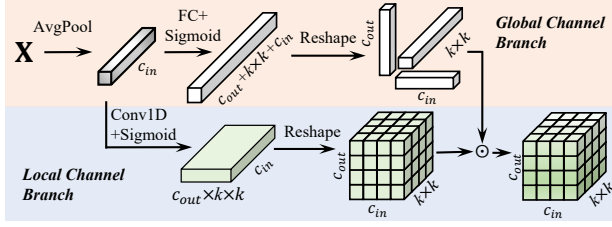


Figure 4. Illustration of Kernel Spatial Modulation (KSM). The KSM consists of two branches: the global channel branch and the local channel branch. The local channel branch employs a very lightweight 1-D convolution to obtain local channel information and predicts a dense modulation matrix of size $k \times k \times C_{in} \times C_{out}$. The global branch uses a fully connected layer to obtain the global channel information and predicts three dimension-wise modulation values along the input channel, output channel, and kernel spatial dimensions. The two branches are fused to obtain the final weight modulation matrix.

tion [26, 32, 63, 73], they are not suitable for predicting dense modulation matrices due to their large parameter and computational costs. To address this, the local channel branch employs a lightweight 1-D convolution, which has been proven to be efficient and effective [67]. It captures local channel information and predicts a dense modulation matrix of size $k \times k \times C_{in} \times C_{out}$. This approach significantly reduces parameters and computational complexity while maintaining the ability to learn fine-grained modulation for each element in the weight.

Global Channel Branch. Though the local channel branch can efficiently predict the modulation matrices, it lacks global information, which is crucial for weight adjustments. To complement the local channel branch, the global channel branch uses a fully connected layer to capture global channel information and predict a sparse modulation vector for efficiency. Specifically, it predicts three dimension-wise modulation values: one for the input channel, one for the output channel, and one for the kernel spatial dimensions, ensuring that both local and global contextual information are incorporated for adaptive modulation.

In this way, the proposed KSM is able to leverage both local and global information, enabling more precise and context-aware modulation of each filter in the weights.

3.3. Frequency Band Modulation

While the proposed Fourier Disjoint Weight (FDW) and Kernel Spatial Modulation (KSM) modules substantially enhance adaptability by ensuring frequency diversity and element-wise adjustments, they remain spatially invariant, as is typical in dynamic convolution [31, 32, 63, 73]. This spatial invariance, where weights are shared across the entire feature map, restricts convolutional layers from dynamically adapting frequency responses to spatially varying content, limiting their ability to fully capture complex struc-

tures across the image.

Natural images and their corresponding features exhibit large spatial variation, which necessitates frequency-specific adaptations for optimal feature extraction. For example, extracting low-frequency components is vital for suppressing feature noise [3], while high-frequency components are essential for capturing fine details and boundaries [2, 51], which are crucial for distinguishing the foreground from the background.

Overview of Frequency Band Modulation. To address the need for spatially dynamic frequency modulation, we propose Frequency Band Modulation (FBM). FBM decomposes the convolutional kernel into multiple frequency bands in the frequency domain and applies spatially specific modulations, adaptively adjusting each frequency component across different spatial locations.

The Frequency Band Modulation operates in the following key steps: 1) Kernel frequency decomposition. Decomposing the frequency response of the convolution weight into different frequency bands. 2) Convolution in the Fourier domain. Performing convolution in the Fourier domain. 3) Spatially variant modulation. Predicting modulation values for each frequency band of the convolution weight across different spatial locations.

The core formulation of FBM is given by:

$$\mathbf{Y} = \sum_{b=0}^{B-1} (\mathbf{A}_b \odot (\mathbf{W}_b * \mathbf{X})), \quad (3)$$

where \mathbf{X} and $\mathbf{Y} \in \mathbb{R}^{h \times w}$ are the input and output feature maps. Note that we omit the channel dimension for simplicity. $\mathbf{A}_b \in \mathbb{R}^{h \times w}$ representing spatial modulation values specific to the b -th frequency band, and \mathbf{W}_b the b -th frequency band of weight. FBM enables adjusting the frequency responses for each spatial location of feature map.

Kernel frequency decomposition. To decompose the convolution kernel into distinct frequency bands, FBM first pads the kernel \mathbf{W} to match the feature map size [52], and then applies a set of binary masks \mathcal{M}_b to isolate specific frequency ranges:

$$\mathbf{W}_b = \mathcal{F}^{-1}(\mathcal{M}_b \odot \mathcal{F}(\mathbf{W})), \quad (4)$$

where \mathcal{F} and \mathcal{F}^{-1} denote the DFT and inverse DFT, and \mathcal{M}_b is a binary mask isolating specific frequency ranges:

$$\mathcal{M}_b(u, v) = \begin{cases} 1 & \text{if } \psi_b \leq \max(|u|, |v|) < \psi_{b+1} \\ 0 & \text{otherwise} \end{cases} \quad (5)$$

Here, ψ_b and ψ_{b+1} are thresholds from the predefined frequency set $\{0, \psi_1, \dots, \psi_{B-1}, \frac{1}{2}\}$, where (u, v) denote the horizontal and vertical frequency indices. By default, we decompose the frequency spectrum into four distinct bands using an octave-based partitioning strategy [60]. The thresholds for dividing the frequency bands are $\{0, \frac{1}{16}, \frac{1}{8}, \frac{1}{4}, \frac{1}{2}\}$.

Convolution in the Fourier Domain. After obtaining the frequency-specific weights \mathbf{W}_b , the output for the corresponding frequency band can be computed as follows:

$$\mathbf{Y}_b = \mathbf{W}_b * \mathbf{X}, \quad (6)$$

where \mathbf{Y}_b represents the output for the b -th frequency band.

However, as discussed in [15, 16], obtaining specific frequency bands of \mathbf{W} directly in the spatial domain is challenging. For instance, \mathbf{W}_b would need to be infinitely large to isolate the low-frequency part of \mathbf{W} , since the ideal low-pass filter *sinc* has infinite support in spatial domain [15].

To overcome this limitation, we perform the convolution in the Fourier domain rather than in the spatial domain. According to the Convolution Theorem [15], convolution in the spatial domain is equivalent to pointwise multiplication of the Fourier transforms in the frequency domain. Therefore, we formulate the frequency-specific convolution as:

$$\mathbf{Y}_b = \mathcal{F}^{-1}((\mathcal{M}_b \odot \mathcal{F}(\mathbf{W})) \odot \mathcal{F}(\mathbf{X})). \quad (7)$$

This formulation enables the efficient computation of convolutions for each frequency band.

Spatially Variant Modulation. After obtaining the output results for each frequency band, a modulation map $\mathbf{A}_b \in \mathbb{R}^{h \times w}$ is generated to control the influence of each frequency band at each spatial location. \mathbf{A}_b can be easily implemented using a standard convolution layer followed by a sigmoid function. Consequently, the output feature map \mathbf{Y} is computed as:

$$\mathbf{Y} = \sum_{b=0}^{B-1} (\mathbf{A}_b \odot \mathbf{Y}_b). \quad (8)$$

This approach allows FBM to adjust frequency responses dynamically at each spatial location, enhancing the model’s ability to capture context-specific features across the image effectively.

Practical Implementation. Mathematically, the obtation of \mathbf{Y}_b in Equation (7) is equivalent to:

$$\mathbf{Y}_b = \mathcal{F}^{-1}(\mathcal{F}(\mathbf{W}) \odot (\mathcal{M}_b \odot \mathcal{F}(\mathbf{X}))) = \mathbf{W} * \mathbf{X}_b. \quad (9)$$

The above derivation reveals that decomposing the convolution kernel into frequency bands (*i.e.*, $\mathbf{W}_b = \mathcal{M}_b \odot \mathbf{W}$) is mathematically equivalent to decomposing the input feature map into corresponding frequency components (*i.e.*, $\mathbf{X}_b = \mathcal{M}_b \odot \mathbf{X}$). This equivalence stems from the commutative property of convolution and the linearity of the Fourier transform. Specifically, modulating and convolving frequency-specific weights with the original feature map can be reinterpreted as first filtering the feature map into sub-bands and then convolving with the full kernel:

$$\mathbf{Y} = \sum_{b=0}^{B-1} (\mathbf{A}_b \odot \mathbf{X}_b) * \mathbf{W}. \quad (10)$$

Table 1. Results comparison on the COCO validation set [42]. The numbers in brackets indicate the parameters of the backbone. Additionally, the notation $n \times$ denotes the convolutional parameter budget of each dynamic convolution relative to the standard convolution in the backbone.

Models	Params	FLOPs	AP ^{box}	AP ^{mask}
<i>Faster R-CNN</i>	43.80 _(23.5) M	207.1G	37.2	
+ CondConv _[NIPS2019] (8×) [73]	+90.0M	+0.01G	38.1	-
+ DY-Conv _[ICLR2022] (4×) [63]	+75.3M	+0.16G	38.3	-
+ DCD _[ICLR2021] [37]	+4.3M	+0.13G	38.1	
+ ODConv _[ICLR2021] (4×) [32]	+65.1M	+0.35G	39.2	-
+ FDConv (Ours)	+3.6M	+1.8G	39.4	
<i>Mask R-CNN</i>	46.5 _(23.5) M	260.1	39.6	36.4
+ DY-Conv _[CVPR2020] (4×) [63]	+75.3M	+0.16G	39.6	36.6
+ ODConv _[ICLR2021] (4×) [32]	+65.1M	+0.35G	42.1	38.6
+ KW _[ICML2024] (1×) [31]	+2.5M	-	41.8	38.4
+ KW _[ICML2024] (4×) [31]	+76.5M	-	42.4	38.9
+ FDConv (Ours)	+3.6M	+1.8G	42.4	38.6

This perspective bridges two seemingly distinct paradigms: frequency-adaptive weight decomposition and multi-band feature processing. The equivalent implementation not only circumvents the impracticality of infinite spatial support in ideal frequency filters but also provides implementation flexibility, one can choose to implement frequency decomposition on either weights or features based on computational constraints, while maintaining strict mathematical equivalence through Fourier duality.

4. Experiment

Datasets and Metrics. We evaluate our methods on challenging semantic segmentation datasets, including Cityscapes [10] and ADE20K [79], using mean Intersection over Union (mIoU) for segmentation [1, 6, 13, 45, 48] and Average Precision (AP) for object detection and instance segmentation [22, 59].

Implementation Details. We follow the settings from the original papers for UPerNet [70], Mask2Former [8], MaskDINO [72], Swin Transformer [46], and ConvNeXt [31, 47]. On COCO [41], we adhere to standard practices [21, 57, 68], training detection and segmentation models for 12 epochs (1× schedule). We empirically set the number of weights to 64 for FDConv. More details are described in the supplementary.

5. Main Results

In this section, we evaluate our FDConv on a range of tasks, including object detection, instance segmentation, and semantic segmentation, using standard benchmarks such as COCO [42], ADE20K [79], and Cityscapes [10].

We compare our FDConv with state-of-the-art dynamic convolutional methods, including CondConv [73], DY-

Table 2. Quantitative comparisons on semantic segmentation tasks with UPerNet [70] on the ADE20K validation set.

Method	Params	FLOPs	mIoU	
			SS	MS
ResNet-50 [23]	66M	947G	40.7	41.8
ResNet-101 [23]	85M	1029G	42.9	44.0
R50 + PEDConv _(BMVC2021) [24]	72M	947G	42.8	43.9
R50 + ODConv _(ICLR2022) (4×) [32]	131M	947G	43.3	44.4
R50 + KW _(ICML2024) (4×) [31]	141M	947G	43.5	44.6
R50 + FDConv (Ours)	70M	949G	43.8	44.9

Table 3. Object detection and instance segmentation performance on the COCO dataset [42] with the Mask R-CNN detector [22]. All models are trained with a 1× schedule [21, 68].

Model	Params	FLOPs	AP ^{box}	AP ^{mask}
ConvNeXt-T [47]	48M	262G	43.4	39.7
+ KW _(ICML2024) [31]	52M	262G	44.8	40.6
+ FDConv (Ours)	51M	263G	45.2	40.8
Swin-T [46]	48M	267G	42.7	39.3
+ FDConv (Ours)	51M	268G	44.5	40.5

Table 4. Semantic segmentation results on Cityscapes [10] using the recent state-of-the-art Mask2Former [8].

Model	Backbone	mIoU
Mask2Former _(CVPR2022) [8]	ResNet-50	79.4
+ FDConv (Ours)	ResNet-50	80.4 (+1.0)

Table 5. Semantic segmentation results with recent state-of-the-art large models Mask2Former [8] and MaskDINO [34] on ADE20K. Backbones pre-trained on ImageNet-22K are marked with †.

Model	Backbone	mIoU
Mask2Former _(CVPR2022) [8]	Swin-B [†]	53.9
+ FDConv (Ours)	Swin-B [†]	54.9 (+1.0)
MaskDINO _(CVPR2023) [34]	Swin-L [†]	56.6
+ FDConv (Ours)	Swin-L [†]	57.2 (+0.5)

Conv [63], DCD [37], ODConv [32] and KW [31]. The experiments demonstrate that FDConv not only achieves the highest performance across detection and segmentation tasks but also does so with large reduced parameter overhead. Moreover, FDConv is highly versatile, it can easily combine with state-of-the-art ConvNet models like ConvNeXt [47] and apply to transformer architectures such as Swin-T [46], Mask2Former [8], and MaskDINO [34]. The experimental results demonstrate that FDConv achieves notable improvements over both conventional competitors and state-of-the-art baselines.

Object Detection. Table 1 shows the results obtained by Faster R-CNN with various dynamic convolutional modules. Our FDConv module, despite adding only +3.6M parameters and +1.8G FLOPs, achieves an AP^{box} of 39.4,

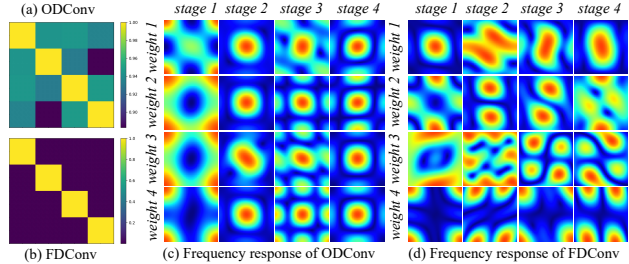


Figure 5. Weight similarity and frequency analyses. (a) demonstrates that existing dynamic convolution methods, such as ODConv [32], exhibit high cosine similarity (>0.88) among their 4 learned weights. The frequency analysis in (c) shows 4 representative ODConv layers from stage 1 to stage 4 of the model, and it demonstrates large homogeneity between the 4 weights. In contrast, the 4 weights of our proposed FDConv show zero similarity in (b), allowing each kernel to learn distinct and complementary features with diversified frequency response, as shown in (d).

2.2% improvement over the baseline and outperforms CondConv [73], DY-Conv [63], and DCD [37], and ODConv [32], which require substantially higher parameter budgets. FDConv not only surpasses other methods in terms of accuracy but also achieves this with a minimal computational footprint, positioning it as a highly efficient enhancement for object detection tasks.

Instance Segmentation. We further evaluate FDConv using Mask R-CNN [22] as the base model, following [31, 32]. FDConv achieves an AP^{box} of 42.4 and AP^{mask} of 38.6, surpassing or matching recent high-performing methods such as ODConv [32] and KW [31]. Notably, while KW [31] achieves marginally higher segmentation performance, it incurs a 4× increase in parameter cost (+76.5M), whereas FDConv adds only 3.6M.

Semantic Segmentation. As shown in Table 2, FDConv achieves the highest mIoU scores, with a single-scale (SS) mIoU of 43.8. Notably, FDConv accomplishes this performance with fewer additional parameters (70M total) compared to ODConv [32] (131M) and KW [31] (141M), underscoring its parameter efficiency while achieving superior segmentation quality.

Combination with Advanced Architectures. Additionally, we test FDConv on object detection and instance segmentation tasks using the COCO [42] to examine its cross-architecture applicability. Table 3 demonstrates that FDConv outperforms other methods, including KW [31], when applied to both ConvNeXt [47] and Swin Transformer [46] backbones. It achieved an AP^{box} of 45.2 with ConvNeXt-T [47] and 44.5 with Swin-T [46], along with enhanced AP^{mask} scores. These results underscore FDConv’s consistent generalization capabilities across various architectures.

Combination with Heavy Models. To assess the adaptability of our FDConv with advanced architectures, we incorporate FDConv into the state-of-the-art Mask2Former [8]

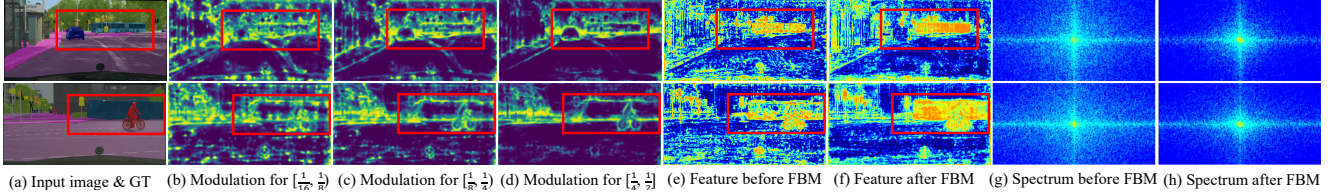


Figure 6. Visualization of Frequency Band Modulation. (a) shows the input images and their corresponding ground truth (GT). (b)–(d) display the modulation maps for different frequency bands, ranging from low to high. (e) and (f) visualize the feature frequency spectrum.

and MaskDINO [34] frameworks. Table 4 shows that Mask2Former-ResNet-50 with FDConv achieves an mIoU improvement of +1.0 (from 79.4 to 80.4) on Cityscapes [10]. On ADE20K [79], Table 5 highlights that with FDConv, Mask2Former-Swin-B [8] achieves an mIoU improvement of +1.0 (from 53.9 to 54.9), while MaskDINO-Swin-L [34] achieves an mIoU improvement of +0.5 (from 56.6 to 57.2). These consistent gains demonstrate that FDConv can effectively enhance heavy architectures.

6. Analyses and Discussion

We use ResNet-50 [76] as the backbone model and conduct a comprehensive analysis of the proposed FDConv. Due to space limitations, more detailed analyses and the results of ablation studies are provided in the *supplementary material*.

Weight Similarity Analysis. We analyze the diversity of learned features in FDConv by comparing weight similarity with existing dynamic convolution methods. As shown in Figure 5(a), traditional methods like ODConv [32] exhibit high cosine similarity (> 0.88) among their learned weights, indicating significant redundancy. This redundancy limits the representational capacity of the model, as each kernel learns overlapping features.

In contrast, FDConv kernels exhibit zero cosine similarity (Figure 5(b)), suggesting that each kernel captures unique, complementary features. This diversity enhances the model’s expressiveness and adaptability.

Weight Frequency Analysis. As shown in Figure 5(c), frequency analysis reveals that ODConv weights exhibit limited frequency diversity across different stages. In contrast, FDConv demonstrates a more diversified frequency response (Figure 5(d)), capturing a broader range of frequency characteristics. This allows FDConv to model a richer set of features, further improving its ability to represent complex input data.

Feature Visualization for FBM. As shown in Figure 6, we visualize the modulation maps for each frequency band. For better performance, we empirically set the modulation map for the lowest frequency band to all 1. We observe that higher modulation values are concentrated around object boundaries, with this effect becoming more pronounced in higher frequency bands (Figure 6(b)–(d)). In contrast, lower

frequency bands exhibit regions of high modulation within the objects themselves (Figure 6(c)).

This selective modulation enables FDConv to suppress high frequencies in regions such as the background and object centers, which do not contribute significantly to accurate predictions. As seen in Figure 6(e)–(f), high-frequency noise in the feature map is largely reduced, and the spectrum in Figure 6(g)–(h) further confirms the suppression of unnecessary high-frequency components. Meanwhile, as shown in Figure 6(e)–(f), foreground features are enhanced, leading to more accurate and complete representations that benefit dense prediction tasks.

7. Conclusion

We introduced Frequency Dynamic Convolution (FDConv), which enhances the frequency adaptability of parallel weights without increasing parameter overhead. By incorporating Fourier Disjoint Weight (FDW), Kernel Spatial Modulation (KSM), and Frequency Band Modulation (FBM), FDConv addresses the limitations of existing dynamic convolution methods, including restricted frequency diversity in parallel weights and high parameter costs.

Our analysis shows that FDConv achieves greater frequency diversity, enabling better feature capture across spatial and frequency domains. Extensive experiments on object detection, segmentation, and classification demonstrate that FDConv outperforms prior state-of-the-art methods, with only a modest increase in parameter cost compared to others that incur much higher overhead. FDConv can be easily integrated into existing architectures, including both ConvNets and vision transformers, making it a versatile and efficient solution for a wide range of computer vision tasks. We hope our analyses and finding would new direction for building more efficient and powerful vision models.

Acknowledgements

This work was supported by the National Key R&D Program of China (2022YFC3300704), the National Natural Science Foundation of China (62331006, 62171038, and 62088101), the Fundamental Research Funds for the Central Universities, and the JST Moonshot R&D Grant Number JPMJMS2011, Japan.

References

- [1] Linwei Chen, Zheng Fang, and Ying Fu. Consistency-aware map generation at multiple zoom levels using aerial image. *IEEE Journal of Selected Topics in Applied Earth Observations and Remote Sensing*, 15:5953–5966, 2022. 1, 6
- [2] Linwei Chen, Ying Fu, Lin Gu, Chenggang Yan, Tatsuya Harada, and Gao Huang. Frequency-aware feature fusion for dense image prediction. *IEEE Transactions Pattern Analysis and Machine Intelligence*, 1(1):1–18, 2024. 2, 3, 5
- [3] Linwei Chen, Ying Fu, Kaixuan Wei, Dezhi Zheng, and Felix Heide. Instance segmentation in the dark. *International Journal of Computer Vision*, 131(8):2198–2218, 2023. 1, 2, 5
- [4] Linwei Chen, Ying Fu, Shaodi You, and Hongzhe Liu. Efficient hybrid supervision for instance segmentation in aerial images. *Remote Sensing*, 13(2):252, 2021.
- [5] Linwei Chen, Ying Fu, Shaodi You, and Hongzhe Liu. Hybrid supervised instance segmentation by learning label noise suppression. *Neurocomputing*, 496:131–146, 2022.
- [6] Linwei Chen, Lin Gu, and Ying Fu. When semantic segmentation meets frequency aliasing. In *Proceedings of International Conference on Learning Representations*, 2024. 1, 3, 6
- [7] Linwei Chen, Lin Gu, Dezhi Zheng, and Ying Fu. Frequency-adaptive dilated convolution for semantic segmentation. In *Proceedings of IEEE Conference on Computer Vision and Pattern Recognition*, pages 3414–3425, 2024. 3
- [8] Bowen Cheng, Ishan Misra, Alexander G Schwing, Alexander Kirillov, and Rohit Girdhar. Masked-attention mask transformer for universal image segmentation. In *Proceedings of IEEE Conference on Computer Vision and Pattern Recognition*, pages 1290–1299, 2022. 6, 7, 8
- [9] Lu Chi, Borui Jiang, and Yadong Mu. Fast fourier convolution. In *Proceedings of Advances in Neural Information Processing Systems*, volume 33, pages 4479–4488, 2020. 3
- [10] Marius Cordts, Mohamed Omran, Sebastian Ramos, Timo Rehfeld, Markus Enzweiler, Rodrigo Benenson, Uwe Franke, Stefan Roth, and Bernt Schiele. The cityscapes dataset for semantic urban scene understanding. In *Proceedings of IEEE Conference on Computer Vision and Pattern Recognition*, pages 3213–3223, 2016. 6, 7, 8
- [11] 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 IEEE Conference on Computer Vision and Pattern Recognition*, pages 5513–5524, 2024. 1
- [12] 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 16x16 words: Transformers for image recognition at scale. In *Proceedings of International Conference on Learning Representations*, pages 1–12, 2020. 4
- [13] Ying Fu, Zheng Fang, Linwei Chen, Tao Song, and Defu Lin. Level-aware consistent multilevel map translation from satellite imagery. *IEEE Transactions on Geoscience and Remote Sensing*, 61:1–14, 2022. 6
- [14] Ying Fu, Hongrong Liu, Yunhao Zou, Shuai Wang, Zhongxiang Li, and Dezhi Zheng. Category-level band learning based feature extraction for hyperspectral image classification. *IEEE Transactions on Geoscience and Remote Sensing*, 62:1–16, 2023. 1
- [15] Rafael C Gonzalez. *Digital image processing*. Pearson education india, 2009. 3, 6
- [16] Julia Grabinski, Steffen Jung, Janis Keuper, and Margret Keuper. Frequencylowcut pooling-plugin and play against catastrophic overfitting. In *Proceedings of European Conference on Computer Vision*, pages 36–57, 2022. 3, 6
- [17] Julia Grabinski, Janis Keuper, and Margret Keuper. Fix your downsampling asap! be natively more robust via aliasing and spectral artifact free pooling. *arXiv preprint arXiv:2307.09804*, 2023. 3
- [18] Qi Guan, Zihao Sheng, and Shibe Xue. Hrpose: Real-time high-resolution 6d pose estimation network using knowledge distillation. *Chinese Journal of Electronics*, 32(1):189–198, 2023. 1
- [19] John Guibas, Morteza Mardani, Zongyi Li, Andrew Tao, Anima Anandkumar, and Bryan Catanzaro. Adaptive fourier neural operators: Efficient token mixers for transformers. In *Proceedings of International Conference on Learning Representations*, pages 1–12, 2022. 3
- [20] David Ha, Andrew Dai, and Quoc V Le. Hypernetworks. In *Proceedings of International Conference on Learning Representations*, 2016. 2
- [21] Ali Hassani and Humphrey Shi. Dilated neighborhood attention transformer. 2022. 6, 7
- [22] Kaiming He, Georgia Gkioxari, Piotr Dollár, and Ross Girshick. Mask r-cnn. In *Proceedings of IEEE International Conference on Computer Vision*, pages 2961–2969, 2017. 6, 7
- [23] Kaiming He, Xiangyu Zhang, Shaoqing Ren, and Jian Sun. Deep residual learning for image recognition. In *Proceedings of IEEE Conference on Computer Vision and Pattern Recognition*, pages 770–778, 2016. 1, 7
- [24] Zejiang Hou and Sun-Yuan Kung. Parameter efficient dynamic convolution via tensor decomposition. In *Proceedings of the British Machine Vision Conference*, page 107, 2021. 3, 7
- [25] Jie Hu, Li Shen, Samuel Albanie, Gang Sun, and Andrea Vedaldi. Gather-excite: Exploiting feature context in convolutional neural networks. *Proceedings of Advances in Neural Information Processing Systems*, 31, 2018. 2
- [26] Jie Hu, Li Shen, and Gang Sun. Squeeze-and-excitation networks. In *Proceedings of IEEE Conference on Computer Vision and Pattern Recognition*, pages 7132–7141, 2018. 2, 3, 5
- [27] Zhipeng Huang, Zhizheng Zhang, Cuiling Lan, Zheng-Jun Zha, Yan Lu, and Baining Guo. Adaptive frequency filters as efficient global token mixers. In *Proceedings of IEEE International Conference on Computer Vision*, pages 1–11, 2023. 3
- [28] Xu Jia, Bert De Brabandere, Tinne Tuytelaars, and Luc V Gool. Dynamic filter networks. *Proceedings of Advances in Neural Information Processing Systems*, 29:1–9, 2016. 2
- [29] Zeqiang Lai, Ying Fu, and Jun Zhang. Hyperspectral image super resolution with real unaligned rgb guidance. *IEEE*

- Transactions on Neural Networks and Learning Systems*, 2024. [1](#)
- [30] HyunJae Lee, Hyo-Eun Kim, and Hyeonseob Nam. Srm: A style-based recalibration module for convolutional neural networks. In *Proceedings of the IEEE/CVF International conference on computer vision*, pages 1854–1862, 2019. [2](#)
- [31] Chao Li and Anbang Yao. Kernelwarehouse: Rethinking the design of dynamic convolution. In *Proceedings of International Conference on Machine Learning*, 2024. [1](#), [2](#), [3](#), [5](#), [6](#), [7](#)
- [32] Chao Li, Aojun Zhou, and Anbang Yao. Omni-dimensional dynamic convolution. In *Proceedings of International Conference on Learning Representations*, 2022. [1](#), [2](#), [3](#), [4](#), [5](#), [6](#), [7](#), [8](#)
- [33] Duo Li, Jie Hu, Changhu Wang, Xiangtai Li, Qi She, Lei Zhu, Tong Zhang, and Qifeng Chen. Involution: Inverting the inheritance of convolution for visual recognition. In *Proceedings of IEEE Conference on Computer Vision and Pattern Recognition*, pages 12321–12330, 2021. [2](#)
- [34] Feng Li, Hao Zhang, Huaizhe Xu, Shilong Liu, Lei Zhang, Lionel M Ni, and Heung-Yeung Shum. Mask dino: Towards a unified transformer-based framework for object detection and segmentation. In *Proceedings of IEEE Conference on Computer Vision and Pattern Recognition*, pages 3041–3050, 2023. [7](#), [8](#)
- [35] Miaoyu Li, Ying Fu, Tao Zhang, and Guanghui Wen. Supervise-assisted self-supervised deep-learning method for hyperspectral image restoration. *IEEE Transactions on Neural Networks and Learning Systems*, 2024. [1](#)
- [36] Xiangtai Li, Jiangning Zhang, Yibo Yang, Guangliang Cheng, Kuiyuan Yang, Yunhai Tong, and Dacheng Tao. Sfnet: Faster and accurate semantic segmentation via semantic flow. *International Journal of Computer Vision*, pages 1–24, 2023. [2](#)
- [37] Yunsheng Li, Yinpeng Chen, Xiyang Dai, Dongdong Chen, Ye Yu, Lu Yuan, Zicheng Liu, Mei Chen, Nuno Vasconcelos, et al. Revisiting dynamic convolution via matrix decomposition. In *Proceedings of International Conference on Learning Representations*, 2021. [3](#), [6](#), [7](#)
- [38] Yutong Li, Miao Ma, Shichang Liu, Chao Yao, and Longjiang Guo. Yolo-drone: a scale-aware detector for drone vision. *Chinese Journal of Electronics*, 33(4):1034–1045, 2024. [1](#)
- [39] Zongyi Li, Nikola Kovachki, Kamyar Azizzadenesheli, Burigede Liu, Kaushik Bhattacharya, Andrew Stuart, and Anima Anandkumar. Fourier neural operator for parametric partial differential equations. In *Proceedings of International Conference on Learning Representations*, pages 1–12, 2021. [3](#)
- [40] Shiqi Lin, Zhizheng Zhang, Zhipeng Huang, Yan Lu, Cuiling Lan, Peng Chu, Quanzeng You, Jiang Wang, Zicheng Liu, Amey Parulkar, et al. Deep frequency filtering for domain generalization. In *Proceedings of IEEE Conference on Computer Vision and Pattern Recognition*, pages 11797–11807, 2023. [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. [6](#)
- [42] 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 *Proceedings of European Conference on Computer Vision*, pages 740–755, 2014. [6](#), [7](#)
- [43] Xudong Lin, Lin Ma, Wei Liu, and Shih-Fu Chang. Context-gated convolution. In *Proceedings of European Conference on Computer Vision*, pages 701–718. Springer, 2020. [3](#)
- [44] Qiankun Liu, Yichen Li, Yuqi Jiang, and Ying Fu. Siamese-detr for generic multi-object tracking. *IEEE Transactions on Image Processing*, 33:3935–3949, 2024. [1](#)
- [45] Songlin Liu, Linwei Chen, Li Zhang, Jun Hu, and Ying Fu. A large-scale climate-aware satellite image dataset for domain adaptive land-cover semantic segmentation. *ISPRS Journal of Photogrammetry and Remote Sensing*, 205:98–114, 2023. [6](#)
- [46] 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 IEEE International Conference on Computer Vision*, pages 10012–10022, 2021. [4](#), [6](#), [7](#)
- [47] Zhuang Liu, Hanzi Mao, Chao-Yuan Wu, Christoph Feichtenhofer, Trevor Darrell, and Saining Xie. A convnet for the 2020s. In *Proceedings of IEEE Conference on Computer Vision and Pattern Recognition*, pages 11976–11986, 2022. [1](#), [6](#), [7](#)
- [48] Jonathan Long, Evan Shelhamer, and Trevor Darrell. Fully convolutional networks for semantic segmentation. In *Proceedings of IEEE International Conference on Computer Vision*, pages 3431–3440, 2015. [6](#)
- [49] Jovita Lukasik, Paul Gavrikov, Janis Keuper, and Margret Keuper. Improving native cnn robustness with filter frequency regularization. *Transactions on Machine Learning Research*, 2023:1–36, 2023. [3](#)
- [50] Ningning Ma, Xiangyu Zhang, Jiawei Huang, and Jian Sun. Weightnet: Revisiting the design space of weight networks. In *Proceedings of European Conference on Computer Vision*, pages 776–792. Springer, 2020. [3](#)
- [51] Salma Abdel Magid, Yulun Zhang, Donglai Wei, Won-Dong Jang, Zudi Lin, Yun Fu, and Hanspeter Pfister. Dynamic high-pass filtering and multi-spectral attention for image super-resolution. In *Proceedings of IEEE International Conference on Computer Vision*, pages 4288–4297, 2021. [2](#), [5](#)
- [52] Michael Mathieu, Mikael Henaff, and Yann LeCun. Fast training of convolutional networks through ffts. In *Proceedings of International Conference on Learning Representations*, pages 1–9, 2014. [5](#)
- [53] Ben Mildenhall, Jonathan T Barron, Jiawen Chen, Dillon Sharlet, Ren Ng, and Robert Carroll. Burst denoising with kernel prediction networks. In *Proceedings of IEEE Conference on Computer Vision and Pattern Recognition*, pages 2502–2510, 2018. [2](#)
- [54] Ioannis Pitas. *Digital image processing algorithms and applications*. John Wiley & Sons, 2000. [3](#)
- [55] Zequn Qin, Pengyi Zhang, Fei Wu, and Xi Li. Fcanet: Frequency channel attention networks. In *Proceedings of IEEE International Conference on Computer Vision*, pages 783–792, 2021. [3](#)

- [56] Niamul Quader, Md Mafijul Islam Bhuiyan, Juwei Lu, Peng Dai, and Wei Li. Weight excitation: Built-in attention mechanisms in convolutional neural networks. In *Proceedings of European Conference on Computer Vision*. Springer, 2020. [3](#)
- [57] Yongming Rao, Wenliang Zhao, Yansong Tang, Jie Zhou, Ser Nam Lim, and Jiwen Lu. Hornet: Efficient high-order spatial interactions with recursive gated convolutions. *Proceedings of Advances in Neural Information Processing Systems*, 35:10353–10366, 2022. [6](#)
- [58] Yongming Rao, Wenliang Zhao, Zheng Zhu, Jiwen Lu, and Jie Zhou. Global filter networks for image classification. In *Proceedings of Advances in Neural Information Processing Systems*, volume 34, pages 980–993, 2021. [3](#)
- [59] Shaoqing Ren, Kaiming He, Ross Girshick, and Jian Sun. Faster r-cnn: Towards real-time object detection with region proposal networks. In *Proceedings of Advances in Neural Information Processing Systems*, pages 91–99, 2015. [6](#)
- [60] Ajay Subramanian, Elena Sizikova, Najib J Majaj, and Dennis G Pelli. Spatial-frequency channels, shape bias, and adversarial robustness. pages 1–10, 2024. [5](#)
- [61] Ye Tian, Ying Fu, and Jun Zhang. Transformer-based under-sampled single-pixel imaging. *Chinese Journal of Electronics*, 32(5):1151–1159, 2023. [1](#)
- [62] Hugo Touvron, Matthieu Cord, Matthijs Douze, Francisco Massa, Alexandre Sablayrolles, and Hervé Jégou. Training data-efficient image transformers & distillation through attention. In *Proceedings of International Conference on Machine Learning*, pages 10347–10357. PMLR, 2021. [1](#)
- [63] Thomas Verelst and Tinne Tuytelaars. Dynamic convolutions: Exploiting spatial sparsity for faster inference. In *Proceedings of IEEE Conference on Computer Vision and Pattern Recognition*, pages 2320–2329, 2020. [1](#), [2](#), [3](#), [4](#), [5](#), [6](#), [7](#)
- [64] Fei Wang, Mengqing Jiang, Chen Qian, Shuo Yang, Cheng Li, Honggang Zhang, Xiaogang Wang, and Xiaoou Tang. Residual attention network for image classification. In *Proceedings of IEEE Conference on Computer Vision and Pattern Recognition*, pages 3156–3164, 2017. [2](#)
- [65] Haohan Wang, Xindi Wu, Zeyi Huang, and Eric P Xing. High-frequency component helps explain the generalization of convolutional neural networks. In *Proceedings of IEEE Conference on Computer Vision and Pattern Recognition*, pages 8684–8694, 2020. [3](#)
- [66] Jiaqi Wang, Kai Chen, Rui Xu, Ziwei Liu, Chen Change Loy, and Dahua Lin. Carafe: Content-aware reassembly of features. In *Proceedings of IEEE International Conference on Computer Vision*, pages 3007–3016, 2019. [2](#)
- [67] Qilong Wang, Banggu Wu, Pengfei Zhu, Peihua Li, Wangmeng Zuo, and Qinghua Hu. Eca-net: Efficient channel attention for deep convolutional neural networks. In *Proceedings of IEEE Conference on Computer Vision and Pattern Recognition*, pages 11534–11542, 2020. [2](#), [5](#)
- [68] 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 IEEE Conference on Computer Vision and Pattern Recognition*, pages 14408–14419, 2023. [6](#), [7](#)
- [69] Sanghyun Woo, Jongchan Park, Joon-Young Lee, and In So Kweon. Cbam: Convolutional block attention module. In *Proceedings of European Conference on Computer Vision*, pages 3–19, 2018. [2](#)
- [70] Tete Xiao, Yingcheng Liu, Bolei Zhou, Yuning Jiang, and Jian Sun. Unified perceptual parsing for scene understanding. In *Proceedings of European Conference on Computer Vision*, pages 418–434, 2018. [6](#), [7](#)
- [71] Yuwen Xiong, Zhiqi Li, Yuntao Chen, Feng Wang, Xizhou Zhu, Jiapeng Luo, Wenhai Wang, Tong Lu, Hongsheng Li, Yu Qiao, et al. Efficient deformable convnets: Rethinking dynamic and sparse operator for vision applications. In *Proceedings of IEEE Conference on Computer Vision and Pattern Recognition*, pages 5652–5661, 2024. [1](#)
- [72] Jiacong Xu, Zixiang Xiong, and Shankar P Bhattacharyya. Pidnet: A real-time semantic segmentation network inspired by pid controllers. In *Proceedings of IEEE Conference on Computer Vision and Pattern Recognition*, pages 19529–19539, 2023. [6](#)
- [73] Brandon Yang, Gabriel Bender, Quoc V Le, and Jiquan Ngiam. Condconv: Conditionally parameterized convolutions for efficient inference. *Proceedings of Advances in Neural Information Processing Systems*, 32, 2019. [1](#), [2](#), [3](#), [4](#), [5](#), [6](#), [7](#)
- [74] Lingxiao Yang, Ru-Yuan Zhang, Lida Li, and Xiaohua Xie. Simam: A simple, parameter-free attention module for convolutional neural networks. In *Proceedings of International Conference on Machine Learning*, pages 11863–11874. PMLR, 2021. [2](#)
- [75] Dong Yin, Raphael Gontijo Lopes, Jon Shlens, Ekin Dogus Cubuk, and Justin Gilmer. A fourier perspective on model robustness in computer vision. In *Proceedings of Advances in Neural Information Processing Systems*, volume 32, 2019. [3](#)
- [76] Fisher Yu, Vladlen Koltun, and Thomas Funkhouser. Dilated residual networks. In *Proceedings of IEEE Conference on Computer Vision and Pattern Recognition*, pages 472–480, 2017. [8](#)
- [77] Tao Zhang, Ying Fu, Jun Zhang, and Chenggang Yan. Deep guided attention network for joint denoising and demosaicing in real image. *Chinese Journal of Electronics*, 33(1):303–312, 2024. [1](#)
- [78] Zhenli Zhang, Xiangyu Zhang, Chao Peng, Xiangyang Xue, and Jian Sun. Exfuse: Enhancing feature fusion for semantic segmentation. In *Proceedings of European Conference on Computer Vision*, pages 269–284, 2018. [2](#)
- [79] Bolei Zhou, Hang Zhao, Xavier Puig, Sanja Fidler, Adela Barriuso, and Antonio Torralba. Scene parsing through ade20k dataset. In *Proceedings of IEEE Conference on Computer Vision and Pattern Recognition*, pages 633–641, 2017. [6](#), [8](#)
- [80] Yunhao Zou, Ying Fu, Tsuyoshi Takatani, and Yinqiang Zheng. Eventhdr: From event to high-speed hdr videos and beyond. *IEEE Transactions on Pattern Analysis and Machine Intelligence*, 47(1):32–50, 2024. [1](#)