# From CNNs to Shift-Invariant Twin Models Based on Complex Wavelets

Hubert Leterme[1]     Kévin Polisano[2]     Valérie Perrier[2]     Karteek Alahari[1]

[1]Univ. Grenoble Alpes, CNRS, Inria, Grenoble INP, LJK, 38000 Grenoble, France
[2]Univ. Grenoble Alpes, CNRS, Grenoble INP, LJK, 38000 Grenoble, France

## Abstract

*We propose a novel antialiasing method to increase shift invariance and prediction accuracy in convolutional neural networks. Specifically, we replace the first-layer combination "real-valued convolutions → max pooling" ($\mathbb{R}$Max) by "complex-valued convolutions → modulus" ($\mathbb{C}$Mod), which is stable to translations. To justify our approach, we claim that $\mathbb{C}$Mod and $\mathbb{R}$Max produce comparable outputs when the convolution kernel is band-pass and oriented (Gabor-like filter). In this context, $\mathbb{C}$Mod can be considered as a stable alternative to $\mathbb{R}$Max. Thus, prior to antialiasing, we force the convolution kernels to adopt such a Gabor-like structure. The corresponding architecture is called mathematical twin, because it employs a well-defined mathematical operator to mimic the behavior of the original, freely-trained model. Our antialiasing approach achieves superior accuracy on ImageNet and CIFAR-10 classification tasks, compared to prior methods based on low-pass filtering. Arguably, our approach's emphasis on retaining high-frequency details contributes to a better balance between shift invariance and information preservation, resulting in improved performance. Furthermore, it has a lower computational cost and memory footprint than concurrent work, making it a promising solution for practical implementation.*

## 1. Introduction

Over the past decade, some progress has been made on understanding the strengths and limitations of convolutional neural networks (CNNs) for computer vision [21, 35, 41]. The ability of CNNs to embed input images into a feature space with linearly separable decision regions is a key factor to achieve high classification accuracy. An important property to reach this linear separability is the ability to discard or minimize non-discriminative image components. In particular, feature vectors are expected to be stable with respect to translations [35]. However, subsampling opera-

tions, typically found in convolution and pooling layers, are an important source of instability—a phenomenon known as aliasing [1]. A few approaches have attempted to address this issue.

**Blurpooled CNNs.** Zhang [42] proposed to apply a low-pass *blurring* filter before each subsampling operation in CNNs. Specifically, 1. max pooling layers (Max → Sub)[1] are replaced by max-blur pooling (Max → Blur → Sub); 2. convolution layers followed by ReLU (Conv → Sub → ReLU) are blurred before subsampling (Conv → ReLU → Blur → Sub).[2] The combination Blur → Sub is referred to as *blur pooling*. This approach follows a well-known practice in signal processing, which involves low-pass filtering a high-frequency signal before subsampling, in order to avoid artifacts in reconstruction. Their approach improved the shift invariance as well as the accuracy of CNNs trained on ImageNet and CIFAR-10 datasets. However, this was achieved with a significant loss of information.

A question then arises: is it possible to design a non-destructive antialiasing method, and if so, does it further improve accuracy? In a more recent work, Zou *et al.* [45] tackled this question through an adaptive antialiasing approach, called *adaptive blur pooling*. In this nonlinear setting, the blurring filter varies across channels and image locations, therefore preserving high-frequency information in strategic zones. Albeit achieving higher prediction accuracy, this approach remains fundamentally based on low-pass filtering. Consequently, features that are not blurred may still be unstable to translations. Furthermore, adaptive blur pooling requires additional memory, computational resources, and trainable parameters.

**Proposed Approach.** In this paper, we propose an alternative antialiasing approach based on complex-valued con-

---

[1]Sub and Conv stand for "convolution" and "subsampling" respectively.

[2]ReLU is computed before blurring; otherwise the network would simply perform on low-resolution images.

volutions, extracting high-frequency features that are stable to translations. We observed improved accuracy for ImageNet and CIFAR-10 classification, compared to the two antialiasing methods based on blur pooling [42, 45]. Furthermore, our approach offers significant advantages in terms of computational efficiency and memory usage, and does not induce any additional training, unlike adaptive blur pooling.

Our proposed method replaces the first layers of a CNN:

$$\text{Conv} \rightarrow \text{Sub} \rightarrow \text{Bias} \rightarrow \text{ReLU} \rightarrow \text{MaxPool}, \quad (1)$$

equivalently written as

$$\text{Conv} \rightarrow \text{Sub} \rightarrow \text{MaxPool} \rightarrow \text{Bias} \rightarrow \text{ReLU}, \quad (2)$$

by the following combination:

$$\mathbb{C}\text{Conv} \rightarrow \text{Sub} \rightarrow \text{Modulus} \rightarrow \text{Bias} \rightarrow \text{ReLU}, \quad (3)$$

where $\mathbb{C}$Conv denotes a convolution operator with a complex-valued kernel, whose real and imaginary parts approximately form a 2D Hilbert transform pair [11]. From (2) and (3), we introduce the two following operators:

$$\mathbb{R}\text{Max} \; : \; \text{Conv} \rightarrow \text{Sub} \rightarrow \text{MaxPool}; \quad (4)$$

$$\mathbb{C}\text{Mod} \; : \; \mathbb{C}\text{Conv} \rightarrow \text{Sub} \rightarrow \text{Modulus}. \quad (5)$$

Our method is motivated by the following theoretical claim. In a recent paper submission [23], we proved that 1. $\mathbb{C}$Mod is nearly invariant to translations, if the convolution kernel is band-pass and clearly oriented; 2. $\mathbb{R}$Max and $\mathbb{C}$Mod produce comparable outputs, except for some filter frequencies regularly scattered across the Fourier domain. We then combined these two properties to establish a stability metric for $\mathbb{R}$Max as a function of the convolution kernel's frequency vector. This work was essentially theoretical, with limited experiments conducted on a deterministic model solely based on the dual-tree complex wavelet packet transform (DT-$\mathbb{C}$WPT). However, it lacked applications to tasks such as image classification. Building upon this theoretical study, in this paper, we consider the $\mathbb{C}$Mod operator as a proxy for $\mathbb{R}$Max, extracting comparable, yet more stable features.

In compliance with the theory, the $\mathbb{R}$Max-$\mathbb{C}$Mod substitution is only applied to the output channels associated with oriented band-pass filters, referred to as *Gabor-like kernels*. This kind of structure is known to arise spontaneously in the first layer of CNNs trained on image datasets such as ImageNet [38]. In this paper, we enforce this property by applying additional constraints to the original model, prior to antialiasing. Specifically, a predefined number of convolution kernels are guided to adopt Gabor-like structures, instead of letting the network learn them from scratch. For this purpose, we rely on the dual-tree complex wavelet packet

transform (DT-$\mathbb{C}$WPT) [2]. Throughout the paper, we refer to this constrained model as a *mathematical twin*, because it employs a well-defined mathematical operator to mimic the behavior of the original model. In this context, replacing $\mathbb{R}$Max by $\mathbb{C}$Mod is straightforward, since the complex-valued filters are provided by DT-$\mathbb{C}$WPT.

**Other Related Work.** Following the ideas developed for antialiasing, Chaman and Dokmanic [5] reached perfect shift invariance by using an adaptive, input-dependent subsampling grid, whereas previous models rely on fixed grids. This idea was harnessed by Xu *et al.* [37] to get shift equivariance in generative models. This approach is not intended to compete with other antialiasing methods, but rather to complement them at the subsampling stages.

Another aspect of shift invariance in CNNs is related to boundary effects. The fact that CNNs can encode the absolute position of an object in the image by exploiting boundary effects was discovered independently by Islam *et al.* [15], and Kayhan and Gemert [17]. This phenomenon is left outside the scope of our paper.

Wavelet scattering networks (ScatterNets), by Bruna and Mallat [4], perform cascading wavelet convolutions and nonlinear operations. They produce shift-invariant image representations that are stable to deformation and preserve high-frequency information. Further extensions include roto-translation invariant ScatterNets [27], hybrid Scatter-Nets combined with fully-trained layers [26] or dictionary learning [40], dual-tree complex wavelet ScatterNets [31], graph ScatterNets [44], learnable ScatterNets via feature map mixing [6] or parametric wavelet filters [9]. As deep learning architectures with well-defined mathematical properties, ScatterNets may be used as explanatory models for standard, freely-trained networks. However, there is no exact correspondence between the two types of architectures [34]. In contrast, our models are enhanced versions of existing networks, rather than ad hoc constructions. This allows to draw comparisons in terms of prediction accuracy and shift invariance.

Finally, we draw the reader's attention to the family of complex-valued CNNs ($\mathbb{C}$VCNNs), which is the focus of a comprehensive survey [22]. These networks are well suited for tasks requiring the phase information to be propagated through the entire network, as done in the context of magnetic resonance imaging [7], polarimetric imaging [43] or audio signals [33]. However, for image recognition tasks, $\mathbb{C}$VCNNs do not seem to perform better than standard CNNs, with equal number of trainable parameters [33]. Conversely, our approach discards the phase information by computing the modulus. In summary, $\mathbb{C}$VCNNs and our models, although both employing complex-valued convolutions, are suited for different contexts.

## 2. Proposed Approach

We first describe the general principles of our antialiasing approach based on complex convolutions. We then provide some details about the mathematical twin based on DT-$\mathbb{C}$WPT, and explain how our method has been benchmarked against blur-pooling-based antialiased models.

We represent feature maps as 2D sequences, defined by straight capital letters: $X \in \mathcal{S}$. Indexing is denoted by square brackets: for any 2D index $\boldsymbol{n} \in \mathbb{Z}^2$, $X[\boldsymbol{n}] \in \mathbb{R}$ or $\mathbb{C}$. The cross-correlation between X and $V \in \mathcal{S}$ is defined by $(X \star V)[\boldsymbol{n}] := \sum_{\boldsymbol{k} \in \mathbb{Z}^2} X[\boldsymbol{n} + \boldsymbol{k}] V[\boldsymbol{k}]$. The down arrow refers to subsampling: for any $m \in \mathbb{N}^*$, $(X \downarrow m)[\boldsymbol{n}] := X[m\boldsymbol{n}]$.

### 2.1. Standard Architectures

A convolution layer with $K$ input channels, $L$ output channels and subsampling factor $m \in \mathbb{N} \setminus \{0\}$ is parameterized by a weight tensor $\mathbf{V} := (V_{lk})_{l \in \{1..L\}, \, k \in \{1..K\}} \in \mathcal{S}^{L \times K}$. For any multichannel input $\mathbf{X} := (X_k)_{k \in \{1..K\}} \in \mathcal{S}^K$, the corresponding output $\mathbf{Y} := (Y_l)_{l \in \{1..L\}} \in \mathcal{S}^L$ is defined such that, for any output channel $l \in \{1..L\}$,

$$Y_l := \sum_{k=1}^{K} (X_k \star V_{lk}) \downarrow m. \qquad (6)$$

For instance, in AlexNet and ResNet, $K = 3$ (RGB input images), $L = 64$, and $m = 4$ and 2, respectively. Next, a bias $\boldsymbol{b} := (b_1, \cdots, b_L)^\top \in \mathbb{R}^L$ is applied to $\mathbf{Y}$, which is then transformed through nonlinear ReLU and max pooling operators. The activated outputs satisfy

$$A_l^{\max} := \text{MaxPool}\left(\text{ReLU}(Y_l + b_l)\right), \qquad (7)$$

where we have defined, for any $Y \in \mathcal{S}$ and any $\boldsymbol{n} \in \mathbb{Z}^2$,

$$\text{ReLU}(Y)[\boldsymbol{n}] := \max(0, Y[\boldsymbol{n}]); \qquad (8)$$

$$\text{MaxPool}(Y)[\boldsymbol{n}] := \max_{\|\boldsymbol{k}\|_\infty \leq 1} Y[2\boldsymbol{n} + \boldsymbol{k}]. \qquad (9)$$

### 2.2. Antialiasing Principle

We consider the first convolution layer of a CNN, as described in (6). As widely discussed in the literature [38], after training with ImageNet, a certain number of convolution kernels $V_{lk}$ spontaneously take the appearance of oriented waveforms with well-defined frequency and orientation (Gabor-like kernels), as illustrated in Fig. 1a. A visual representation of trained convolution kernels is provided in Appendix E. In this paper, we refer to the corresponding output channels $l \in \mathcal{G} \subset \{1..L\}$ as *Gabor channels*. The main idea is to substitute, for any $l \in \mathcal{G}$, $\mathbb{R}$Max by $\mathbb{C}$Mod, as explained hereafter. Following (2), expression (7) can be rewritten

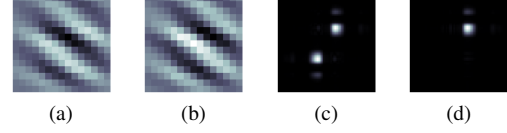$$A_l^{\max} = \text{ReLU}\left(Y_l^{\max} + b_l\right), \qquad (10)$$



Figure 1. (a), (b): Real and imaginary parts of a Gabor-like convolution kernel $W_{lk} := V_{lk} + i\mathcal{H}(V_{lk})$, forming a 2D Hilbert transform pair. (c), (d): Power spectra (energy of the Fourier transform) of $V_{lk}$ and $W_{lk}$, respectively.

where $Y_l^{\max}$ is the output of an $\mathbb{R}$Max operator as introduced in (4). More formally,

$$Y_l^{\max} := \text{MaxPool}\left(\sum_{k=1}^{K}(X_k \star V_{lk}) \downarrow m\right). \qquad (11)$$

Then, following (3), the $\mathbb{R}$Max-$\mathbb{C}$Mod substitution yields

$$A_l^{\text{mod}} = \text{ReLU}\left(Y_l^{\text{mod}} + b_l\right), \qquad (12)$$

where $Y_l^{\text{mod}}$ is the output of a $\mathbb{C}$Mod operator (5), satisfying

$$Y_l^{\text{mod}} := \left| \sum_{k=1}^{K}(X_k \star W_{lk}) \downarrow (2m) \right|. \qquad (13)$$

In the above expression, $W_{lk}$ is a complex-valued analytic kernel satisfying

$$W_{lk} := V_{lk} + i\mathcal{H}(V_{lk}), \qquad (14)$$

where $\mathcal{H}$ denotes the two-dimensional *Hilbert transform* as introduced by Havlicek *et al.* [11]. It is defined, for any real-valued filter $V \in \mathcal{S}$ and any frequency vector $\boldsymbol{\xi} := (\xi_1, \xi_2)^\top \in [-\pi, \pi]^2$, by

$$\widehat{\mathcal{H}V}(\boldsymbol{\xi}) := -i \, \text{sgn}(\xi_1) \, \widehat{V}(\boldsymbol{\xi}). \qquad (15)$$

The Hilbert transform is designed such that the Fourier transform of $W_{lk}$ is entirely supported in the half-plane of nonnegative $x$-values. Therefore, since $V_{lk}$ is Gabor-like, the energy of $W_{lk}$ is concentrated within a small window in the Fourier domain, as depicted in Fig. 1d. Due to this property, the modulus operator provides a smooth envelope for complex-valued cross-correlations with $W_{lk}$ [19]. This leads to the output $Y_l^{\text{mod}}$ (13) being nearly invariant to translations. Additionally, the subsampling factor in (13) is twice that in (11), to account for the factor-2 subsampling achieved through max pooling (9).

### 2.3. Wavelet-Based Twin Models (WCNNs)

As explained in Sec. 2.2, introducing an imaginary part to the Gabor-like convolution kernels improves shift invariance. Our antialiasing method therefore restricts to the Gabor channels $l \in \mathcal{G} \subset \{1..L\}$. However, $\mathcal{G}$ is unknown a

priori: for a given output channel $l \in \{1 .. L\}$, whether $V_{lk}$ will become band-pass and oriented after training is unpredictable. Thus, we need a way to automatically separate the set $\mathcal{G}$ of Gabor channels from the set of remaining channels, denoted by $\mathcal{F} := \{1 .. L\} \setminus \mathcal{G}$. To this end, we built "mathematical twins" of standard CNNs, based on the dual-tree wavelet packet transform (DT-$\mathbb{C}$WPT). These models, which we call WCNNs in short, reproduce the behavior of freely-trained architectures with a higher degree of control and fewer trainable parameters. In this section, we present their general structure; a more detailed description is provided in Appendix A. For the purpose of readability, we assume that $K = 3$ (RGB input images).

We denote by $L_{\text{gabor}} := \text{card}(\mathcal{G})$ and $L_{\text{free}} := \text{card}(\mathcal{F})$ the number of Gabor and remaining channels, respectively. They are determined empirically from the trained CNNs—see Tab. 1. In a twin WCNN architecture, the two groups of output channels are organized such that $\mathcal{F} = \{1 .. L_{\text{free}}\}$ and $\mathcal{G} = \{(L_{\text{free}} + 1) .. L\}$. The first $L_{\text{free}}$ channels, which are outside the scope of our antialiasing approach, remain freely-trained as in the standard architecture. Regarding the $L_{\text{gabor}}$ Gabor channels, the convolution kernels $V_{lk}$ are constrained to satisfy the following requirements. First, all three RGB input channels are processed with the same filter, up to a multiplicative constant. More formally, there exists a *luminance* weight vector $\boldsymbol{\mu} := (\mu_1,\, \mu_2,\, \mu_3)^\top$, with $\mu_k \in [0,\, 1]$ and $\sum_{k=1}^{3} \mu_k = 1$, such that,

$$\forall k \in \{1 .. 3\},\; V_{lk} = \mu_k \widetilde{V}_l, \tag{16}$$

where $\widetilde{V}_l := \sum_{k=1}^{3} V_{lk}$ denotes the mean kernel. Furthermore, $\widetilde{V}_l$ must be band-pass and oriented (Gabor-like filter). The following paragraphs explain how these two constraints are implemented in our WCNN architecture.

**Monochrome Filters.** Expression (16) is actually a property of standard CNNs: the oriented band-pass RGB kernels generally appear monochrome (see kernel visualization of freely-trained CNNs Appendix E). A numerical assessment of this property can by found in [23]. In WCNNs, this constraint is implemented with a trainable $1 \times 1$ convolution layer [24], parameterized by $\boldsymbol{\mu}$, computing the following luminance image:

$$X^{\text{lum}} := \sum_{k=1}^{3} \mu_k X_k. \tag{17}$$

**Gabor-Like Kernels.** To guarantee the Gabor-like property on $\widetilde{V}_l$, we implemented DT-$\mathbb{C}$WPT, which is achieved through a series of subsampled convolutions. The number of decomposition stages $J \in \mathbb{N} \setminus \{0\}$ was chosen such that $m = 2^{J-1}$, where, as a reminder, $m$ denotes the subsampling factor as introduced in (6). DT-$\mathbb{C}$WPT generates a set

of filters $\left(W_{k'}^{\text{dt}}\right)_{k' \in \{1 .. 4 \times 4^J\}}$, which tiles the Fourier domain $[-\pi,\, \pi]^2$ into $4 \times 4^J$ overlapping square windows.[3] Their real and imaginary parts approximately form a 2D Hilbert transform pair such as defined in (15).

The WCNN architecture is designed such that, for any Gabor channel $l \in \mathcal{G}$, $\widetilde{V}_l$ is the real part of one such filter:

$$\exists k' \in \left\{1 .. 4 \times 4^J\right\} : \widetilde{V}_l = \text{Re}\left(W_{k'}^{\text{dt}}\right). \tag{18}$$

The output $Y_l$ introduced in (6) then becomes

$$Y_l = \left(X^{\text{lum}} \star \widetilde{V}_l\right) \downarrow 2^{J-1}. \tag{19}$$

To summarize, a WCNN substitutes the freely-trained convolution (6) with a combination of (17) and (19), for any Gabor output channels $l \in \mathcal{G}$. This combination is wrapped into a *wavelet block*, also referred to as WBlock in short. Technical details about its exact design are provided Appendix A. Note that the Fourier resolution of $V_{lk}$ increases with the subsampling factor $m$. This property is consistent with what is observed in freely-trained CNNs: in AlexNet, where $m = 4$, the Gabor-like filters are more localized in frequency (and less spatially localized) than in ResNet, where $m = 2$.

Figure 2 displays the kernels $\mathbf{V} \in \mathcal{S}^{L \times K}$, with $K = 3$ and $L = 64$, for the WCNN architectures based on AlexNet and ResNet, which we refer to as WAlexNet and WResNet, respectively. In these models, $m = 4$ and 2, which implies $J = 3$ and 2, respectively. The kernels are shown as RGB color images, after training with ImageNet, for both freely-trained and Gabor channels.

**Antialiased WCNNs.** Using the antialiasing principles presented in Sec. 2.2, we replace $\mathbb{R}$Max (11) by $\mathbb{C}$Mod (13) for all Gabor channels $l \in \mathcal{G}$. In the corresponding model, referred to as $\mathbb{C}$WCNN, the wavelet block is replaced by a *complex wavelet block* ($\mathbb{C}$WBlock), in which (19) becomes

$$Z_l = \left(X^{\text{lum}} \star \widetilde{W}_l\right) \downarrow 2^J, \tag{20}$$

where $\widetilde{W}_l$ is obtained by considering both real and imaginary parts of the DT-$\mathbb{C}$WPT filter (18). Then, a modulus is applied to $Z_l$, which yields $Y_l^{\text{mod}}$ such as defined in (13), with $W_{lk} := \mu_k \widetilde{W}_l$ for any RGB channel $k \in \{1 .. 3\}$. Finally, we apply a bias and ReLU to $Y_l^{\text{mod}}$, following (12).

A schematic representation of WAlexNet and its antialiased version, referred to as $\mathbb{C}$WAlexNet, is provided in Figs. 3b and 3c (top).

### 2.4. WCNNs with Blur Pooling

We benchmark our approach against the antialiasing methods proposed by Zhang [42] and Zou *et al.* [45].

---

[3]Figure 1 actually illustrates one of these filters, with $J = 3$.

(a) WAlexNet (32 Gabor channels)
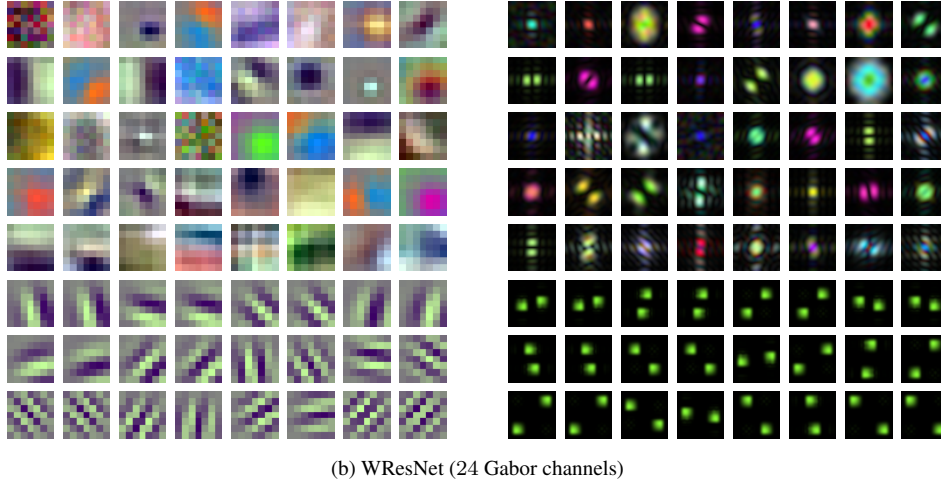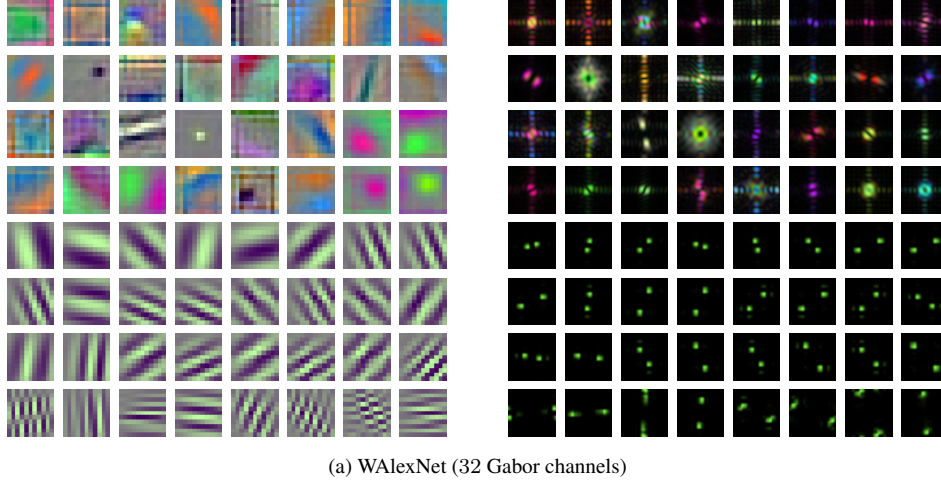


(b) WResNet (24 Gabor channels)

Figure 2. Convolution kernels in the first layer of WAlexNet (a) and WResNet-34 (b), after training with ImageNet-1K. For each output channel $l \in \{1 \dots 64\}$, the corresponding convolution kernel $(V_{lk})_{k \in \{1..3\}}$ is displayed as an RGB image in the spatial domain (left), and its associated magnitude spectrum in the Fourier domain (right). The convolution kernels associated to the Gabor channels are displayed on the 4 and 3 last rows for WAlexNet and WResNet, respectively. For the sake of visual rendering, they have been respectively cropped to $(11 \times 11)$ and $(7 \times 7)$ to match the size of the freely-trained kernels.

To this end, we first consider a WCNN antialiased with static or adaptive blur pooling, respectively referred to as BlurWCNN and ABlurWCNN. A schematic representation of BlurWAlexNet is provided in Fig. 3b (bottom). Then, we substitute the blurpooled Gabor channels (right branch of the diagram) with our own $\mathbb{C}$Mod-based approach. The corresponding models are respectively referred to as $\mathbb{C}$BlurWCNN and $\mathbb{C}$ABlurWCNN. Again, a schematic representation of $\mathbb{C}$BlurWAlexNet can be found in Fig. 3c (bottom). Note that, for a fair comparison, all models use blur pooling in the freely-trained channels as well as deeper layers. Therefore, several antialiasing methods are employed in different parts of the network.

## 2.5. Adaptation to ResNet: Batch Normalization

In many architectures including ResNet, the bias is computed after an operation called *batch normalization* (BN) [14]. In this context, (1) becomes

$$\text{Conv} \rightarrow \text{Sub} \rightarrow \text{BN} \rightarrow \text{Bias} \rightarrow \text{ReLU} \rightarrow \text{MaxPool}. \quad (21)$$

As detailed in Appendix B, the $\mathbb{R}$Max-$\mathbb{C}$Mod substitution yields, analogously to (3),

$$\mathbb{C}\text{Conv} \rightarrow \text{Sub} \rightarrow \text{Modulus} \rightarrow \text{BN0} \rightarrow \text{Bias} \rightarrow \text{ReLU}, \quad (22)$$

where BN0 refers to a special type of batch normalization without mean centering. A schematic representation of ResNet-based models, as done in Fig. 3 for AlexNet, is provided in Fig. 8.

**(a) AlexNet**

**(b) WAlexNet (baseline)**

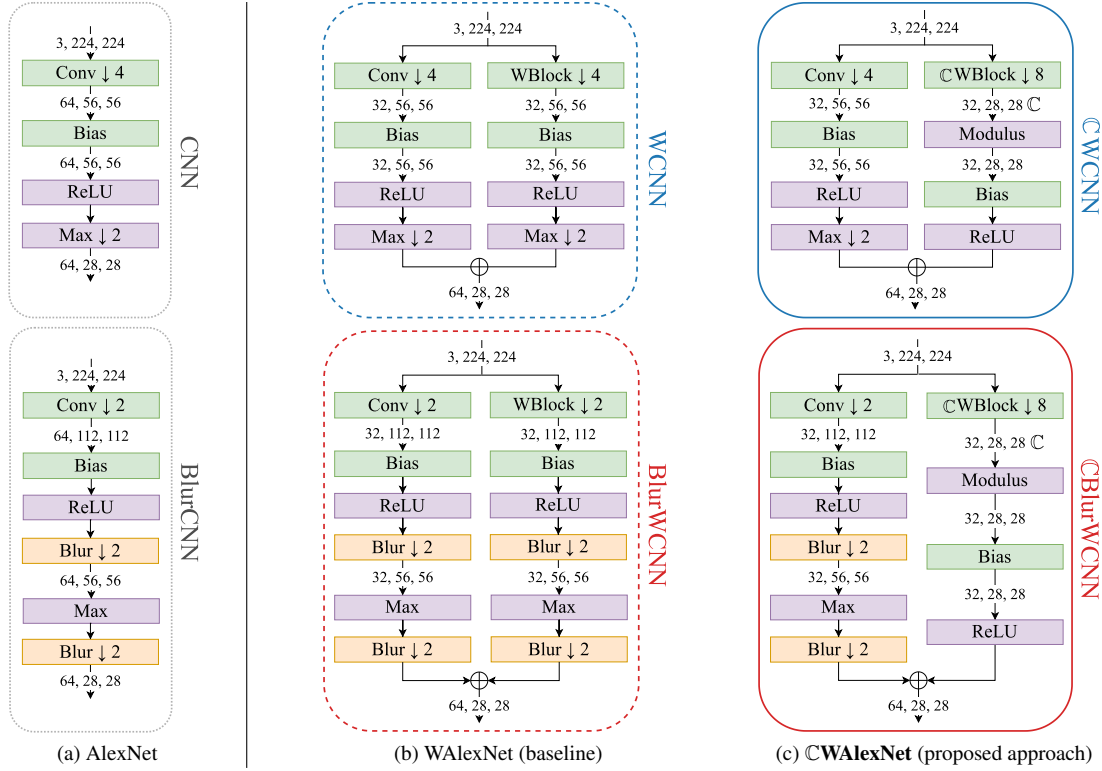**(c) $\mathbb{C}$WAlexNet (proposed approach)**

Figure 3. First layers of AlexNet and its variants, corresponding to a convolution layer followed by ReLU and max pooling (1). The models are framed according to the same colors and line styles as in Figs. 4 and 5. The green modules are the ones containing trainable parameters; the orange and purple modules represent static linear and nonlinear operators, respectively. The numbers between each module represent the depth (number of channels), height and width of each output. Fig. 3a: freely-trained models. Top: standard AlexNet. Bottom: Zhang's "blurpooled" AlexNet. Fig. 3b: mathematical twins (WAlexNet) reproducing the behavior of standard (top) and blurpooled (bottom) AlexNet. The left side of each diagram corresponds to the $L_{\text{free}} := 32$ freely-trained output channels, whereas the right side displays the $L_{\text{gabor}} := 32$ remaining channels, where freely-trained convolutions have been replaced by a wavelet block (WBlock) as described in Sec. 2.3. Fig. 3c: $\mathbb{C}$Mod-based antialiased WAlexNet, where WBlock has been replaced by $\mathbb{C}$WBlock, and max pooling by a modulus. The bias and ReLU are placed after the modulus, following (3). In the bottom models, we compare Zhang's antialiasing approach (Fig. 3b) with ours (Fig. 3c) in the Gabor channels.

|  | WAlexNet | WResNet |
|---|---|---|
| $m$ (subsampling factor) | 4 | 2 |
| $J$ (decomposition depth) | 3 | 2 |
| $L_{\text{free}}$, $L_{\text{gabor}}$ (output channels) | 32, 32 | 40, 24 |

Table 1. Experimental settings for our WCNN twin models. Other details are provided in Appendix C.

## 3. Experiments

### 3.1. Experiment Details

**ImageNet.** We built our WCNN and $\mathbb{C}$WCNN twin models based on AlexNet [20] and ResNet-34 [12]. Their overall design is described in Sec. 2, along with setting details in Tab. 1. The values of $L_{\text{free}}$ and $L_{\text{gabor}}$ were determined empirically from the freely-trained AlexNet and ResNet-34; further details are provided in Appendix C. Zhang's static blur pooling approach is tested on both AlexNet and ResNet, whereas Zou *et al.*'s adaptive approach is only tested on ResNet. The latter was indeed not implemented on AlexNet in the original paper, and we could not make it work on this architecture.

As mentioned above, we compare blur-pooling-based antialiasing approach (Fig. 3b, bottom) with ours (Fig. 3c, bottom). To apply static or adaptive blur pooling to the WC-NNs, we proceed as follows. Following Zhang's implementation, the wavelet block is not antialiased if $m = 2$ as in ResNet, for computational reasons. However, when $m = 4$ as in AlexNet, a blur pooling layer is placed after ReLU, and the wavelet block's subsampling factor is divided by 2. Moreover, max pooling is replaced by max-blur pooling. The size of the blurring filters is set to 3, as recommended by Zhang [42]. Besides, DT-$\mathbb{C}$WPT decompositions are performed with Q-shift orthogonal filters of length

10 as introduced by Kingsbury [18].

Our models were trained on the ImageNet ILSVRC2012 dataset [29], following the standard procedure provided by PyTorch [28].[4] Moreover, we set aside 100K images from the training set—100 per class—in order to compute the top-1 error rate after each training epoch ("validation set").

**CIFAR-10.** We also trained ResNet-18, ResNet-34 and their variants on the CIFAR-10 dataset. Training was performed on 300 epochs, with an initial learning rate equal to 0.1, decreased by a factor of 10 every 100 epochs. As for ImageNet, we set aside $5\,000$ images out of 50K to compute accuracy during the training phase. Given the images of small size in this dataset ($32 \times 32$ pixels), feature extraction can be performed efficiently with a reduced number of layers. For this reason, the first layers (1) arguably have a higher influence on the overall predictive power. We therefore expect to clearly highlight the benefits of our approach on this specific task.

## 3.2. Evaluation Metrics

**Classification Accuracy.** Classification accuracy was computed on the standard ImageNet evaluation set (50K images). We followed the *ten-crops* procedure [20]: predictions are made over 10 patches extracted from each input image, and the softmax outputs are averaged to get the overall prediction. We also considered center crops of size 224 for *one-crop* evaluation. In both cases, we used top-1-5 error rates. For CIFAR-10 evaluation (10K images), we measured the top-1 error rate with one- and ten-crops.

**Measuring Shift Invariance.** For each image in the ImageNet evaluation set, we extracted several patches of size 224, each of which being shifted by $0.5$ pixel along a given axis. We then compared their outputs in order to measure the model's robustness to shifts. This was done by computing the Kullback-Leibler (KL) divergence between output vectors—which, under certain hypotheses, can be interpreted as probability distributions [3, pp. 205-206]. This metric is intended for visual representation.

In addition, we measured the mean flip rate (mFR) between predictions [13], as done by Zhang [42] in its blur-pooled models. For each direction (vertical, horizontal and diagonal), we measured the mean frequency upon which two shifted input images yield different top-1 predictions, for shift distances varying from 1 to 8 pixels. We then normalized the results with respect to AlexNet's mFR, and averaged over the three directions. This metric is also referred to as *consistency*.

We repeated the procedure for the models trained on CIFAR-10. This time, we extracted patches of size $32 \times 32$
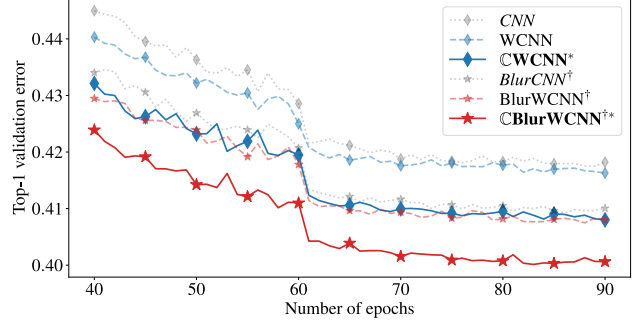
Figure 4. Evolution of the top-1 validation error (one-crop) along training with ImageNet, for AlexNet-based models, as described in Fig. 3. The freely-trained models, upon which the mathematical twins are built, appear in faint gray. Legend: [†]blur pooling; [*]ℂMod-based antialiasing (our approach).
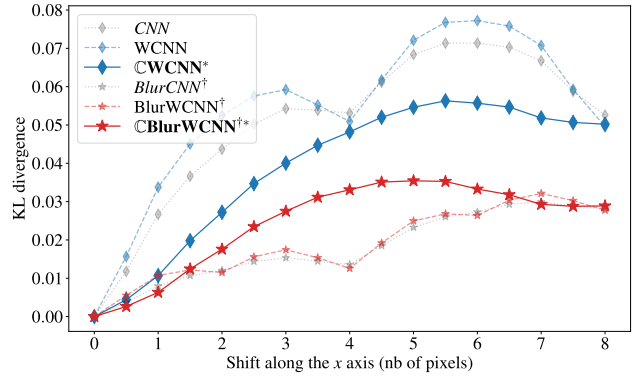


Figure 5. AlexNet-based models: mean KL divergence between the outputs of a reference image versus shifted images. Legend: [†]blur pooling; [*]ℂMod-based antialiasing (our approach).

from the evaluation set, and computed mFR for shifts varying from 1 to 4 pixels. Besides, normalization was performed with respect to ResNet-18's mFR.

## 3.3. Results and Discussion

**Validation and Test Accuracy.** Top-1 accuracy of AlexNet-based models along training with ImageNet is plotted in Fig. 4. Similar plots are provided for ResNet in Fig. 10. In addition, error rates of AlexNet- and ResNet-based architectures, computed on the evaluation sets, are provided in Tab. 2 for ImageNet and Tab. 3 for CIFAR-10.

Our ℂMod-based approach significantly outperforms the baselines for AlexNet: ℂWCNN vs WCNN (blue diamonds), and ℂBlurWCNN vs BlurWCNN (red stars). Remarkably, the exclusive application of ℂMod-based antialiasing to the Gabor channels (ℂWCNN, solid blue line) is sufficient to match the performance of blur-pooling-based antialiasing (BlurWCNN, dashed red line), which, in contrast, is implemented throughout the entire network. Positive results are also obtained for ResNet-based models

| Model | One-crop | | Ten-crops | | Shifts |
|---|---|---|---|---|---|
| | top-1 | top-5 | top-1 | top-5 | mFR |
| **AlexNet** | | | | | |
| *CNN* | *45.3* | *22.2* | *41.3* | *19.3* | *100.0* |
| WCNN | 44.9 | 21.8 | 40.8 | 19.0 | 101.4 |
| ℂWCNN* | **44.3** | **21.3** | **40.2** | **18.5** | **88.0** |
| *BlurCNN†* | *44.4* | *21.6* | *40.7* | *18.7* | *63.8* |
| BlurWCNN† | 44.3 | 21.4 | 40.5 | 18.5 | **63.1** |
| ℂBlurWCNN†* | **43.3** | **20.5** | **39.6** | **17.9** | 69.4 |
| **ResNet-34** | | | | | |
| *CNN* | *27.6* | *9.2* | *24.8* | *7.7* | *78.1* |
| WCNN | 27.4 | 9.2 | 24.7 | 7.6 | 77.2 |
| ℂWCNN* | **27.2** | **9.0** | **24.4** | **7.4** | **73.1** |
| *BlurCNN†* | *26.7* | *8.6* | *24.0* | *7.2* | *61.2* |
| BlurWCNN† | 26.7 | 8.6 | 24.1 | 7.3 | 65.2 |
| ℂBlurWCNN†* | **26.5** | **8.4** | **23.7** | **7.0** | 62.5 |
| *ABlurCNN‡* | *26.1* | *8.3* | *23.5* | *7.0* | *60.8* |
| ABlurWCNN‡ | 26.0 | 8.2 | 23.6 | 6.9 | 62.1 |
| ℂABlurWCNN‡* | 26.1 | **8.2** | 23.7 | 7.0 | 63.1 |

Table 2. Evaluation metrics on ImageNet (%): the lower the better. Models: †static and ‡adaptive blur pooling; *ℂMod-based antialiasing (our approach).

| Model | ResNet-18 | | | ResNet-34 | | |
|---|---|---|---|---|---|---|
| | 1crp | 10crp | shft | 1crp | 10crps | shft |
| *CNN* | *14.9* | *10.8* | *100.0* | *15.2* | *10.9* | *100.3* |
| WCNN | 14.2 | 10.3 | 92.4 | 14.5 | 10.5 | 99.2 |
| ℂWCNN* | **13.8** | **9.6** | **88.8** | **12.9** | **9.2** | **93.0** |
| *BlurCNN†* | *14.2* | *10.4* | *87.7* | *15.7* | *11.6* | *88.2* |
| BlurWCNN† | 13.1 | 9.7 | **84.6** | 13.2 | 9.9 | 85.6 |
| ℂBlurWCNN†* | **12.3** | **8.9** | 85.7 | **12.4** | **9.1** | **83.7** |
| *ABlurCNN‡* | *14.6* | *11.0* | *90.9* | *16.3* | *12.8* | *91.9* |
| ABlurWCNN‡ | 14.5 | 11.0 | 86.5 | 14.0 | 10.4 | 93.3 |
| ℂABlurWCNN‡* | **12.8** | **9.7** | **81.7** | **12.8** | **9.2** | **86.6** |

Table 3. Evaluation metrics on CIFAR-10 (%): top-1 error rate using one- and ten-crops methods ("1crp" and "10crp"); and mFR measuring consistency ("shft"). Models: †static and ‡adaptive blur pooling; *ℂMod-based antialiasing (our approach).

**Shift Invariance (KL Divergence).** The mean KL divergence between outputs of shifted images are plotted in Fig. 5 for AlexNet trained on ImageNet. Moreover, the mean flip rate for shifted inputs (consistency) is reported in Tab. 2 for ImageNet (AlexNet and ResNet-34) and Tab. 3 for CIFAR-10 (ResNet-18 and 34).

In both AlexNet and WAlexNet, the initial convolution layer's output undergoes a one-pixel shift for every four-pixel shift in the input image. Consequently, any divergence between the output vectors is due to the instability of subsequent layers to one-pixel shifts. In contrast, instabilities which are accountable to the initial layer are observed for shifts that are *not* multiples of $4$. Likewise, for input shifts of eight pixels, the max pooling's output is shifted by exactly one pixel, resulting in even higher stability. In ℂWAlexNet, the same eight-to-one-pixel ratio occurs to the modulus layer's output, which explains why the two curves meet for 8-pixel shifts. However, the ℝMax-ℂMod substitution has greatly reduced first-layer instabilities, resulting in a flattened curve and avoiding the "bumps" observed for non-antialiased models.

On the other hand, BlurAlexNet and BlurWAlexNet exhibit considerably flattened curves compared to non-antialiased models, as well as ℂWAlexNet. This demonstrates the effectiveness of blur-pooling-based method in enhancing shift invariance. Applying ℂMod-based antialiasing instead of blur pooling on the Gabor channels (ℂBlurWAlexNet) actually degrades shift invariance (except for shifts smaller than $1.5$ pixels), as evidenced by the bell-shaped curve. Nevertheless, the corresponding classifier is significantly more accurate. This is not surprising, as our approach prioritizes the conservation of high-frequency details, which are important for classification. An extreme reduction of shift variance using a large blur pooling filter would indeed result in a significant loss of accuracy. Therefore, our work achieves a better balance between shift invariance and information preservation. To gain further insights into this phenomenon, we conducted experiments by varying the size of the blurring filters in AlexNet-based models. Figure 6 shows the relationship between consistency and prediction accuracy on ImageNet (custom validation set), for different filter sizes ranging from 1 (no blur pooling) to 7 (heavy loss of high-frequency information).[5] We find that a near-optimal trade-off is achieved when the filter size is set to 2 or 3. Furthermore, at equivalent consistency levels, ℂBlurWCNN outperforms BlurWCNN in terms of accuracy. Note however that the optimal version of ℂBlurWCNN is not more consistent than the optimal version of BlurWCNN, despite achieving higher accuracy.

**Computational Resources.** Table 4 compares the computational resources and memory footprint required for each

trained on ImageNet (Tab. 2, bottom). However, adaptive blur pooling, when applied to the Gabor channels (ABlurWCNN), yields similar or marginally higher accuracy than our approach (ℂABlurWCNN). Nevertheless, our method is computationally more efficient, requires less memory (see "Computational Resources" below for more details), and do not demand additional training, unlike adaptive blur pooling. To better support this claim, we conducted ablation studies which we present in Appendix D.3. Finally, when trained on CIFAR-10 (see Tab. 3), our ℂMod-based antialiased models built upon ResNet-18 and 34 achieve significant gains in accuracy over non-antialiased models, as well as models antialiased with both blur-pooling-based methods.

As a side note, the training curves for WCNNs (colored dashed lines) closely follow those of standard CNNs (gray dotted lines). This is an expected result, since the former models are designed to mimic the behavior of the latter.

---

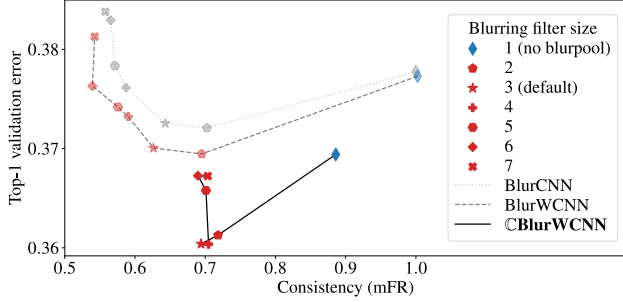[5]Similar plots can be found in Zhang's paper [42].

Figure 6. Classification accuracy (ten-crops) vs consistency, measuring the stability of predictions to small input shifts, for AlexNet-based models (the lower the better for both axes). For each of the three architectures, we increased the blurring filter size from 1 (*i.e.*, no blur pooling) to 7. The blue diamonds (no blur pooling) and red stars (blur pooling with filters of size 3) correspond to the models from Figs. 4 and 5. At equivalent consistency levels, our ℂMod-based approach (solid line) yields higher accuracy.

| Method | Computational cost | | Memory footprint | |
|---|---|---|---|---|
| | AlexNet | ResNet | AlexNet | ResNet |
| *No antialiasing (ref)* | *1.0* | *1.0* | *1.0* | *1.0* |
| BlurPool [42] | 4.0 | 1.0 | 4.7 | 1.9 |
| ABlurPool [45] | – | 2.1 | – | 2.0 |
| **ℂMod (ours)** | **0.5** | **0.5** | **0.6** | **0.4** |

Table 4. Computational cost and memory footprint required for each antialiasing method, per Gabor channel. The values are normalized relative to non-antialiased AlexNet or ResNet. Computational cost: FLOPs for computing $Y_l^{\max}$ (11) or $Y_l^{\mod}$ (13). Memory footprint: size of the intermediate and output tensors saved by PyTorch for the backward pass.

antialiasing method. To achieve this, we considered models with freely-trained convolutions, because the goal was to evaluate the computational performances of the various antialiasing approaches, excluding implementation tricks based on DT-ℂWPT from the scope of analysis. More details are provided in Appendices F and G.

## 4. Conclusion

The mathematical twins introduced in this paper serve a proof of concept for our ℂMod-based antialiasing approach. However, its range of application extends well beyond DT-ℂWPT filters. While we focused on the first convolution layer, it is important to note that such initial layers play a critical role in CNNs by extracting low-level geometric features such as edges, corners or textures. Therefore, a specific attention is required for their design. In contrast, deeper layers are more focused on capturing high-level structures that conventional image processing tools are poorly suited for, as pointed out by Oyallon *et al.* in the context of hybrid scattering networks [26].

Furthermore, while our approach is tailored for CNN architectures, which were chosen to make a fair comparison to related methods, it has potential for broader applicability. There is a growing interest in using self-attention mechanisms in computer vision [8] to capture complex, long-range dependencies among image representations. Recent work on vision transformers has proposed using the first layers of a CNN as a "convolutional token embedding" [10, 36, 39], effectively reintroducing inductive biases to the architecture, such as locality and weight sharing. By applying our antialiasing method to this embedding, we can provide self-attention modules with nearly shift-invariant inputs, which could be highly beneficial in improving the performance of vision transformers, especially when the amount of available data is limited.

## References

[1] Aharon Azulay and Yair Weiss. Why do deep convolutional networks generalize so poorly to small image transformations? *Journal of Machine Learning Research*, 20(184):1–25, 2019. 1

[2] Ilker Bayram and Ivan W. Selesnick. On the Dual-Tree Complex Wavelet Packet and M-Band Transforms. *IEEE Transactions on Signal Processing*, 56(6):2298–2310, June 2008. 2, 14

[3] Christopher M. Bishop and Tom M. Mitchell. *Pattern Recognition and Machine Learning*. Springer, 2014. 7

[4] Joan Bruna and Stéphane Mallat. Invariant Scattering Convolution Networks. *IEEE Transactions on Pattern Analysis and Machine Intelligence*, 35(8):1872–1886, May 2013. 2

[5] Anadi Chaman and Ivan Dokmanic. Truly Shift-Invariant Convolutional Neural Networks. In *Proc. IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 3773–3783, 2021. 2

[6] Fergal Cotter and Nick Kingsbury. A Learnable Scatternet: Locally Invariant Convolutional Layers. In *Proc. IEEE International Conference on Image Processing (ICIP)*, pages 350–354, Sept. 2019. 2

[7] Muneer Ahmad Dedmari, Sailesh Conjeti, Santiago Estrada, Phillip Ehses, Tony Stöcker, and Martin Reuter. Complex Fully Convolutional Neural Networks for MR Image Reconstruction. In Florian Knoll, Andreas Maier, and Daniel

Rueckert, editors, *Machine Learning for Medical Image Reconstruction*, Lecture Notes in Computer Science, pages 30–38, Cham, 2018. Springer International Publishing. 2

[8] Alexey Dosovitskiy, Lucas Beyer, Alexander Kolesnikov, Dirk Weissenborn, Xiaohua Zhai, Thomas Unterthiner, Mostafa Dehghani, Matthias Minderer, Georg Heigold, Sylvain Gelly, Jakob Uszkoreit, and Neil Houlsby. An Image is Worth 16x16 Words: Transformers for Image Recognition at Scale. In *Proc. International Conference on Learning Representations*, 2021. 9

[9] Shanel Gauthier, Benjamin Thérien, Laurent Alsène-Racicot, Muawiz Chaudhary, Irina Rish, Eugene Belilovsky, Michael Eickenberg, and Guy Wolf. Parametric Scattering Networks. In *Proc. IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 5749–5758, 2022. 2

[10] Ali Hassani, Steven Walton, Nikhil Shah, Abulikemu Abuduweili, Jiachen Li, and Humphrey Shi. Escaping the Big Data Paradigm with Compact Transformers. *arXiv:2104.05704*, June 2022. 9

[11] J.P. Havlicek, J.W. Havlicek, and A.C. Bovik. The analytic image. In *Proc. International Conference on Image Processing*, volume 2, pages 446–449 vol.2, Oct. 1997. 2, 3

[12] Kaiming He, Xiangyu Zhang, Shaoqing Ren, and Jian Sun. Deep residual learning for image recognition. In *Proc. IEEE Computer Society Conference on Computer Vision and Pattern Recognition*, 2016. 6

[13] Dan Hendrycks and Thomas Dietterich. Benchmarking Neural Network Robustness to Common Corruptions and Perturbations. In *Proc. International Conference on Learning Representations*, Mar. 2019. 7

[14] Sergey Ioffe and Christian Szegedy. Batch Normalization: Accelerating Deep Network Training by Reducing Internal Covariate Shift. In *Proc. 32nd International Conference on Machine Learning*, pages 448–456. PMLR, June 2015. 5

[15] Md Amirul Islam, Sen Jia, and Neil D. B. Bruce. How Much Position Information Do Convolutional Neural Networks Encode? In *Proc. International Conference on Learning Representations*, Jan. 2020. 2

[16] Bernd Jahne. *Practical Handbook on Image Processing for Scientific and Technical Applications*. CRC Press, 2004. 15

[17] Osman Semih Kayhan and Jan C. van Gemert. On Translation Invariance in CNNs: Convolutional Layers Can Exploit Absolute Spatial Location. In *Proc. IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 14274–14285, 2020. 2

[18] Nick Kingsbury. Design of Q-shift complex wavelets for image processing using frequency domain energy minimization. In *Proceedings International Conference on Image Processing*, volume 1, pages I–1013, 2003. 7

[19] Nick Kingsbury and Julian Magarey. Wavelet Transforms in Image Processing. In Ales Procházka, Jan Uhlíř, P. W. J. Rayner, and N. G. Kingsbury, editors, *Signal Analysis and Prediction*, Applied and Numerical Harmonic Analysis, pages 27–46. Birkhäuser, Boston, MA, 1998. 3

[20] Alex Krizhevsky, Ilya Sutskever, and Geoffrey E. Hinton. ImageNet classification with deep convolutional neural networks. *Communications of the ACM*, 60(6):84–90, May 2017. 6, 7

[21] Yann LeCun, Yoshua Bengio, and Geoffrey Hinton. Deep learning. *Nature*, 521(7553):436–444, 2015. 1

[22] ChiYan Lee, Hideyuki Hasegawa, and Shangce Gao. Complex-Valued Neural Networks: A Comprehensive Survey. *IEEE/CAA Journal of Automatica Sinica*, 9(8):1406–1426, Aug. 2022. 2

[23] Hubert Leterme, Kévin Polisano, Valérie Perrier, and Karteek Alahari. On the Shift Invariance of Max Pooling Feature Maps in Convolutional Neural Networks. *arXiv:2209.11740*, Sept. 2022. 2, 4, 12, 14

[24] Min Lin, Qiang Chen, and Shuicheng Yan. Network In Network. *arXiv:1312.4400 [cs]*, 2014. 4

[25] Jun Liu and Jieping Ye. Efficient L1/Lq Norm Regularization. *arXiv:1009.4766*, Sept. 2010. 12

[26] Edouard Oyallon, Eugene Belilovsky, and Sergey Zagoruyko. Scaling the Scattering Transform: Deep Hybrid Networks. In *Proc. IEEE International Conference on Computer Vision*, pages 5618–5627, 2017. 2, 9

[27] Edouard Oyallon and Stéphane Mallat. Deep Roto-Translation Scattering for Object Classification. In *Proc. IEEE Conference on Computer Vision and Pattern Recognition*, pages 2865–2873, 2015. 2

[28] Adam Paszke, Sam Gross, Soumith Chintala, Gregory Chanan, Edward Yang, Zachary DeVito, Zeming Lin, Alban Desmaison, Luca Antiga, and Adam Lerer. Automatic differentiation in PyTorch. Oct. 2017. 7

[29] Olga Russakovsky, Jia Deng, Hao Su, Jonathan Krause, Sanjeev Satheesh, Sean Ma, Zhiheng Huang, Andrej Karpathy, Aditya Khosla, Michael Bernstein, Alexander C. Berg, and Li Fei-Fei. ImageNet Large Scale Visual Recognition Challenge. *International Journal of Computer Vision*, 115(3):211–252, Apr. 2015. 7

[30] Ivan W. Selesnick, Richard Baraniuk, and Nick Kingsbury. The dual-tree complex wavelet transform. *IEEE Signal Processing Magazine*, 22(6):123–151, Nov. 2005. 11

[31] Amarjot Singh and Nick Kingsbury. Dual-Tree wavelet scattering network with parametric log transformation for object classification. In *ICASSP, IEEE International Conference on Acoustics, Speech and Signal Processing - Proceedings*, 2017. 2

[32] Antonio Torralba and Aude Oliva. Statistics of natural image categories. *Network: Computation in Neural Systems*, 14(3):391–412, Jan. 2003. 12

[33] Chiheb Trabelsi, Olexa Bilaniuk, Ying Zhang, Dmitriy Serdyuk, Sandeep Subramanian, João Felipe Santos, Soroush Mehri, Negar Rostamzadeh, Yoshua Bengio, and Christopher J. Pal. Deep Complex Networks. In *International Conference on Learning Representations*, Feb. 2018. 2

[34] Mark Tygert, Joan Bruna, Soumith Chintala, Yann LeCun, Serkan Piantino, and Arthur Szlam. A Mathematical Motivation for Complex-Valued Convolutional Networks. *Neural Computation*, 28(5):815–825, May 2016. 2

[35] Thomas Wiatowski and Helmut Bölcskei. A Mathematical Theory of Deep Convolutional Neural Networks for Feature Extraction. *IEEE Transactions on Information Theory*, 64(3):1845–1866, Mar. 2018. 1

[36] Haiping Wu, Bin Xiao, Noel Codella, Mengchen Liu, Xiyang Dai, Lu Yuan, and Lei Zhang. CvT: Introducing Convolutions to Vision Transformers. In *Proc. IEEE/CVF International Conference on Computer Vision*, pages 22–31, 2021. 9

[37] Jin Xu, Hyunjik Kim, Thomas Rainforth, and Yee Teh. Group Equivariant Subsampling. In *Advances in Neural Information Processing Systems*, volume 34, pages 5934–5946, 2021. 2

[38] Jason Yosinski, Jeff Clune, Yoshua Bengio, and Hod Lipson. How transferable are features in deep neural networks? In *Advances in Neural Information Processing Systems*, pages 3320–3328, 2014. 2, 3

[39] Kun Yuan, Shaopeng Guo, Ziwei Liu, Aojun Zhou, Fengwei Yu, and Wei Wu. Incorporating Convolution Designs Into Visual Transformers. In *Proc. IEEE/CVF International Conference on Computer Vision*, pages 579–588, 2021. 9

[40] John Zarka, Louis Thiry, Tomás Angles, and Stéphane Mallat. Deep Network Classification by Scattering and Homotopy Dictionary Learning. In *International Conference on Learning Representations*, 2020. 2

[41] Matthew D. Zeiler and Rob Fergus. Visualizing and understanding convolutional networks. In *Lecture Notes in Computer Science (Including Subseries Lecture Notes in Artificial Intelligence and Lecture Notes in Bioinformatics)*, 2014. 1

[42] Richard Zhang. Making Convolutional Networks Shift-Invariant Again. In *International Conference on Machine Learning*, pages 7324–7334. PMLR, May 2019. 1, 2, 4, 6, 7, 8, 9, 13, 15, 18

[43] Zhimian Zhang, Haipeng Wang, Feng Xu, and Ya-Qiu Jin. Complex-Valued Convolutional Neural Network and Its Application in Polarimetric SAR Image Classification. *IEEE Transactions on Geoscience and Remote Sensing*, 55(12):7177–7188, Dec. 2017. 2

[44] Dongmian Zou and Gilad Lerman. Graph convolutional neural networks via scattering. *Applied and Computational Harmonic Analysis*, 49(3):1046–1074, Nov. 2020. 2

[45] Xueyan Zou, Fanyi Xiao, Zhiding Yu, Yuheng Li, and Yong Jae Lee. Delving Deeper into Anti-Aliasing in ConvNets. *International Journal of Computer Vision*, 131(1):67–81, Jan. 2023. 1, 2, 4, 9, 13, 18

## A. Design of WCNNs: Technical Details

In this section, we provide complements to the description of the mathematical twin (WCNN) introduced in Sec. 2.3. We explained that, for each Gabor channel $l \in \mathcal{G}$, the average kernel $\widetilde{V}_l$ is the real part of a DT-$\mathbb{C}$WPT filter, as described in (18). We now explain how the filter selection is done; in other words, how $k'$ is chosen among $\{1 \ldots 4 \times 4^J\}$. Since input images are real-valued, we restrict to the filters with bandwidth located in the half-plane of positive $x$-values. For the sake of concision, we denote by $K_{\mathrm{dt}} := 2 \times 4^J$ the number of such filters.

For any RGB image $\mathbf{X} \in \mathcal{S}^3$, a luminance image $\mathrm{X}^{\mathrm{lum}} \in \mathcal{S}$ is computed following (17), using a $1 \times 1$ convolution layer. Then, DT-$\mathbb{C}$WPT is performed on $\mathrm{X}^{\mathrm{lum}}$. We denote

by $\mathbf{D} := (\mathrm{D}_k)_{k \in \{1 .. K_{\mathrm{dt}}\}}$ the tensor containing the real part of the DT-$\mathbb{C}$WPT feature maps:

$$\mathrm{D}_k = \left(\mathrm{X}^{\mathrm{lum}} \star \operatorname{Re} \mathrm{W}_k^{(J)}\right) \downarrow 2^{J-1}. \tag{23}$$

For the sake of computational efficiency, DT-$\mathbb{C}$WPT is performed with a succession of subsampled separable convolutions and linear combinations of real-valued wavelet packet feature maps [30]. To match the subsampling factor $m := 2^{J-1}$ of the standard model, the last decomposition stage is performed without subsampling.

**Filter Selection.** The number of dual-tree feature maps $K_{\mathrm{dt}}$ may be greater than the number of Gabor channels $L_{\mathrm{gabor}}$. In that case, we therefore want to select filters that contribute the most to the network's predictive power. First, the low-frequency feature maps $\mathrm{D}_0$ and $\mathrm{D}_{(4^J+1)}$ are discarded. Then, a subset of $K'_{\mathrm{dt}} < K_{\mathrm{dt}}$ feature maps is manually selected and permuted in order to form clusters in the Fourier domain. Considering a (truncated) permutation matrix $\boldsymbol{\Sigma} \in \mathbb{R}^{K'_{\mathrm{dt}} \times K_{\mathrm{dt}}}$, the output of this transformation, denoted by $\mathbf{D}' \in \mathcal{S}^{K'_{\mathrm{dt}}}$, is defined by:

$$\mathbf{D}' := \boldsymbol{\Sigma} \mathbf{D}. \tag{24}$$

The feature maps $\mathbf{D}'$ are then sliced into $Q$ groups of channels $\mathbf{D}^{(q)} \in \mathcal{S}^{K_q}$, each of them corresponding to a cluster of band-pass dual-tree filters with neighboring frequencies and orientations. On the other hand, the output of the wavelet block, $\mathbf{Y}^{\mathrm{gabor}} := (\mathrm{Y}_l)_{l \in \{L_{\mathrm{free}}+1 .. L\}} \in \mathcal{S}^{L_{\mathrm{gabor}}}$, where $\mathrm{Y}_l$ has been introduced in (6), is also sliced into $Q$ groups of channels $\mathbf{Y}^{(q)} \in \mathcal{S}^{L_q}$. Then, for each group $q \in \{1 .. Q\}$, an affine mapping between $\mathbf{D}^{(q)}$ and $\mathbf{Y}^{(q)}$ is performed. It is characterized by a trainable matrix $\boldsymbol{A}^{(q)} := \left(\boldsymbol{\alpha}_1^{(q)}, \cdots, \boldsymbol{\alpha}_{L_q}^{(q)}\right)^\top \in \mathbb{R}^{L_q \times K_q}$ such that, for any $l \in \{1 .. L_q\}$,

$$\mathrm{Y}_l^{(q)} := \boldsymbol{\alpha}_l^{(q)\top} \cdot \mathbf{D}^{(q)}. \tag{25}$$

As in the color mixing stage, this operation is implemented as a $1 \times 1$ convolution layer.

A schematic representation of the real- and complex-valued wavelet blocks can be found in Fig. 7.

**Sparse Regularization.** For any group $q \in \{1 .. Q\}$ and output channel $l \in \{1 .. L_q\}$, we want the model to select one and only one wavelet packet feature map within the $q$-th group. In other words, each row vector $\boldsymbol{\alpha}_l^{(q)} := \left(\alpha_{l,1}^{(q)}, \cdots, \alpha_{l,K_q}^{(q)}\right)^\top$ of $\boldsymbol{A}^{(q)}$ contains no more than one nonzero element, such that (25) becomes
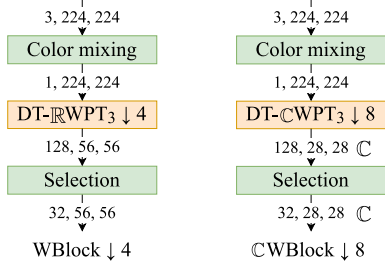
$$\mathrm{Y}_l^{(q)} = \alpha_{lk}^{(q)} \mathrm{D}_k^{(q)} \tag{26}$$

Figure 7. Detail of a wavelet block with $J = 3$ as in AlexNet, in its $\mathbb{R}$Max (left) and $\mathbb{C}$Mod (right) versions. DT-$\mathbb{R}$WPT corresponds to the real part of DT-$\mathbb{C}$WPT.

for some (unknown) value of $k \in \{1 . . K_q\}$. To enforce this property during training, we add a mixed-norm $l^1/l^\infty$-regularizer [25] to the loss function to penalize non-sparse feature map mixing as follows:

$$\mathcal{L} := \mathcal{L}_0 + \sum_{q=1}^{Q} \lambda_q \sum_{l=1}^{L_q} \left( \frac{\|\boldsymbol{\alpha}_l^{(q)}\|_1}{\|\boldsymbol{\alpha}_l^{(q)}\|_\infty} - 1 \right), \quad (27)$$

where $\mathcal{L}_0$ denotes the standard cross-entropy loss and $\boldsymbol{\lambda} \in \mathbb{R}^Q$ denotes a vector of regularization hyperparameters. Note that the unit bias in (27) serves for interpretability of the regularized loss ($\mathcal{L} = \mathcal{L}_0$ in the desired configuration) but has no impact on training.

## B. Batch Normalization in ResNet

This section complements Sec. 2.5 in the main paper. In many recent architectures including ResNet, the bias (see Fig. 3) is replaced by an affine batch normalization layer (BN). In this section, we show how to adapt our approach to this context.

A BN layer is parameterized by trainable weight and bias vectors, respectively denoted by $\boldsymbol{a}$ and $\boldsymbol{b} \in \mathbb{R}^L$. In the remaining of the section, we consider input images $\mathbf{X}$ as a stack of discrete stochastic processes. Then, expression (7) is replaced by

$$A_l := \text{MaxPool}\left\{ \text{ReLU}\left( a_l \cdot \frac{Y_l - \mathbb{E}_m[Y_l]}{\sqrt{\mathbb{V}_m[Y_l] + \varepsilon}} + b_l \right) \right\}, \quad (28)$$

with $Y_l$ satisfying (6) (output of the first convolution layer). In the above expression, we have introduced $\mathbb{E}_m(Y_l) \in \mathbb{R}$ and $\mathbb{V}_m(Y_l) \in \mathbb{R}_+$, which respectively denote the mean expected value and variance of $Y_l[\boldsymbol{n}]$, for indices $\boldsymbol{n}$ contained in the support of $Y_l$, denoted by $\text{supp}(Y_l)$. Let us denote by $N \in \mathbb{N} \setminus \{0\}$ the support size of input images. Therefore, if the filter's support size $N_{\text{filt}}$ is much smaller that $N$, then $\text{supp}(Y_l)$ is roughly of size $N/m$. We thus define the above

quantities as follows:

$$\mathbb{E}_m[Y_l] := \frac{m^2}{N^2} \sum_{\boldsymbol{n} \in \mathbb{Z}^2} \mathbb{E}[Y_l[\boldsymbol{n}]]; \quad (29)$$

$$\mathbb{V}_m[Y_l] := \frac{m^2}{N^2} \sum_{\boldsymbol{n} \in \mathbb{Z}^2} \mathbb{V}[Y_l[\boldsymbol{n}]]. \quad (30)$$

In practice, estimators are computed over a minibatch of images, hence the layer's denomination. Besides, $\varepsilon > 0$ is a small constant added to the denominator for numerical stability. For the sake of concision, we now assume that $\boldsymbol{a} = \mathbf{1}$. Extensions to other multiplicative factors is straightforward.

Let $l \in \mathcal{G}$ denote a Gabor channel. Then, recall that $Y_l$ satisfies (19) (output of the WBlock), with

$$\widetilde{V}_l := \text{Re } \widetilde{W}_l, \quad (31)$$

where $\widetilde{W}_l$ denotes one of the Gabor-like filters spawned by DT-$\mathbb{C}$WPT. The following proposition states that, if the kernel's bandwidth is small enough, then the output of the convolution layer sums to zero.

**Proposition 1.** *We assume that the Fourier transform of* $\widetilde{W}_l$ *is supported in a region of size* $\kappa \times \kappa$ *which does not contain the origin (Gabor-like filter). If, moreover,* $\kappa \leq \frac{2\pi}{m}$, *then*

$$\sum_{\boldsymbol{n} \in \mathbb{Z}^2} Y_l[\boldsymbol{n}] = 0. \quad (32)$$

*Proof.* This proposition takes advantage of Shannon's sampling theorem. A similar reasoning can be found in the proof of Theorem 1 in [23]. $\square$

In practice, the power spectrum of DT-$\mathbb{C}$WPT filters cannot be exactly zero on regions with nonzero measure, since they are finitely supported. However, we can reasonably assume that it is concentrated within a region of size $\pi/2^{J-1} = \pi/m$, as explained in [23]. Therefore, since we have discarded low-pass filters, the conditions of Prop. 1 are approximately met for $\widetilde{W}_l$.

We now assume that (32) is satisfied. Moreover, we assume that $\mathbb{E}[Y_l[\boldsymbol{n}]]$ is constant for any $\boldsymbol{n} \in \text{supp}(Y_l)$. Aside from boundary effects, this is true if $\mathbb{E}[X^{\text{lum}}[\boldsymbol{n}]]$ is constant for any $\boldsymbol{n} \in \text{supp}(X^{\text{lum}})$.[7] We then get, for any $\boldsymbol{n} \in \mathbb{Z}^2, \mathbb{E}[Y_l[\boldsymbol{n}]] = 0$. Therefore, interchanging max pooling and ReLU yields the normalized version of (10):

$$A_l^{\max} = \text{ReLU}\left( \frac{Y_l^{\max}}{\sqrt{\mathbb{E}_m[Y_l^2] + \varepsilon}} + b_l \right). \quad (33)$$

---

[7]This property is a rough approximation for images of natural scenes or man-made objects. In practice, the main subject is generally located at the center, the sky at the top, *etc.* These are sources of variability for color and luminance distributions across images, as discussed by Torralba and Oliva [32].
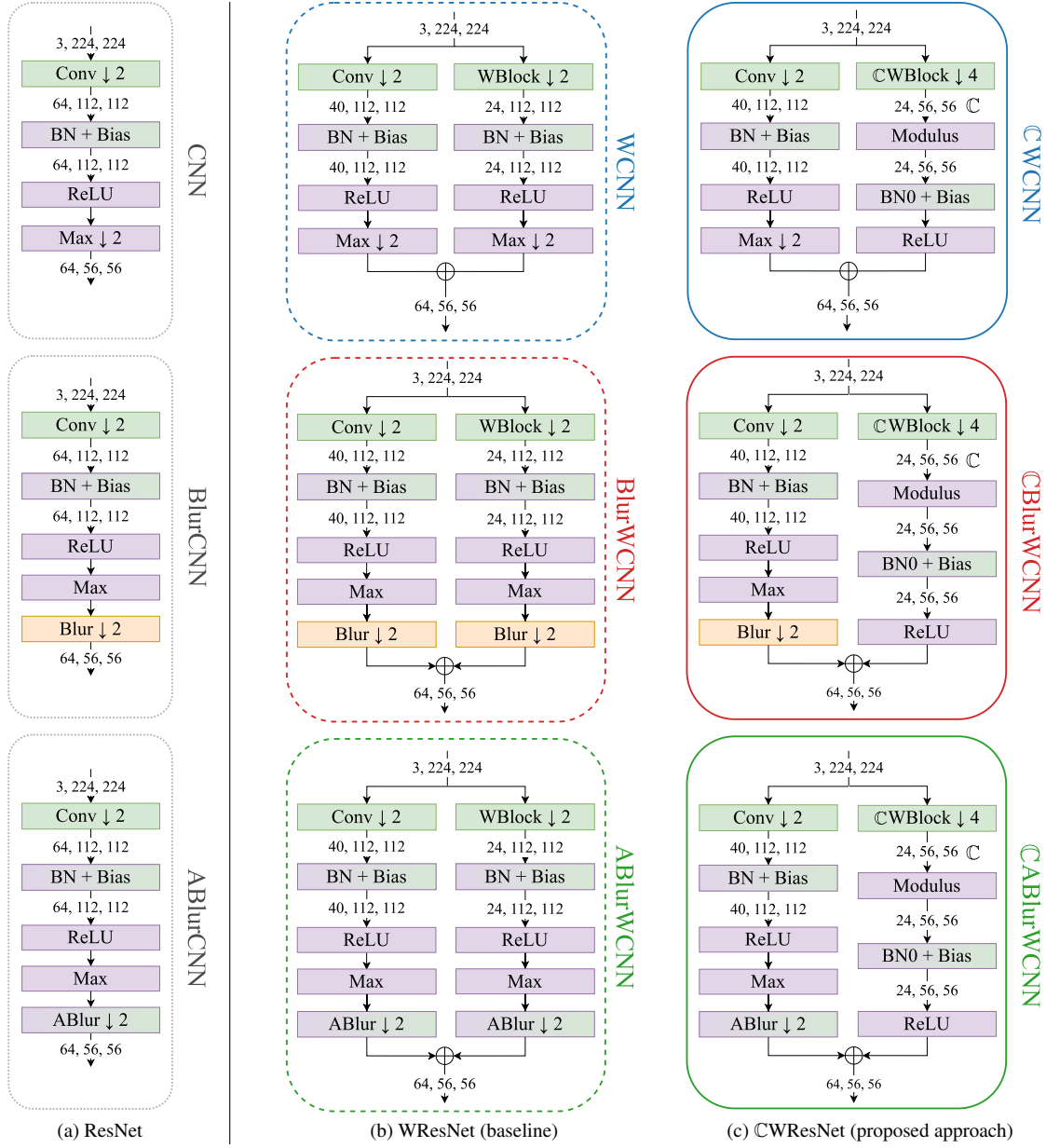
Figure 8. First layers of ResNet and its variants, corresponding to a convolution layer followed by ReLU and max pooling. The models are framed according to the same colors and line styles as in Fig. 10. The bias module from Fig. 3 has been replaced by an affine batch normalization layer ("BN + Bias", or "BN0 + Bias" when placed after Modulus—see Appendix B). Top: ResNet without blur pooling. Middle: Zhang's "blurpooled" models [42]. Bottom: Zou *et al.*'s approach, using adaptive blur pooling [45].

As in Sec. 2.2, we replace $Y_l^{\max}$ by $Y_l^{\mathrm{mod}}$ for any Gabor channel $l \in \mathcal{G}$, which yields the normalized version of (12):

$$A_l^{\mathrm{mod}} := \mathrm{ReLU}\left( \frac{Y_l^{\mathrm{mod}}}{\sqrt{\mathbb{E}_m[Y_l^2] + \varepsilon}} + b_l \right). \qquad (34)$$

Implementing (34) as a deep learning architecture is cumbersome because $Y_l$ needs to be explicitly computed and kept in memory, in addition to $Y_l^{\mathrm{mod}}$. Instead, we want

to express the second-order moment $\mathbb{E}_m[Y_l^2]$ (in the denominator) as a function of $Y_l^{\mathrm{mod}}$. To this end, we state the following proposition.

**Proposition 2.** *If we restrict the conditions of Prop. 1 to $\kappa \leq \pi/m$, we have*

$$\|Y_l\|_2^2 = 2 \left\| Y_l^{\mathrm{mod}} \right\|_2^2. \qquad (35)$$

*Proof.* This result, once again, takes advantage of Shannon's sampling theorem. The proof of our Proposition 3 in [23] is based on similar arguments. □

As for Prop. 1, the conditions of Prop. 2 are approximately met. We therefore assume that (35) is satisfied, and (34) becomes

$$A_l^{\mathrm{mod}} := \mathrm{ReLU}\left(\frac{Y_l^{\mathrm{mod}}}{\sqrt{\frac{1}{2}\mathbb{E}_{2m}[Y_l^{\mathrm{mod}\,2}] + \varepsilon}} + b_l\right). \quad (36)$$

In the case of ResNet, the bias layer (Bias) is therefore preceded by a batch normalization layer without mean centering satisfying (36), which we call BN0. The second-order moment of $Y_l^{\mathrm{mod}}$ is computed on feature maps which are twice smaller than $Y_l$ in both directions—hence the index "$2m$" in (36), which is the subsampling factor for the $\mathbb{C}$Mod operator.

Schematic representations of $\mathbb{R}$Max- and $\mathbb{C}$Mod-based $\mathbb{C}$WResNet are provided in Fig. 8.
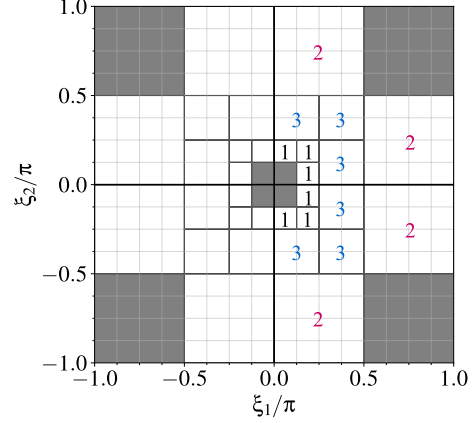
## C. Experimental Settings

In this section, we provide further information that complements the experimental details presented in Sec. 3.1 and Tab. 1.

As explained in Sec. 2.3, the decomposition depth $J$ is chosen such that $m = 2^{J-1}$ (subsampling factor). Since $m = 4$ in AlexNet and 2 in ResNet, we get $J = 3$ and 2, respectively. Therefore, the number of dual-tree filters $K_{\mathrm{dt}} := 2 \times 4^J$ is equal to 128 and 32, respectively.
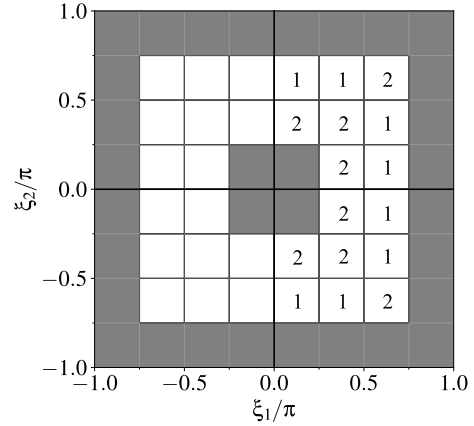
We then manually selected $K'_{\mathrm{dt}} < K_{\mathrm{dt}}$ filters. In particular, we removed the two low-pass filters, which are outside the scope of our theoretical study. Besides, for computational reasons, in WAlexNet we removed 32 "extremely" high-frequency filters which are clearly absent from the standard model (see Fig. 9a). Finally, in WResNet we removed the 14 filters whose bandwidths outreach the boundaries of the Fourier domain $[-\pi,\,\pi]^2$ (see Fig. 9b). These filters indeed have a poorly-defined orientation, since a small fraction of their energy is located at the far end of the Fourier domain [2, see Fig. 1, "Proposed DT-$\mathbb{C}$WPT"]. Therefore, they somewhat exhibit a checkerboard pattern.[8]

As explained in Appendix A, once the DT-$\mathbb{C}$WPT feature maps have been manually selected, the output $\mathbf{D}' \in \mathcal{S}^{K'_{\mathrm{dt}}}$ is sliced into $Q$ groups of channels $\mathbf{D}^{(q)} \in \mathcal{S}^{K_q}$. For each group $q$, a depthwise linear mapping from $\mathbf{D}^{(q)}$ to a bunch of output channels $\mathbf{Y}^{(q)} \in \mathcal{S}^{L_q}$ is performed. Finally, the wavelet block's output feature maps $\mathbf{Y}^{\mathrm{gabor}} \in \mathcal{S}^{L_{\mathrm{gabor}}}$ are obtained by concatenating the outputs $\mathbf{Y}^{(q)}$ depthwise, for

---

[8]Note that the same procedure could have been applied to WAlexNet, but it was deemed unnecessary because the boundary filters were spontaneously discarded during training.



(a) WAlexNet ($J = 3$)



(b) WResNet ($J = 2$)

Figure 9. Mapping scheme from DT-$\mathbb{C}$WPT feature maps $\mathbf{D} \in \mathcal{S}^{K_{\mathrm{dt}}}$ to the wavelet block's output $\mathbf{Y}^{\mathrm{gabor}} \in \mathcal{S}^{L_{\mathrm{gabor}}}$. Each wavelet feature map is symbolized by a small square in the Fourier domain, where its energy is mainly located. The gray areas show the feature maps which have been manually removed. Elsewhere, each group of feature maps $\mathbf{D}^{(q)} \in \mathcal{S}^{K_q}$ is symbolized by a dark frame—in (b), $K_q$ is always equal to 1. For each group $q \in \{1..Q\}$, a number indicates how many output channels $L_q$ are assigned to it. The colored numbers in (a) refer to groups on which we have applied $l^\infty/l^1$-regularization. Note that, when inputs are real-valued, only the half-plane of positive $x$-values is considered.

any $q \in \{1..Q\}$. Figure 9 shows how the above grouping is made, and how many output channels $L_q$ each group $q$ is assigned to.

During training, the above process aims at selecting one single DT-$\mathbb{C}$WPT feature map among each group. This is achieved through mixed-norm $l^\infty/l^1$ regularization, as introduced in (27). The regularization hyperparameters $\lambda_q$ have been chosen empirically. If they are too small, then regularization will not be effective. On the contrary, if they

| Model | Filt. frequency | Reg. param. |
|---|---|---|
| WAlexNet | $[\pi/8,\ \pi/4[$ | − |
|  | $[\pi/4,\ \pi/2[$ | $4.1 \cdot 10^{-3}$ |
|  | $[\pi/2,\ \pi[$ | $3.2 \cdot 10^{-4}$ |
| WResNet | any | − |

Table 5. Regularization hyperparameters $\lambda_q$ for each group $q$ of DT-ℂWPT feature maps. The groups with only one feature map do not need any regularization since this feature map is automatically selected. The second and third rows of WAlexNet correspond to the blue and magenta groups in Fig. 9a, respectively.

| Model | One-crop | | Ten-crops | | Shifts |
|---|---|---|---|---|---|
|  | top-1 | top-5 | top-1 | top-5 | mFR |
| WCNN♦ | 27.4 | 9.2 | 24.7 | 7.6 | 77.2 |
| BlurWCNN† | 26.7 | 8.6 | 24.1 | 7.3 | 65.2 |
| ℂBlurWCNN†* | **26.5** | **8.4** | **23.7** | **7.0** | 62.5 |
| → BlurWCNN†◇ | 26.8 | 8.6 | 24.3 | 7.1 | 69.7 |
| ABlurWCNN‡ | **26.0** | **8.2** | **23.6** | **6.9** | **62.1** |
| ℂABlurWCNN‡* | 26.1 | **8.2** | 23.7 | 7.0 | 63.1 |
| → ABlurWCNN†◇ | 26.4 | 8.4 | 23.9 | 7.0 | 70.3 |

Table 6. Ablation study on WResNet. Evaluation metrics on ImageNet (%): the lower the better. Models: †static and ‡adaptive blur pooling; *ℂMod-based antialiasing on the Gabor channels (our approach); ◇ablated model: no antialiasing on the Gabor channels (neither blur pooling nor ℂMod); ♦no antialiasing anywhere.

are too large, then the regularization term will become predominant, forcing the trainable parameter vector $\boldsymbol{\alpha}_l^{(q)}$ to randomly collapse to 0 except for one element. The chosen values of $\lambda_q$ are displayed in Tab. 5.

Finally, the split $L_{\text{free}}$-$L_{\text{gabor}}$ between the freely-trained and Gabor channels, provided in the last row of Tab. 1 (main document), have been empirically determined from the standard models. More specifically, considering standard AlexNet and ResNet-34 trained on ImageNet (see Fig. 12), we determined the characteristics of each convolution kernel: frequency, orientation, and coherence index (which indicates whether an orientation is clearly defined). This was done by computing the *tensor structure* [16]. Then, by applying proper thresholds, we isolated the Gabor-like kernels from the others, yielding the approximate values of $L_{\text{free}}$ and $L_{\text{gabor}}$. Furthermore, this procedure allowed us to draw a rough estimate of the distribution of the Gabor-like filters in the Fourier domain, which was helpful to design the mapping scheme shown in Fig. 9.

# D. Additional Figures and Tables

## D.1. Validation Error Along Training

As for AlexNet in Fig. 4, Top-1 validation curves of ResNet-based models along training with ImageNet and CIFAR-10 are plotted in Fig. 10. In many cases, accuracy increases when applying our ℂMod-based antialiasing method (dashed versus solid lines). However, in some cases, the gain on the validation set (a subset of images put aside from the training set) is marginal, for instance BlurWResNet-34 versus ℂBlurWResNet-34 trained on ImageNet (Fig. 10a, red stars), or non-existent, for instance WResNet-18 versus ℂWResNet-18 trained on CIFAR-10 (Fig. 10b, blue diamonds). Yet, notable gains are observed on the evaluation set (a separate dataset provided by ImageNet or CIFAR), as reported in Tabs. 2 and 3. This indicates a better generalization capability of our approach, compared to the static blur pooling-based method, or no antialiasing. In particular, for ImageNet classification (Tab. 2, bottom), simply replacing the blurpooled Gabor channels in

the first layer with our ℂMod-based approach (BlurWCNN versus ℂBlurWCNN) produces improvements on a similar order of magnitude as adaptive blur pooling, which is applied throughout the whole network (ABlurWCNN).
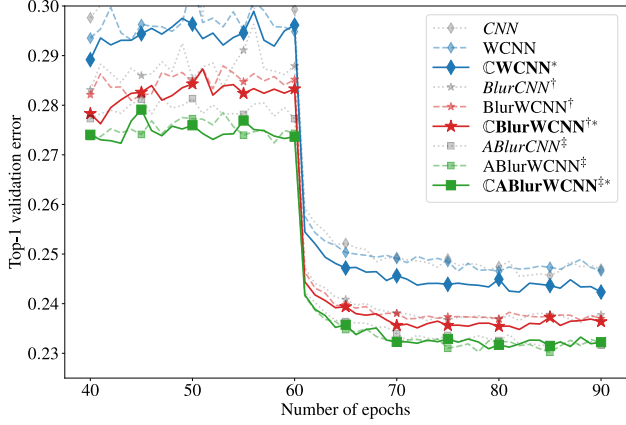
## D.2. Accuracy vs Consistency

Figure 11 shows the relationship between consistency and prediction accuracy of AlexNet and ResNet-based models on ImageNet, for different filter sizes ranging from 1 (no blur pooling) to 7 (heavy loss of high-frequency information). The data for AlexNet on the validation set are displayed in the main document, Fig. 6. As recommended by Zhang [42], the optimal trade-off is generally achieved when the blurring filter size is equal to 3. Moreover, in either case, at equivalent level of consistency, replacing blur pooling by our ℂMod-based antialiasing approach in the Gabor channels increases accuracy.
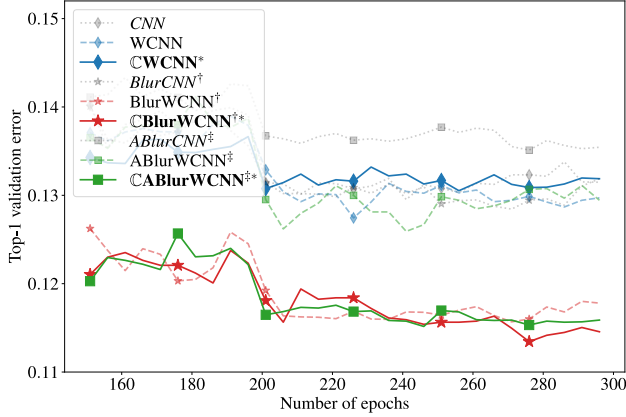
## D.3. Ablation Study

In Sec. 3.3, we showed that our models outperform the baselines, except when adaptive blur pooling is used as an antialiasing method. In that case, although applying our ℂMod-based approach on the Gabor channels instead of adaptive blur pooling leads to slightly degraded accuracy, it comes with far less computational resources and memory footprint, as shown in Tab. 4. However, one may wonder: what happens if we do not employ any antialiasing on the Gabor channels (and keep on using blur pooling in the rest of the network)? Do we reach similar accuracy as well? To answer this question, we trained such ablated models, in both static and adaptive blur-pooling situations.

We found that removing antialiasing from the Gabor channels (either blur pooling or ℂMod) generally leads to decreased accuracy and consistency, as shown in Tab. 6. However, when using adaptive blur pooling, the primary improvements in performance come from other parts of the network—as evidenced when comparing the first and last rows. This limits the influence of antialiasing on the Gabor channels. Yet, although adaptive blur pooling yields the
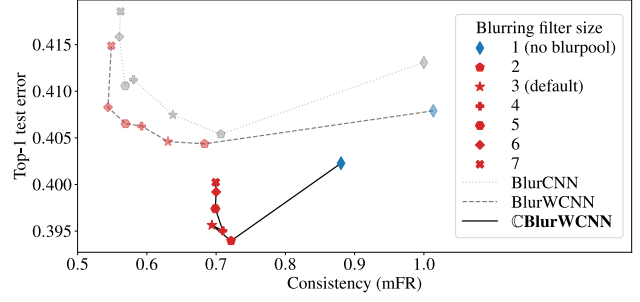
(a) ResNet-34, trained on ImageNet-1K
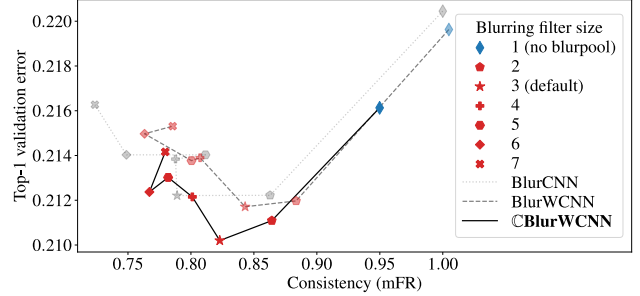


(b) ResNet-18, trained on CIFAR-10



(c) ResNet-34, trained on CIFAR-10

Figure 10. Evolution of the top-1 validation error (one-crop) along training, for ResNet-based models, as described in Fig. 8. The freely-trained models, upon which the mathematical twins are built, appear in faint gray. Legend: $^\dagger$static and $^\ddagger$adaptive blur pooling; $^*\mathbb{C}$Mod-based antialiasing (our approach).
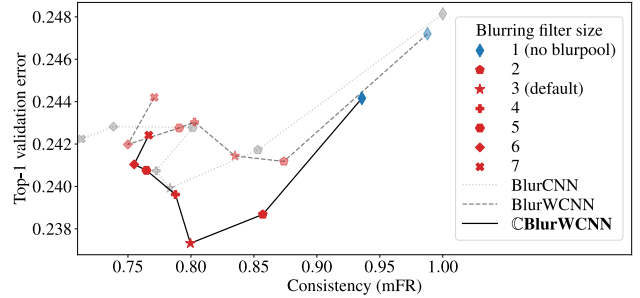
most important gains in performance, it comes at a substantial increase in computational resources due to the need for



(a) AlexNet, test set (50K images)



(b) ResNet-34, validation set (100K images)



(c) ResNet-34, test set (50K images)

Figure 11. Classification accuracy (ten-crops) vs consistency, measuring the stability of predictions to small input shifts (the lower the better for both axes). The metrics have been computed on ImageNet-1K, on both validation set (100K images set aside from the training set) and test set (50K images provided as a separate dataset). For each model (BlurCNN, BlurWCNN and $\mathbb{C}$BlurWCNN), we increased the blurring filter size from 1 (*i.e.*, no blur pooling) to 7. The blue diamonds (no blur pooling) and red stars (blur pooling with filters of size 3) correspond to the models from Fig. 4 for AlexNet and Fig. 10a for ResNet (models trained after 90 epochs).

multiple such antialiasing layers throughout the network.

# E. Kernel Visualization (Standard Models)

The convolution kernels $\mathbf{V} := (\mathbf{V}_{lk})_{l \in \{1..64\}, k \in \{1..3\}} \in \mathcal{S}^{64 \times 3}$ satisfying (11) are shown in Fig. 12 for freely-trained AlexNet and ResNet-34 trained on ImageNet. The kernels are shown as RGB color images. When comparing

$$
\begin{aligned}
t_{\mathrm{s}} &= 1.0 && \text{(addition)} \\
t_{\mathrm{p}} &= 1.0 && \text{(multiplication)} \\
t_{\mathrm{e}} &= 0.75 && \text{(exponential)} \\
t_{\mathrm{mod}} &= 3.5 && \text{(modulus)} \\
t_{\mathrm{relu}} &= 0.75 && \text{(ReLU)} \\
t_{\max} &= 12.0 && \text{(max pooling)}
\end{aligned}
$$

Table 7. Computation time per operation, determined experimentally using PyTorch (CPU computations). They have been normalized with respect to the computation time of an addition.

them with WCNN kernels in Fig. 2, we can notice that, up to a few exceptions, the freely-trained channels (4 and 5 first rows for AlexNet and ResNet, respectively) have been specialized to lower-frequency filters (mono- or bi-color blobs).

## F. Computational Cost

This section provides technical details about our estimation of the computational cost (FLOPs), such as reported in Tab. 4, for *one input image* and *one Gabor channel*. This metric was estimated in the case of standard 2D convolutions.

### F.1. Average Computation Time per Operation

The estimated computation time for each type of operation has been reported in Tab. 7.

### F.2. Computational Cost per Layer

In the following paragraphs, $L \in \mathbb{N} \setminus \{0\}$ denotes the number of output channels (depth) and $N' \in \mathbb{N} \setminus \{0\}$ denotes the size of output feature maps (height and width). However, note that $N'$ is not necessary the same for all layers. For instance, in standard ResNet, the output of the first convolution layer is of size $N' = 112$, whereas the output of the subsequent max pooling layer is of size $N' = 56$. For each type of layer, we calculate the number of FLOPs required to produce a single output channel $l \in \{1 \mathinner{..} L\}$. Moreover, we assume, without loss of generality, that the model processes one input image at a time.

**Convolution Layers.** Inputs of size $(K \times N \times N)$ (input channels, height and width); outputs of size $(L \times N' \times N')$. For each output unit, a convolution layer with kernels of size $\left(N_{\mathrm{filt}} \times N_{\mathrm{filt}}\right)$ requires $K N_{\mathrm{filt}}^2$ multiplications and $K N_{\mathrm{filt}}^2 - 1$ additions. Therefore, the computational cost per output channel is equal to

$$
T_{\mathrm{conv}} = {N'}^2 \left( \left(K N_{\mathrm{filt}}^2 - 1\right) \cdot t_{\mathrm{s}} + K N_{\mathrm{filt}}^2 \cdot t_{\mathrm{p}} \right). \quad (37)
$$

**Complex Convolution Layers.** Inputs of size $(K \times N \times N)$; complex-valued outputs of size $(L \times N' \times N')$. For each output unit, a complex-valued convolution layer requires $2 \times K N_{\mathrm{filt}}^2$ multiplications and $2 \times \left(K N_{\mathrm{filt}}^2 - 1\right)$ additions. Computational cost per output channel:

$$
T_{\mathbb{C}\,\mathrm{conv}} = 2{N'}^2 \left( \left(K N_{\mathrm{filt}}^2 - 1\right) \cdot t_{\mathrm{s}} + K N_{\mathrm{filt}}^2 \cdot t_{\mathrm{p}} \right). \quad (38)
$$

Note that, in our implementations, the complex-valued convolution layers are less expensive than the real-valued ones, because the output size $N'$ is twice smaller, due to the larger subsampling factor.

**Bias and ReLU.** Inputs and outputs of size $(L \times N' \times N')$. One evaluation for each output unit:

$$
T_{\mathrm{bias}} = {N'}^2\, t_{\mathrm{s}} \qquad \text{and} \qquad T_{\mathrm{relu}} = {N'}^2\, t_{\mathrm{relu}}. \quad (39)
$$

**Max Pooling.** Outputs of size $(L \times N' \times N')$, with $N'$ depending on whether subsampling is performed at this stage (no subsampling when followed by a blur pooling layer). One evaluation for each output unit:

$$
T_{\max} = {N'}^2\, t_{\max}. \quad (40)
$$

**Modulus Pooling.** Complex-valued inputs and real-valued outputs of size $(L \times N' \times N')$. One evaluation for each output unit:

$$
T_{\mathrm{mod}} = {N'}^2\, t_{\mathrm{mod}}. \quad (41)
$$

**Batch Normalization.** Inputs and outputs of size $(L \times N' \times N')$. A batch normalization (BN) layer, described in (28), can be split into several stages.

1. Mean: ${N'}^2$ additions.

2. Standard deviation: ${N'}^2$ multiplications, ${N'}^2$ additions (second moment), ${N'}^2$ additions (subtract squared mean).

3. Final value: ${N'}^2$ additions (subtract mean), $2{N'}^2$ multiplications (divide by standard deviation and multiplicative coefficient).

Overall, the computational cost per image and output channel of a BN layer is equal to

$$
T_{\mathrm{bn}} = {N'}^2 \left(4\,t_{\mathrm{s}} + 3\,t_{\mathrm{p}}\right). \quad (42)
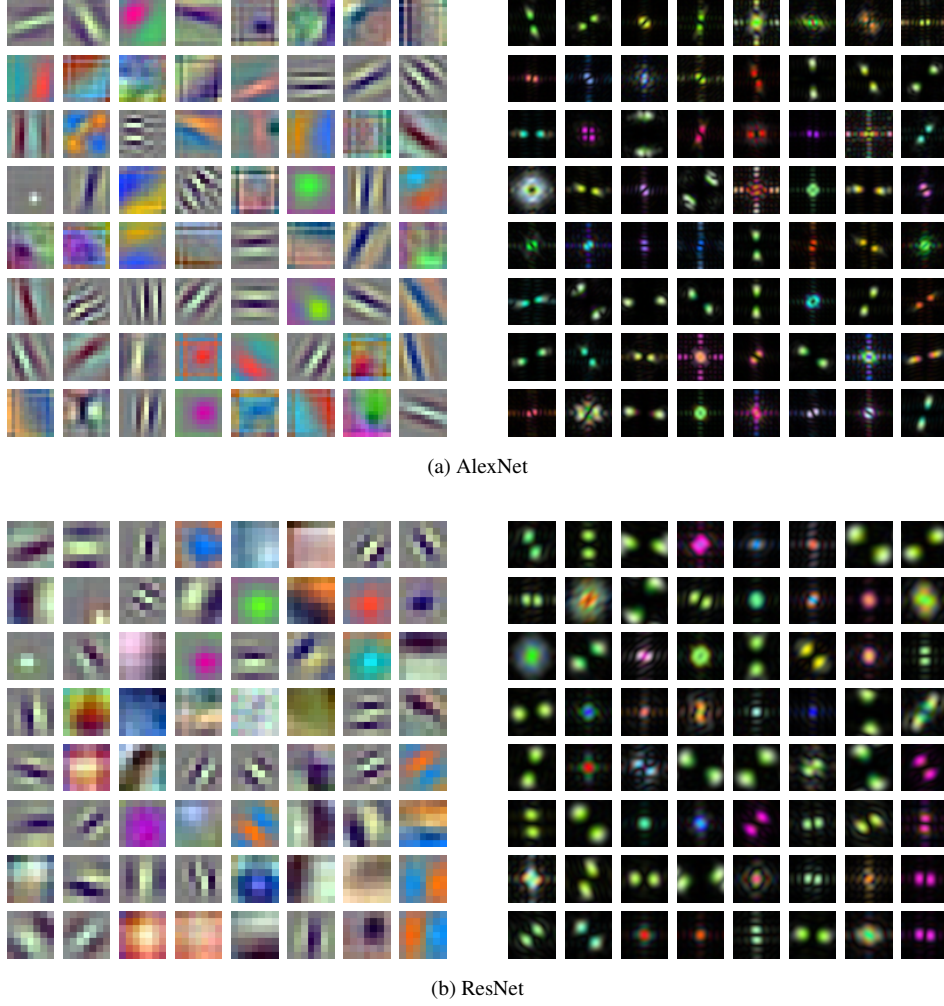$$

(a) AlexNet



(b) ResNet

Figure 12. Convolution kernels in the first layer of freely-trained AlexNet (a) and ResNet-34 (b), after training with ImageNet-1K. For each output channel $l \in \{1 \dots 64\}$, the corresponding convolution kernel $(V_{lk})_{k \in \{1..3\}}$ is displayed as an RGB image in the spatial domain (left), and its associated magnitude spectrum in the Fourier domain (right).

**Static Blur Pooling.** Inputs of size $(L \times 2N' \times 2N')$; outputs of size $(L \times N' \times N')$. For each output unit, a static blur pooling layer [42] with filters of size $(N_\mathrm{b} \times N_\mathrm{b})$ requires $N_\mathrm{b}^2$ multiplications and $N_\mathrm{b}^2 - 1$ additions. The computational cost per output channel is therfore equal to

$$T_{\mathrm{blur}} = N'^2 \left( (N_\mathrm{b}^2 - 1) \cdot t_\mathrm{s} + N_\mathrm{b}^2 \cdot t_\mathrm{p} \right). \qquad (43)$$

**Adaptive Blur Pooling.** Inputs of size $(L \times 2N' \times 2N')$; outputs of size $(L \times N' \times N')$. An adaptive blur pooling layer [45] with filters of size $(N_\mathrm{b} \times N_\mathrm{b})$ splits the $L$ output channels into $Q := L/L_\mathrm{g}$ groups of $L_\mathrm{g}$ channels that share the same blurring filters. The adaptive blur pooling layer can be decomposed into the following stages.

1. Generation of blurring filters using a convolution layer with trainable kernels of size $(N_\mathrm{b} \times N_\mathrm{b})$: inputs of size

$(L \times 2N' \times 2N')$, outputs of size $(QN_\mathrm{b}^2 \times N' \times N')$. For each output unit, this stage requires $LN_\mathrm{b}^2$ multiplications and $LN_\mathrm{b}^2 - 1$ additions. The computational cost divided by the number $L$ of channels is therefore equal to

$$T_{\mathrm{conv\,ablur}} = N'^2 \, \frac{N_\mathrm{b}^2}{L_\mathrm{g}} \times \\ \left( (LN_\mathrm{b}^2 - 1) \cdot t_\mathrm{s} + LN_\mathrm{b}^2 \cdot t_\mathrm{p} \right). \quad (44)$$

Note that, despite being expressed on a per-channel basis, the above computational cost depends on the number $L$ of output channels. This is due to the asymptotic complexity of this stage in $O(L^2)$.

2. Batch normalization, inputs and outputs of size

$(Q N_{\mathrm{b}}^2 \times N' \times N')$:

$$T_{\mathrm{bn\,ablur}} = N'^2 \frac{N_{\mathrm{b}}^2}{L_{\mathrm{g}}} \left( 4\, t_{\mathrm{s}} + 3\, t_{\mathrm{p}} \right). \qquad (45)$$

3. Softmax along the depthwise dimension:

$$T_{\mathrm{sftmx\,ablur}} = N'^2 \frac{N_{\mathrm{b}}^2}{L_{\mathrm{g}}} (t_{\mathrm{e}} + t_{\mathrm{s}} + t_{\mathrm{p}}). \qquad (46)$$

4. Blur pooling of input feature maps, using the filter generated at stages (1)–(3): inputs of size $(L \times 2N' \times 2N')$, outputs of size $(L \times N' \times N')$. The computational cost per output channel is identical to the static blur pooling layer, even though the weights may vary across channels and spatial locations:

$$T_{\mathrm{blur}} = N'^2 \left( (N_{\mathrm{b}}^2 - 1) \cdot t_{\mathrm{s}} + N_{\mathrm{b}}^2 \cdot t_{\mathrm{p}} \right). \qquad (47)$$

Overall, the computational cost of an adaptive blur pooling layer per input image and output channel is equal to

$$\begin{aligned} T_{\mathrm{ablur}} = N'^2 \frac{N_{\mathrm{b}}^2}{L_{\mathrm{g}}} \big[ &\left( (L+1)N_{\mathrm{b}}^2 + 3 \right) \cdot t_{\mathrm{s}} \\ &+ \left( (L+1)N_{\mathrm{b}}^2 + 4 \right) \cdot t_{\mathrm{p}} + t_{\mathrm{e}} \big]. \end{aligned} \qquad (48)$$

We notice that an adaptive blur pooling layer has an asymptotic complexity in $O(N_{\mathrm{b}}^4)$, versus $O(N_{\mathrm{b}}^2)$ for static blur pooling.

### F.3. Application to AlexNet- and ResNet-based Models

Since they are normalized by the computational cost of standard models, the FLOPs reported in Tab. 4 only depend on the size of the convolution kernels and blur pooling filters, respectively denoted by $N_{\mathrm{filt}}$ and $N_{\mathrm{b}} \in \mathbb{N} \setminus \{0\}$. In addition, the computational cost of the adaptive blur pooling layer depend on the number of output channels $L$ as well as the number of output channels per group $L_{\mathrm{g}}$.

In practice, $N_{\mathrm{filt}}$ is respectively equal to 11 and 7 for AlexNet- and ResNet-based models. Moreover, $N_{\mathrm{b}} = 3$, $L = 64$ and $L_{\mathrm{g}} = 8$. Actually, the computational cost is largely determined by the convolution layers, including step (1) of adaptive blur pooling.

### G. Memory Footprint

This section provides technical details about our estimation of the memory footprint for *one input image* and *one output channel*, such as reported in Tab. 4. This metric is generally difficult to estimate, and is very implementation-dependent. Hereafter, we consider the size of the output tensors, as well as intermediate tensors saved by `torch.autograd` for the backward pass. However, we didn't take into account the tensors containing the trainable parameters. To get the size of intermediate tensors, we used the Python package PyTorchViz.[9] These tensors are saved according to the following rules.

- Convolution (Conv), batch normalization (BN), Bias, max pooling (MaxPool or Max), blur pooling (Blur-Pool), and Modulus: the input tensors are saved, not the output. When Bias follows Conv or BN, no intermediate tensor is saved.

- ReLU, Softmax: the output tensors are saved, not the input.

- If an intermediate tensor is saved at both the output of a layer and the input of the next layer, its memory is not duplicated. An exception is Modulus, which stores the input feature maps as complex numbers.

- MaxPool or Max: a tensor of indices is kept in memory, indicating the position of the maximum values. The tensors are stored as 64-bit integers, so they weight twice as much as conventional float-32 tensors.

- BN: four 1D tensors of length $L$ are kept in memory: computed mean and variance, and running mean and variance. For BN0 (see Appendix B), where variance is not computed, only two tensors are kept in memory.

In the following paragraphs, we denote by $L$ the number of output channels, $N$ the size of input images (height and width), $m$ the subsampling factor of the baseline models (4 for AlexNet, 2 for ResNet), $N_{\mathrm{b}}$ the blurring filter size (set to 3 in practice). For each model, a table contains the size of all saved intermediate or output tensors. For example, the values associated to "Layer1 → Layer2" correspond to the depth (number of channel), height and width of the intermediate tensor between Layer1 and Layer2.

### G.1. AlexNet-based Models

**No Antialiasing.** Conv → Bias → ReLU → MaxPool.

| | | | |
|---|---|---|---|
| ReLU → MaxPool | $L$ | $\frac{N}{m}$ | $\frac{N}{m}$ |
| MaxPool → *output* | $L$ | $\frac{N}{2m}$ | $\frac{N}{2m}$ |
| MaxPool indices (×2) | $L$ | $\frac{N}{2m}$ | $\frac{N}{2m}$ |

Then, the memory footprint for each output channel is equal to

$$\implies S_{\mathrm{std}} = \frac{7}{4} \frac{N^2}{m^2}.$$

**Static Blur Pooling.** Conv $\to$ Bias $\to$ ReLU $\to$ BlurPool $\to$ Max $\to$ BlurPool.

| | | | |
|---|---|---|---|
| ReLU $\to$ BlurPool | $L$ | $\frac{2N}{m}$ | $\frac{2N}{m}$ |
| BlurPool $\to$ Max | $L$ | $\frac{N}{m}$ | $\frac{N}{m}$ |
| Max $\to$ BlurPool | $L$ | $\frac{N}{m}$ | $\frac{N}{m}$ |
| Max indices ($\times 2$) | $L$ | $\frac{N}{m}$ | $\frac{N}{m}$ |
| BlurPool $\to$ *output* | $L$ | $\frac{N}{2m}$ | $\frac{N}{2m}$ |

$$\implies S_{\text{blur}} = \frac{33}{4}\frac{N^2}{m^2}.$$

**$\mathbb{C}$Mod-based Approach.** $\mathbb{C}$Conv $\to$ Modulus $\to$ Bias $\to$ ReLU.

| | | | |
|---|---|---|---|
| $\mathbb{C}$Conv $\to$ Modulus | $2L$ | $\frac{N}{2m}$ | $\frac{N}{2m}$ |
| Modulus $\to$ Bias | $L$ | $\frac{N}{2m}$ | $\frac{N}{2m}$ |
| ReLU $\to$ *output* | $L$ | $\frac{N}{2m}$ | $\frac{N}{2m}$ |

$$\implies S_{\text{mod}} = \frac{N^2}{m^2}.$$

## G.2. ResNet-based Models

**No Antialiasing.** Conv $\to$ BN $\to$ Bias $\to$ ReLU $\to$ MaxPool.

| | | | |
|---|---|---|---|
| Conv $\to$ BN | $L$ | $\frac{N}{m}$ | $\frac{N}{m}$ |
| BN metrics | $4L$ | – | – |
| ReLU $\to$ MaxPool | $L$ | $\frac{N}{m}$ | $\frac{N}{m}$ |
| MaxPool $\to$ *output* | $L$ | $\frac{N}{2m}$ | $\frac{N}{2m}$ |
| MaxPool indices ($\times 2$) | $L$ | $\frac{N}{2m}$ | $\frac{N}{2m}$ |

$$\implies S_{\text{std}} = \frac{11}{4}\frac{N^2}{m^2} + 4 \approx \frac{11}{4}\frac{N^2}{m^2}.$$

**Static Blur Pooling.** Conv $\to$ BN $\to$ Bias $\to$ ReLU $\to$ Max $\to$ BlurPool.

| | | | |
|---|---|---|---|
| Conv $\to$ BN | $L$ | $\frac{N}{m}$ | $\frac{N}{m}$ |
| BN metrics | $4L$ | – | – |
| ReLU $\to$ Max | $L$ | $\frac{N}{m}$ | $\frac{N}{m}$ |
| Max $\to$ BlurPool | $L$ | $\frac{N}{m}$ | $\frac{N}{m}$ |
| Max indices ($\times 2$) | $L$ | $\frac{N}{m}$ | $\frac{N}{m}$ |
| BlurPool $\to$ *output* | $L$ | $\frac{N}{2m}$ | $\frac{N}{2m}$ |

$$\implies S_{\text{blur}} = \frac{21}{4}\frac{N^2}{m^2} + 4 \approx \frac{21}{4}\frac{N^2}{m^2}.$$

**Adaptive Blur Pooling.** Conv $\to$ BN $\to$ Bias $\to$ ReLU $\to$ Max $\to$ ABlurPool.

| | | | |
|---|---|---|---|
| Conv $\to$ BN | $L$ | $\frac{N}{m}$ | $\frac{N}{m}$ |
| BN metrics | $4L$ | – | – |
| ReLU $\to$ Max | $L$ | $\frac{N}{m}$ | $\frac{N}{m}$ |
| Max $\to$ ABlurPool | $L$ | $\frac{N}{m}$ | $\frac{N}{m}$ |
| Max indices ($\times 2$) | $L$ | $\frac{N}{m}$ | $\frac{N}{m}$ |
| ABlurPool $\to$ *output* | $L$ | $\frac{N}{2m}$ | $\frac{N}{2m}$ |
| **Generate adaptive blurring filter** | | | |
| Conv $\to$ BN $\to$ Bias $\to$ Softmax | | | |
| Conv $\to$ BN | $\frac{LN_{\text{b}}^2}{L_{\text{g}}}$ | $\frac{N}{2m}$ | $\frac{N}{2m}$ |
| BN metrics | $4\frac{LN_{\text{b}}^2}{L_{\text{g}}}$ | – | – |
| Softmax $\to$ *output* | $\frac{LN_{\text{b}}^2}{L_{\text{g}}}$ | $\frac{N}{2m}$ | $\frac{N}{2m}$ |

$$\implies S_{\text{ablur}} = \frac{21}{4}\frac{N^2}{m^2} + 4 + \frac{N_{\text{b}}^2}{L_{\text{g}}}\left(\frac{N^2}{2m^2} + 4\right)$$
$$\approx \frac{21}{4}\frac{N^2}{m^2} + \frac{N_{\text{b}}^2}{L_{\text{g}}}\frac{N^2}{2m^2}.$$

**$\mathbb{C}$Mod-based Approach.** $\mathbb{C}$Conv $\to$ Modulus $\to$ BN0 $\to$ Bias $\to$ ReLU.

| | | | |
|---|---|---|---|
| $\mathbb{C}$Conv $\to$ Modulus | $2L$ | $\frac{N}{2m}$ | $\frac{N}{2m}$ |
| Modulus $\to$ BN0 | $L$ | $\frac{N}{2m}$ | $\frac{N}{2m}$ |
| BN0 metrics | $2L$ | – | – |
| ReLU $\to$ *output* | $L$ | $\frac{N}{2m}$ | $\frac{N}{2m}$ |

$$\implies S_{\text{mod}} = \frac{N^2}{m^2} + 2 \approx \frac{N^2}{m^2}.$$