# Gabor Layers Enhance Network Robustness

Juan C. Pérez [*1], Motasem Alfarra [*2], Guillaume Jeanneret [*1], Adel Bibi[2], Ali Thabet[2], Bernard Ghanem[2], and Pablo Arbeláez[1]

[1] Universidad de los Andes, Colombia
[2] King Abdullah University of Science and Technology (KAUST), Saudi Arabia

**Abstract.** We revisit the benefits of merging classical vision concepts with deep learning models. In particular, we explore the effect on robustness against adversarial attacks of replacing the first layers of various deep architectures with Gabor layers, *i.e.* convolutional layers with filters that are based on learnable Gabor parameters. We observe that architectures enhanced with Gabor layers gain a consistent boost in robustness over regular models and preserve high generalizing test performance, even though these layers come at a negligible increase in the number of parameters. We then exploit the closed form expression of Gabor filters to derive an expression for a Lipschitz constant of such filters, and harness this theoretical result to develop a regularizer we use during training to further enhance network robustness. We conduct extensive experiments with various architectures (LeNet, AlexNet, VGG16 and WideResNet) on several datasets (MNIST, SVHN, CIFAR10 and CIFAR100) and demonstrate large empirical robustness gains. Furthermore, we experimentally show how our regularizer provides consistent robustness improvements.

**Keywords:** Gabor, Robustness, Adversarial Attacks, Regularizer.

## 1 Introduction

Deep neural networks (DNNs) have provided outstanding gains in performance in several fields, from computer vision [21,16] to machine learning [24] and natural language processing [17]. However, despite this success, powerful DNNs are still highly susceptible to small perturbations in their input, known as adversarial attacks [15]. Their accuracy on standard benchmarks can be drastically reduced in the presence of perturbations that are imperceptible to the human eye. Furthermore, the construction of such perturbations is rather undemanding and, in some cases, as simple as performing a single gradient ascent step of a loss function with respect to the image [15].

The brittleness of DNNs in the presence of adversarial attacks has spurred much interest in the machine learning community, as evidenced by the emerging corpus of recent methods that focus on designing adversarial attacks [15,6,33,48]. This phenomenon is far-reaching and widespread, and is of particular importance in real-world scenarios, *e.g.*, autonomous cars [5,9] and devices for the visually

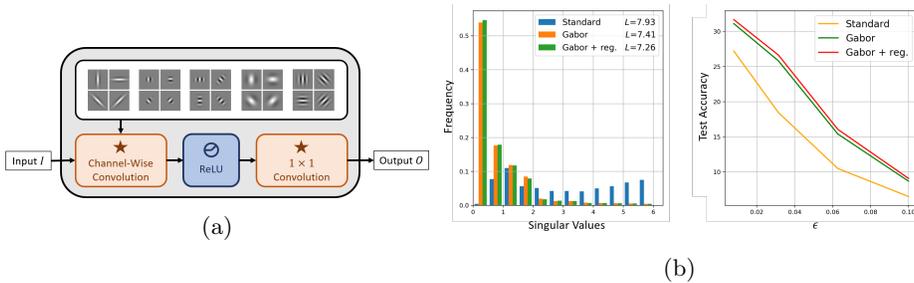---

[*] denotes equal contribution

Fig. 1: **Gabor layers and their effect on network robustness.** (a): Gabor layers convolve each channel of the input with a set of learned Gabor filters. As low-level filters, Gabor filters offer a natural approach to represent local signals. (b): Replacing standard convolutional layers with Gabor layers imposes structure in the distribution of singular values of the filters, reduces the Lipschitz constant of the filters (shown as $L$ in the legend of the left plot), and improves accuracy under adversarial attacks (right figure). The results shown here are for VGG16 on CIFAR100.

impaired [38]. The risks that this degenerate behavior poses underscore the need for models that are not only accurate, but also robust to adversarial attacks.

Despite the complications that adversarial examples raise in modern computer vision, such inconveniences were not a major concern in the pre-DNN era. Many classical computer vision methods drew inspiration from the animal visual system, and so were designed to extract and use features that were meaningful to humans [27,28,35,26]. As such, these methods were structured, generally comprehensible and, hence, better understood when compared to DNNs. In many cases, these methods even exhibited rigorous stability properties under robustness analysis [14]. However, mainly due to large performance gaps on several tasks, classical methods were overshadowed by DNNs. It is precisely in the frontier between classical computer vision and DNNs that a stream of works arose to combine tools and insights from both worlds to improve performance. For instance, the works of [46,52] showed that introducing structured layers inspired by the classical compressed sensing literature can outperform pure learning-based DNNs. Moreover, Bai *et al.* [3] achieved large gains in performance in instance segmentation by introducing intuitions from the classical watershed transform into DNNs.

In this paper, searching for robustness in computer vision, we draw inspiration from biological vision, as the survival of species strongly depends on both the accuracy and robustness of the animal visual system. We note that Marr's and Julesz' work [31,19] argues that the visual cortex initially processes low-level agnostic information, in which the input of the system is segmented according to blobs, edges, bars, curves and boundaries. Furthermore, Hubel and Wiesel [18] demonstrated that individual cells on the primary visual cortex respond to wave textures with different angles in an animal model, providing evidence that supports Marr's theory. Since Gabor filters [13] are based on mathemat-

ical functions that are capable of modeling elements that resemble those that the animal visual cortices respond to, these filters became of customary use in the computer vision community, and have been used for texture characterization [19,26], character recognition [45], edge detection [34], and face recognition [22,8]. While several works examine their integration into DNNs [41,29,1], none investigate the effect of introducing parameterized Gabor filters into DNNs on the robustness of these networks. Our work fills this gap in the literature, as we provide experimental results demonstrating the significant impact that such an architectural change has on improving robustness. Figure 1 shows an overview of our work and results.

**Contributions**: Our main contributions are two-fold: **(1)** We propose a *parameterized* Gabor-structured convolutional layer as a replacement for early convolutional layers in DNNs. We observe that such layers can have a remarkable impact on robustness. Thereafter, we analyze and derive an analytical expression for a Lipschitz constant of the Gabor filters, and propose a new training regularizer to further boost robustness. **(2)** We empirically validate our claims with a large number of experiments on different architectures (LeNet [25], AlexNet [21], VGG16 [44] and Wide-ResNet [50]) and over several datasets (MNIST [23], SVHN [32], CIFAR10 and CIFAR100 [20]). We show that introducing our proposed Gabor layers in DNNs induces a consistent boost in robustness at negligible cost, while preserving high generalizing test performance. In addition, we experimentally show that our novel regularizer based on the Lipschitz constant we derive can further improve adversarial robustness. For instance, we improve adversarial robustness on certain networks by almost 18% with $\ell_\infty$ bounded noise of $8/255$. Lastly, we show empirically that combining this architectural change with adversarial training [30,43] can further improve robustness.

## 2   Related Work

**Integrating Gabor Filters with DNNs.** Several works attempted to combine Gabor filters and DNNs. For instance, the work of [41] showed that replacing the first convolutional layers in DNNs with Gabor filters speeds up the training procedure, while [29] demonstrated that introducing Gabor layers reduces the parameter count without hurting generalization accuracy. Regarding large scale datasets, Alekseev and Bobe [1] showed that the standard classification accuracy of AlexNet [21] on ImageNet [40] can be attained even when the first convolutional filters are replaced with Gabor filters. Moreover, other works have integrated Gabor filters with DNNs for various applications, *e.g.*, pedestrian detection [36], object recognition [49], hyper-spectral image classification [7], and Chinese optical character recognition [53]. Likewise, in this work, we study the effects of introducing Gabor filters into various DNNs by means of a *Gabor layer*, a convolution-based layer we propose in which the convolutional filters are constructed by a parameterized Gabor function with learnable parameters. Furthermore, and based on the well-defined spatial structure of these filters, we study the effect of these layers on robustness, and find encouraging results.

**Robust Neural Networks.** Recent work demonstrated that DNNs are vulnerable to adversarial perturbations. This susceptibility incited a stream of research that aimed to develop not only accurate but also robust DNNs. A straightforward approach to this nuisance is the direct augmentation of data corrupted with adversarial examples in the training set [15]. However, the performance of this approach can be computationally limited, since the amount of augmentation needed for a high dimensional input space is computationally prohibitive. Moreover, Papernot *et al.* [37] showed that distilling DNNs into smaller networks can improve robustness. Another approach to robustness is through the functional perspective lens. For instance, Parseval Networks [11] showed that robustness can be achieved by regularizing the Lipschitz constant of each layer in a DNN to be smaller than 1. In this work, along the lines of Parseval Networks [11], and since Gabor filters can be generated by sampling from a continuous Gabor function, we derive an analytical closed form expression for the Lipschitz constant of the filters of the proposed Gabor layer. This derivation allows us to propose well-motivated regularizers that can encourage Lipschitz constant minimization, and then harness such regularizers to improve the robustness of networks with Gabor layers.

**Adversarial Training.** An orthogonal direction for obtaining robust models is through optimization of a saddle point problem, in which an adversary, whose aim is to maximize the objective, is introduced into the traditional optimization objective. In other words, instead of the typical training scheme, one can minimize the worst adversarial loss over all bounded energy (often measured in $\ell_\infty$ norm) perturbations around every given input in the training data. This approach is one of the most celebrated for training robust networks, and is now popularly known as adversarial training [30]. However, this training comes at an inconvenient computational cost. To this regard, several works [47,43,51] proposed faster and computationally-cheaper versions of adversarial training capable of achieving similar robustness levels. In this work, we use "free" adversarial training [43] in our experiments to further study Gabor layers and adversarial training as orthogonal approaches to achieve robustness. Our results show how adversarial training and Gabor layers interact positively and can, hence, be jointly used for enhancing network robustness.

## 3   Methodology

As demonstrated by Hubel and Wiesel [18], the first layers of visual processing in the animal brain are responsible for detecting low-level visual information. Since Gabor filters have the capacity to capture low-level representations, and inspired by the robust properties of the animal visual system, we hypothesize that Gabor filters possess inherent robustness properties that are transferable to other systems, perhaps even DNNs. In this section, we discuss our proposed Gabor layer and its implementation. Then, we derive a Lipschitz constant to the Gabor filter, and design a regularizer with aims at controlling the robustness properties of the layer by controlling the Lipschitz constant.
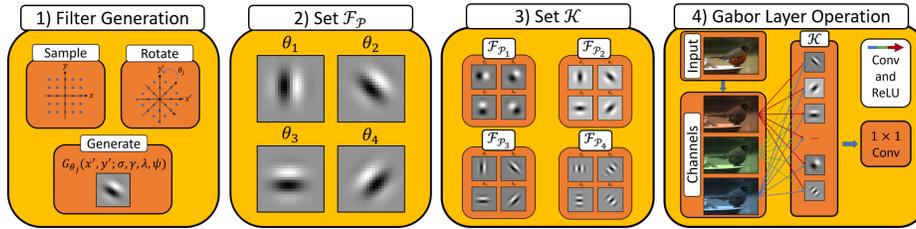
Fig. 2: **Gabor layer operations.** *(1)* We generate filters by rotating a sampled grid over multiple orientations and then evaluating the Gabor function according to set of parameters $\mathcal{P}$, to yield the set of filters $\mathcal{F}_\mathcal{P}$ *(2)*. Then, we construct the total set of filters $\mathcal{K}$ *(3)* by joining multiple sets $\mathcal{F}_{\mathcal{P}_i}$. Finally, the Gabor Layer operation *(4)* separately convolves every filter in $\mathcal{K}$ with every channel from the input, applies ReLU non-linearity, and then applies a $1 \times 1$ convolution to the output features to get the desired number of output channels.

### 3.1 Convolutional Gabor Filter as a Layer

We start by introducing the Gabor functions, defined as follows:

$$
G_\theta(x', y'; \sigma, \gamma, \lambda, \psi) := e^{-\sigma^2 \, (x'^2 + \gamma^2 y'^2)} \, \cos(\lambda x' + \psi)
$$
$$
x' = x\cos\theta - y\sin\theta \quad y' = x\sin\theta + y\cos\theta. \tag{1}
$$

To construct a discrete Gabor filter, we discretize $x$ and $y$ in (1) uniformly on a grid, where the number of grid samples determines the filter size. Given a fixed set of parameters $\{\sigma, \gamma, \lambda, \psi\}$, a grid $\{(x_i, y_i)\}_{i=1}^{k^2}$ of size $k \times k$, a rotation angle $\theta_j$ and a filter scale $\alpha_j$, computing Equation (1) with a scale over the grid yields a single surface $\alpha_j G_{\theta_j}(x', y'; \sigma, \gamma, \lambda, \psi) \in \mathbb{R}^{1 \times k \times k}$, that we interpret as a filter for a convolutional layer. The learnable parameters [39] for such a function are given by the set $\mathcal{P} = \{\alpha_j, \sigma, \gamma, \lambda, \psi; \forall j = 1, \dots, r\}$, where the rotations $\theta_j$ are restricted to be $r$ angles uniformly sampled from the interval $[0, 2\pi]$. Evaluating these functions results in $r$ rotated filters, each with a scale $\alpha_j$ defined by the set $\mathcal{F}_\mathcal{P} = \{\alpha_j G_{\theta_j}\}_{j=1}^r$. In this work, we consider several sets of learnable parameters $\mathcal{P}$, say $p$ of them, thus, the set of all Gabor filters (totaling to $rp$ filters) is given by the set $\mathcal{K} = \{\mathcal{F}_{\mathcal{P}_i}\}_{i=1}^p$. Refer to Figure 2 for a graphical guide on the construction of $\mathcal{K}$.

### 3.2 Implementation of the Gabor Layer

Given an input tensor $I$ with $m$ channels, $I \in \mathbb{R}^{m \times h \times w}$, the Gabor layers follow a depth-wise separable convolution-based [10] approach to convolve the Gabor filters in $\mathcal{K}$ with $I$. The tensor $I$ is first separated into $m$ individual channels. Next, each channel is convolved with each filter in $\mathcal{K}$, and then a ReLU activation is applied to the output. Formally, the Gabor layer with filters in the set $\mathcal{K}$ operating on some input tensor $I$ is presented as $\mathcal{R} = \{\text{ReLU}(I_i \star f_j), \; I_i \in$

$\mathcal{I}$, $f_j \in \mathcal{K}$, $\forall i, j\}$ where $I_i = I(i, :, :) \in \mathbb{R}^{1 \times h \times w}$ and $\star$ denotes the convolution operation. This operation produces $|\mathcal{R}| = mrp$ responses. Finally, the responses are stacked and convolved with a $1 \times 1$ filter with $n$ filters. Thus, the final response is of size $n \times h' \times w'$. Figure 2 shows an overview of this operation.

### 3.3   Regularization

A function $f : \mathbb{R}^n \to \mathbb{R}$ is $L$-Lipschitz if $\|f(x) - f(y)\| \le L\|x - y\|$, $\forall x, y \in \mathbb{R}^n$, where $L$ is the Lipschitz constant. Studying the Lipschitz constant of a DNN is essential for exploring the its robustness properties, since DNNs with a small Lipschitz constant enjoy a better-behaving backpropagated signal of gradients, improved computational stability [42], and an enhanced robustness to adversarial attacks [11].

Therefore, to train networks that are robust, Cisse *et al.* [11] proposed a regularizer that encourages the weights of the networks to be tight frames, which are extensions of orthogonal matrices to non-square matrices. The training procedure is, however, nontrivial to implement. Following the intuition of [11], we study the continuity properties of the Gabor layer and derive an expression for a Lipschitz constant as a function of the parameters of the filters, $\mathcal{P}$. This expression allows for direct network regularization on the parameters of the Gabor layer corresponding to the fastest decrease in the Lipschitz constant of the layer. To this end, we present our main theoretical result, which allows us to apply regularization on the Lipschitz constant of the filters of a Gabor layer.

**Theorem 1.** *Given a Gabor filter $G_\theta(m, n; \sigma, \gamma, \lambda, \psi)$, a Lipschitz constant $L$ of the convolutional layer that has $G_\theta$ as its filter, with circular boundary conditions for the convolution, is given by:*

$$L = \left(1 + |X'|e^{-\sigma^2 m_*^2}\right)\left(1 + |Y'|e^{-\sigma^2\gamma^2 n_*^2}\right),$$

*where $X' = X \setminus \{0\}$, $Y' = Y \setminus \{0\}$, $X = \{x_i\}_{i=1}^{k^2}$ and $Y = \{y_i\}_{i=1}^{k^2}$ are sets of sampled values of the rotated $(x', y')$ grid where $\{0\} \in X, Y$, $m_* = \mathrm{argmin}_{x \in X'}|x|$ and $n_* = \mathrm{argmin}_{y \in Y'}|y|$.*

*Proof.* To compute the Lipschitz constant of a convolutional layer, one must compute the largest singular value of the underlying convolutional matrix of the filter. For separable convolutions, this computation is equivalent to the maximum magnitude of the 2D Discrete Fourier Transform (DFT) of the Gabor filter $G_\theta$ [42,4]. Thus, the Lipschitz constant of the convolutional layer is given by $L = \max_{u,v}|\mathrm{DFT}\left(G_\theta(m, n; \sigma, \gamma, \lambda, \psi)\right)|$, where DFT is the 2D DFT over the coordinates $m$ and $n$ in the spatial domain, $u$ and $v$ are the coordinates in the frequency domain, and $|\cdot|$ is the magnitude operator. Note that $G_\theta$ can be expressed as a product of two functions that are independent of the sampling sets $X$ and $Y$ as follows:

$$G_\theta(m, n; \sigma, \gamma, \lambda, \psi) := \underbrace{e^{-\sigma^2 m^2}\cos(\lambda m + \psi)}_{f(m;\sigma,\lambda,\psi)}\underbrace{e^{-\sigma^2\gamma^2 n^2}}_{g(n;\sigma,\gamma)}.$$

Thus, we have

$$
\begin{aligned}
L &= \max_{u,v} \left| \mathrm{DFT}\left( G_\theta(m,n;\sigma,\gamma,\lambda,\psi) \right) \right| \\
&= \max_{u,v} \left| \sum_{m \in X} e^{-\omega_m u m} f(m;\sigma,\lambda,\psi) \sum_{n \in Y} e^{-\omega_n v n} g(n;\sigma,\gamma) \right| \\
&\leq \max_{u,v} \sum_{m \in X} \left| f(m;\sigma,\lambda,\psi) \right| \sum_{n \in Y} \left| g(n;\sigma,\gamma) \right|.
\end{aligned}
$$

Note that $\omega_m = \frac{j2\pi}{|X|}$, $\omega_n = \frac{j2\pi}{|Y|}$ and $j^2 = -1$. The last inequality follows from Cauchy–Schwarz and the fact that $|e^{-\omega_m u m}| = |e^{-\omega_n v n}| = 1$. Note that since $|g(n;\sigma,\gamma)| = g(n;\sigma,\gamma)$, and $|f(m;\sigma,\lambda,\psi)| \leq e^{-\sigma^2 m^2}$ we have that:

$$
L \leq \sum_{m \in X} e^{-\sigma^2 m^2} \sum_{n \in Y} e^{-\sigma^2 \gamma^2 n^2} \leq \left( 1 + |X'| e^{-\sigma^2 m_*^2} \right) \left( 1 + |Y'| e^{-\sigma^2 \gamma^2 n_*^2} \right).
$$

The last inequality follows by construction, since we have $\{0\} \in X, Y$, *i.e.*, the choice of uniform grid contains the 0 element in both $X$ and $Y$, regardless of the orientation $\theta$, where we define $m_* = \mathrm{argmin}_{x \in X'} |x|$, and $n_* = \mathrm{argmin}_{y \in Y'} |y|$.

<div align="right">□</div>

### 3.4   Lipschitz Constant Regularization

Theorem 1 provides an explicit expression for a Lipschitz constant of the Gabor filter as a function of its parameters. Note that the expression we derived decreases exponentially fast with $\sigma$. In particular, we note that, as $\sigma$ increases, $G_\theta$ converges to a scaled Dirac-like surface. Hence, this Lipschitz constant is minimized when the filter resembles a Dirac-delta. Therefore, to train DNNs with improved robustness, one can minimize the Lipschitz constant we find in Theorem 1. Note that the Lipschitz constant of the network can be upper bounded by the product of the Lipschitz constants of individual layers. Thus, decreasing the Lipschitz constant we provide in Theorem 1 can aid in decreasing the overall Lipschitz constant of a DNN, and thereafter enhance the network's robustness. To this end, we propose the following regularized loss:

$$
\mathcal{L} = \mathcal{L}_{\mathrm{ce}} - \beta \sum_i \sigma_i^2, \tag{2}
$$

where $\mathcal{L}_{\mathrm{ce}}$ is the typical cross-entropy loss and $\beta > 0$ is a trade-off parameter. The loss in Equation (2) can lead to unbounded solutions for $\sigma_i$. To alleviate this behavior, we also propose the following loss:

$$
\mathcal{L} = \mathcal{L}_{\mathrm{ce}} - \beta \sum_i \left( \mu \, \tanh \sigma_i \right)^2, \tag{3}
$$

where $\mu$ is a scaling constant for $\tanh \sigma$. In the following section, we present experiments showing the effect of our Gabor layers on network robustness. Specifically, we show the gains obtained from the architectural modification of introducing Gabor layers, the introduction of our proposed regularizer, and the addition of adversarial training to the overall pipeline.

## 4   Experiments

To demonstrate the benefits and impact on robustness of integrating our Gabor layer to DNNs, we conduct extensive experiments with LeNet [25], AlexNet [21], VGG16 [44], and Wide-ResNet [50] on the MNIST [23], CIFAR10, CIFAR100 [20] and SVHN [32] datasets. In each of the aforementioned networks, we replace up to the first three convolutional layers with Gabor layers, and measure the impact of the Gabor layers in terms of accuracy, robustness, and the distribution of singular values of the layers. Moreover, we perform experiments demonstrating that the robustness of the Gabor-layered networks can be enhanced furthermore by using the regularizer we propose in Equations (2) and (3), and even when jointly employing regularization and adversarial training [43].

### 4.1   Implementation Details

We train all networks with stochastic gradient descent with weight decay of $5 \times 10^{-4}$, momentum of 0.9, and batch size of 128. For MNIST, we train the networks for 90 epochs with a starting learning rate of $10^{-2}$, which is multiplied by a factor of $10^{-1}$ at epochs 30 and 60. For SVHN, we train models for 160 epochs with a starting learning rate of $10^{-2}$ that is multiplied by a factor of $10^{-1}$ at epochs 80 and 120. For CIFAR10 and CIFAR100, we train the networks for 300 epochs with a starting learning rate of $10^{-2}$ that is multiplied by a factor of $10^{-1}$ every 100 epochs.

### 4.2   Robustness Assessment

Following common practice in the literature for empirically evaluating robustness properties [30,43], we assess the robustness of a DNN by measuring its prediction accuracy when the input is probed with adversarial attacks, which is widely referred to in the literature as "adversarial accuracy". We also measure the "flip rate", which is defined as the percentage of instances of the test set for which the predictions of the network changed when under adversarial attacks.

Formally, if $x \in \mathbb{R}^d$ is some input to a classifier $C : \mathbb{R}^d \to \mathbb{R}^k$, $C(x)$ is the prediction of $C$ at input $x$. Then, $x^{\mathrm{adv}} = x + \eta$ is an adversarial example if the prediction of the classifier has changed, *i.e.* $C(x^{\mathrm{adv}}) \neq C(x)$. Both $\eta$ and $x^{\mathrm{adv}}$ must adhere to constraints, namely: *(i)* the $\ell_p$-norm of $\eta$ must be bounded by some $\epsilon$, *i.e.*, $\|\eta\|_p \leq \epsilon$, and *(ii)* $x^{\mathrm{adv}}$ must lie in the space of valid instances $X$, *i.e.*, $x^{\mathrm{adv}} \in [0,1]^d$. A standard approach to constructing $x^{\mathrm{adv}}$ for some input $x$ is by running Projected Gradient Descent (PGD) [30] with $x$ as an initialization for several iterations. For some loss function $\mathcal{L}$, a PGD iteration projects a step of the Fast Gradient Sign Method [15] onto the valid set $\mathcal{S}$ which is defined by the constraints on $\eta$ and $x^{\mathrm{adv}}$. Formally, one iteration of PGD attack is:

$$x^{k+1} = \prod_{\mathcal{S}} \left( x^k + \delta \, \mathrm{sign} \left( \nabla_{x^k} \mathcal{L}(x^k, y) \right) \right),$$

Table 1: **Test set accuracies on different datasets of various baselines, and their Gabor-layered versions.** Gabor-layered architectures can recover the accuracies of their standard counterparts while providing robustness. $\Delta$ is the absolute difference between the baselines and the Gabor-layered architectures.

| Dataset | Architecture | Baseline | Gabor | $\Delta$ |
|---------|--------------|----------|-------|----------|
| MNIST | LeNet | 99.36 | 99.03 | 0.33 |
| SVHN | WideResNet | 96.62 | 96.70 | 0.08 |
| SVHN | VGG16 | 96.52 | 96.18 | 0.34 |
| CIFAR10 | VGG16 | 92.03 | 91.35 | 0.68 |
| CIFAR100 | AlexNet | 46.48 | 45.15 | 1.33 |
| CIFAR100 | WideResNet | 77.68 | 76.86 | 0.82 |
| CIFAR100 | VGG16 | 67.54 | 64.49 | 3.05 |

where $\prod_{\mathcal{S}}$ is the projection operator onto $\mathcal{S}$ and $y$ is the label. In our experiments, we consider attacks where $\eta$ is $\epsilon$-$\ell_\infty$ bounded. For each image, we run PGD for 200 iterations and perform 10 random restarts inside the $\epsilon$-$\ell_\infty$ ball centered in the image. Following prior art [47], we set $\epsilon \in \{0.1, 0.2, 0.3\}$ for MNIST and $\epsilon \in \{2/255, 8/255, 16/255\}$ for all other datasets. Throughout our experiments, we assess robustness by measuring the adversarial accuracies and flip rates when under PGD attacks.

### 4.3   Performance of Gabor-Layered Architectures

The Gabor function in Equation (1) restricts the space of patterns attainable by the Gabor filters. However, this set of patterns is aligned with what is observed in practice in the early layers of many standard architectures [2,21]. This observation follows the intuition that DNNs learn hierarchical representations, with early layers detecting lines and blobs, and deeper layers learning semantic information [12]. By experimenting with Gabor layers on various DNNs, we find that Gabor-layered DNNs recover close-to-standard, and sometimes better, test-set accuracies on several datasets. In Table 1, we report the test-set accuracies of several dataset-network pairs for standard DNNs and their Gabor-layered counterparts. We show the absolute difference in performance in the last column.

Moreover, in Figure 3, we provide a visual comparison between the patterns learned by AlexNet in its original implementation [21] and those learned in the Gabor-layered version of AlexNet (trained on CIFAR100). We observe that filters in the Gabor layer converge to filters that are similar to those found in the original implementation of AlexNet, where we observe blob-like structures and oriented edges and bars of various sizes. Note that both sets of filters are, in turn, similar to filter banks traditionally used in computer vision, as those proposed by Leung and Malik [26]. Next, we show that the Gabor-layered networks highlighted in Table 1 not only achieve test set accuracies as high as those of standard DNNs, but also enjoy better robustness properties for free.

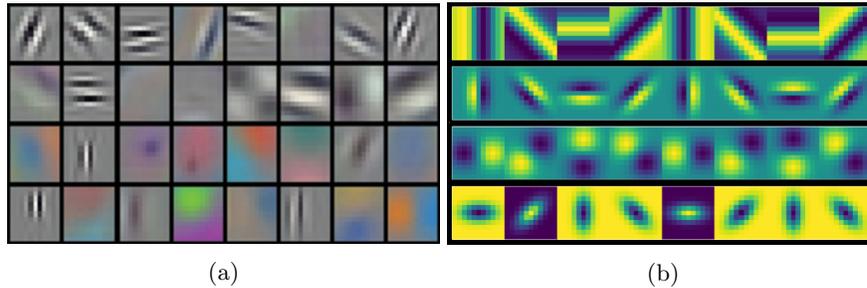(a)                                                    (b)

Fig. 3: **Comparison between filters learned by AlexNet and by Gabor-layered AlexNet.** (a) Various filters learned in the first convolutional layer of AlexNet in its original implementation [21]. (b) Several filters learned in the Gabor-layered version of AlexNet. Each column in (b) is a different orientation of the same filter, while each row represents a different set of parameters of the Gabor function. Note that standard convolutional layers use multiple-channeled filters, while Gabor layers use one-dimensional filters. Compellingly, both sets of filters present blobs and oriented edges and bars of various sizes and intensities. Best viewed in color.

### 4.4   Distribution of Singular Values

The Lipschitz constant of a network is an important quantity in the study of the network's robustness properties, since it is a measure of the variation of the network's predictions under input perturbations. Hence, in this work we study the distribution of singular values and, as a consequence, the Lipschitz constant of the filters of the layers in which we introduced Gabor layers instead of regular convolutional layers, in a similar fashion to [11]. In Figure 4 we report box-plots of the distributions of singular values for the first layer of LeNet trained on MNIST, and the first three layers of VGG16 trained on CIFAR100. Each plot shows the distribution of singular values of the standard architectures (S), Gabor-layered architectures (G), and Gabor-layered architectures trained *with* the regularizers (G+r) proposed in Equations (2) and (3).

Figure 4 demonstrates that the singular values of the filters in Gabor layers tend to be concentrated around smaller values, while also being distributed in smaller ranges than those of their standard convolutional counterparts, as shown by the interquartile ranges. Additionally, in most cases, the Lipschitz constant of the filters of Gabor layers, *i.e.* the top notch of each box-plot, is smaller than that of standard convolutional layers.

Moreover, we find that training Gabor-layered networks with the regularizer we introduced in Equations (2) and (3) incites further reduction in the singular values of the Gabor filters, as shown in Figure 4. For instance, the Gabor-layered version of LeNet trained on MNIST has a smaller interquartile range of the singular values, but still suffers from a large Lipschitz constant. However, by jointly introducing both the Gabor layer *and* the regularizer, the Lipschitz constant
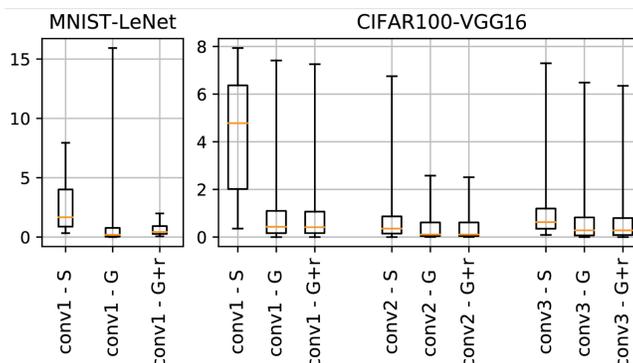
Fig. 4: **Box-plot representation of the distribution of singular values in layers of LeNet and VGG16.** Left: LeNet on MNIST. Right: VGG16 on CIFAR100. S: Standard; G: Gabor-layered; G+r: Gabor-layered with regularization. The top notch of each box-plot corresponds to the maximum value of the distribution, *i.e.* the Lipschitz constant of the layer.

decreases by a factor of almost 5. This reduction in the Lipschitz constant is consistent in all layers of VGG16 trained on CIFAR100.

In the following section, we show the bulk of our experiments, in which we investigate the effect on robustness of introducing Gabor layers into DNNs, by conducting a study across different architectures and datasets.

### 4.5    Robustness in Gabor-Layered Architectures

After observing significant differences in the distribution of singular values between standard convolutional layers and Gabor layers, we now study the impact on robustness that Gabor layers introduce. We study the robustness properties of different architectures trained on various datasets when Gabor layers are introduced in the first layers of each network. The modifications that we perform on each architecture are:

- **LeNet.** We replace the first layer with a Gabor layer with $p = 2$.
- **AlexNet.** We replace the first layer with a Gabor layer with $p = 7$.
- **WideResNet.** We replace the first layer with a Gabor layer with $p = 4$ for SVHN, and $p = 3$ for CIFAR100.
- **VGG16.** We replace the first three layers with Gabor layers with parameters $p = 3$, $p = 1$ and $p = 3$, respectively.

We present the details for the choice of $p$ in the **supplementary material**.

**Standard Architectures *vs.* Gabor-Layered Architectures.** We report the adversarial accuracies on both standard architectures and Gabor-layered architectures, where we refer to each with "S" and "G", respectively. Table 2

Table 2: **Adversarial accuracy comparison**. We compare Standard (S), Gabor-layered (G), and *regularized* Gabor-layered (G+r) architectures. For each attack strength ($\epsilon$), the highest performance is in **bold**; second-highest is underlined.

| $\epsilon$ | | $2/255$ | | | $8/255$ | | | $16/255$ | | |
|---|---|---|---|---|---|---|---|---|---|---|
| Dataset | Network | S | G | G+r | S | G | G+r | S | G | G+r |
| SVHN | WRN | 40.27 | 49.35 | **53.36** | 1.02 | 1.03 | **1.13** | 1.32 | **1.36** | 1.19 |
| SVHN | VGG16 | 57.86 | 62.88 | **64.03** | 5.84 | 14.57 | **15.99** | 2.33 | 7.98 | **8.88** |
| CIFAR10 | VGG16 | 34.22 | 37.60 | **38.07** | 23.63 | 30.11 | **30.69** | 13.88 | 19.50 | **19.93** |
| CIFAR100 | AN | **15.05** | 14.77 | 14.68 | 4.80 | 7.71 | **7.88** | 5.37 | 6.25 | **6.67** |
| CIFAR100 | WRN | 4.52 | 8.06 | **9.33** | 2.38 | 2.92 | **3.07** | 1.65 | 2.4 | **2.59** |
| CIFAR100 | VGG16 | 27.22 | 31.12 | **31.68** | 18.46 | 25.82 | **26.64** | 10.49 | 15.40 | **16.06** |

presents results on SVHN, CIFAR10 and CIFAR100, and Table 3 presents results on MNIST. We observe that Gabor-layered architectures consistently outperform their standard counterparts across datasets, and can provide up to a 9% boost in adversarial robustness. For instance, with $\epsilon = 8/255$ attacks, introducing Gabor layers into VGG16 boosts adversarial accuracy from 23.63 to 30.11 (6.48% relative increment) and from 5.84 to 14.57 (8.73% relative increment) on CIFAR10 and SVHN (Table 2), respectively. For LeNet on MNIST, under an $\epsilon = 0.2$ attack, introducing Gabor layers can boost adversarial accuracy from 4.39% to 7.94% (80% relative increment). Additionally, we report the flip rates for these experiments in the **supplementary material**. On the flip rates, we conduct a similar analysis, which yields equivalent conclusions: introducing Gabor layers leads to boosts in robustness.

Moreover, to explore whether the robustness effect of using Gabor-layers holds in the regime of large-scale datasets, we conduct an experiment on ImageNet [40], which we report in the **supplementary material**. Due to limitations in computational resources, we conduct attacks only for $\epsilon \in \{8/255, 16/255\}$. In this experiment, we again measure adversarial accuracy and flip rates, and observe similar gains in robustness in terms of flip rates.

It is worthwhile to note that the increase in robustness we observe in the Gabor-layered networks came *solely* from an architectural change, *i.e.* replacing convolutional layers with Gabor layers, without there being any other modification. Our experimental results demonstrate that: (**1**) Simply introducing Gabor filters, in the form of Gabor layers, as low-level feature extractors in DNNs provides beneficial effects to robustness. (**2**) Such robustness improvements are consistent across datasets and architectures. Inspired by these results, we now investigate the robustness effects of using our proposed regularization approach, as proposed in Equations (2) and (3).

**Robustness Effects of Introducing Regularization.** To better control robustness and to regularize the Lipschitz constant we derived in Theorem 1,

Table 3: **Adversarial accuracy comparison on MNIST**. We compare Standard (S), Gabor-layered (G), and *regularized* Gabor-layered (G+r) architectures. For each attack strength ($\epsilon$), the highest performance is in **bold**; second-highest is <u>underlined</u>.

| $\epsilon$ | | 0.1 | | | 0.2 | | | 0.3 | |
|---|---|---|---|---|---|---|---|---|---|
| Dataset Network | S | G | G+r | S | G | G+r | S | G | G+r |
| MNIST   LeNet | 80.04 | <u>80.58</u> | **88.42** | 4.39 | <u>7.94</u> | **22.69** | 0.44 | **0.78** | <u>0.76</u> |

we proposed two regularizers for the loss function as per Equations (2) and (3). As we noted in Subsection 4.4 (refer to Figure 4), Gabor-layered architectures inherently tend to have lower singular values than their standard counterparts. Additionally, upon training the same architectures *with* the proposed regularizer, we observe that Gabor layers tend to enjoy further reduction in the Lipschitz constant. These results suggest that such architectures may have enhanced robustness properties as a consequence.

To assess the role of the proposed regularizer on robustness, we train Gabor-layered architectures from scratch following the same parameters from Subsection 4.1 and include the regularizer. We present the results in Tables 2 and 3, where we refer to these regularized architectures as "G+r". We observe that, in most cases, adding the regularizer improves adversarial accuracy. For instance, for LeNet on MNIST, the regularizer improves adversarial accuracy over the Gabor-layered architecture without any regularization by 8% and 14% with $\epsilon = {}^2/255, {}^8/255$ attacks, respectively. This improvement is still present in more challenging datasets. For instance, for VGG16 on SVHN under attacks with $\epsilon = {}^2/255$ and $\epsilon = {}^8/255$, we observe increments of over 1% from the regularized architecture with respect to its non-regularized equivalent. For the rest of the architectures and datasets, we observe modest but, nonetheless, sustained increments in performance. We report the flip rates for these experiments in the **supplementary material**. The conclusions obtained from analyzing the flip rates are analogous to what we conclude from the adversarial accuracies: applying regularization on these layers provides minor but consistent improvements in robustness.

It is worthy to note that, although the implementation of the regularizer is trivial and that we perform optimization including the regularizer for the same number of epochs as regular training, our results still show that it is possible to achieve desirable robustness properties. We expect that substantial modifications and optimization heuristics can be applied in the training procedure, with aims at stronger exploitation of the insights we have provided here, and most likely resulting in more significant boosts in robustness.

### 4.6   Effects of Adversarial Training

Adversarial training [30] has become the standard approach for tackling robustness. In these experiments, we investigate how Gabor layers interact with

Table 4: Adversarial accuracy with $\epsilon = {}^8\!/_{255}$.

| | | | CIFAR10 | | CIFAR100 | |
|---|---|---|---|---|---|---|
| Gabor | Reg. | Adv. training | AlexNet | VGG16 | AlexNet | VGG16 |
| | | ✓ | 19.24 | 41.95 | 13.48 | 18.49 |
| ✓ | | ✓ | 20.26 | 42.41 | 10.94 | **19.66** |
| ✓ | ✓ | ✓ | **22.15** | **44.02** | **13.62** | 19.41 |

adversarial training. We study whether the increments in robustness we observed when introducing Gabor layers can still be observed in the regime of adversarially-trained models. To study such interaction, we adversarially train standard, Gabor-layered and regularized Gabor-layered architectures, and then compare their robustness properties. We use the adversarial training described in [43], with 8 mini-batch replays, and $\epsilon = {}^8\!/_{255}$.

In Table 4, we report adversarial accuracies for AlexNet on CIFAR10 and CIFAR100 under $\epsilon = {}^8\!/_{255}$ attacks. Even in the adversarially trained networks-regime, our experiments show that **(1)** Gabor-layered architectures outperform their standard counterparts, and **(2)** regularization of Gabor-layered architectures provides substantial improvements in robustness with respect to their non-regularized equivalents. Such results demonstrate that Gabor layers represent an orthogonal approach towards robustness and, hence, that Gabor layers and adversarial training can be jointly harnessed for enhancing the robustness of DNNs.

The results we present here are empirical evidence that using closed form expressions for filter-generating functions in convolutional layers can be exploited for the purpose of increasing robustness in DNNs. We refer the interested reader to the **supplementary material** for the rest of the experimental results.

## 5    Conclusions

In this work, we study the effects in robustness of architectural changes in convolutional neural networks. We show that introducing Gabor layers consistently improves the robustness across various neural network architectures and datasets. We also show that the Lipschitz constant of the filters in these Gabor layers tends to be lower than that of traditional filters, which was theoretically and empirically shown to be beneficial to robustness [11]. Furthermore, theoretical analysis allows us to find a closed form expression for a Lipschitz constant of the Gabor filters. We then leverage this expression as a regularizer in the pursuit of enhanced robustness, and validate its usefulness experimentally. Finally, we study the interaction between Gabor layers, our regularizer, and adversarial training, and show that the benefits of using Gabor layers are still observed when deep learning models are specifically trained for the purpose of adversarial robustness, showing that Gabor layers can be jointly used with adversarial training for further enhancements in robustness.

# References

1. Alekseev, A., Bobe, A.: Gabornet: Gabor filters with learnable parameters in deep convolutional neural networks. arXiv:1904.13204 (2019)
2. Atanov, A., Ashukha, A., Struminsky, K., Vetrov, D., Welling, M.: The deep weight prior (2018)
3. Bai, M., Urtasun, R.: Deep watershed transform for instance segmentation. In: IEEE Conference on Computer Vision and Pattern Recognition (CVPR) (2017)
4. Bibi, A., Ghanem, B., Koltun, V., Ranftl, R.: Deep layers as stochastic solvers. In: International Conference on Learning Representations (ICLR) (2019)
5. Cao, Y., Xiao, C., Cyr, B., Zhou, Y., Park, W., Rampazzi, S., Chen, Q.A., Fu, K., Mao, Z.M.: Adversarial sensor attack on lidar-based perception in autonomous driving. Proceedings of the 2019 ACM SIGSAC Conference on Computer and Communications Security (2019)
6. Carlini, N., Wagner, D.: Towards evaluating the robustness of neural networks. In: IEEE Symposium on Security and Privacy (SP) (2017)
7. Chen, Y., Zhu, L., Ghamisi, P., Jia, X., Li, G., Tang, L.: Hyperspectral images classification with gabor filtering and convolutional neural network. IEEE Geoscience and Remote Sensing Letters (2017)
8. Chengjun Liu, Wechsler, H.: Independent component analysis of gabor features for face recognition. IEEE Transactions on Neural Networks (2003)
9. Chernikova, A., Oprea, A., Nita-Rotaru, C., Kim, B.: Are self-driving cars secure? evasion attacks against deep neural networks for steering angle prediction. IEEE Security and Privacy Workshops (SPW) (2019)
10. Chollet, F.: Xception: Deep learning with depthwise separable convolutions. IEEE Conference on Computer Vision and Pattern Recognition (CVPR) (2017)
11. Cisse, M., Bojanowski, P., Grave, E., Dauphin, Y., Usunier, N.: Parseval networks: Improving robustness to adversarial examples. In: International Conference on Machine Learning (ICML) (2017)
12. Farabet, C., Couprie, C., Najman, L., LeCun, Y.: Learning hierarchical features for scene labeling. IEEE transactions on pattern analysis and machine intelligence **35**(8), 1915–1929 (2012)
13. Gabor, D.: Theory of communication. part 1: The analysis of information. Journal of the Institution of Electrical Engineers-Part III: Radio and Communication Engineering (1946)
14. Goldstein, T., Osher, S.: The split bregman method for l1-regularized problems. SIAM Journal on Imaging Sciences (2009)
15. Goodfellow, I.J., Shlens, J., Szegedy, C.: Explaining and harnessing adversarial examples. International Conference on Learning Representations (ICLR) (2015)
16. He, K., Zhang, X., Ren, S., Sun, J.: Deep residual learning for image recognition. IEEE Conference on Computer Vision and Patter Recognition (CVPR) (2016)
17. Hinton, G., Deng, L., Yu, D., Dahl, G.E., Mohamed, A., Jaitly, N., Senior, A., Vanhoucke, V., Nguyen, P., Sainath, T.N., Kingsbury, B.: Deep neural networks for acoustic modeling in speech recognition: The shared views of four research groups. IEEE Signal Processing Magazine (2012)

18. Hubel, D.H., N., W.T.: Receptive fields of single neurons in the cat's striate cortex. The Journal of Physiology (1959)
19. Julesz, B.: Textons, the elements of texture perception, and their interactions. Nature (1981)
20. Krizhevsky, A., Hinton, G.: Learning multiple layers of features from tiny images. Tech. rep., Citeseer (2009)
21. Krizhevsky, A., Sutskever, I., Hinton, G.E.: Imagenet classification with deep convolutional neural networks. In: Conference on Neural Information Processing Systems (NeurIPS12) (2012)
22. Lampinen, J., Oja, E.: Distortion tolerant pattern recognition based on self-organizing feature extraction. IEEE Transactions on Neural Networks (1995)
23. LeCun, Y.: The mnist database of handwritten digits. http://yann. lecun. com/exdb/mnist/ (1998)
24. LeCun, Y., Bengio, Y., Hinton, G.: Deep learning. nature (2015)
25. LeCun, Y., Haffner, P., Bottou, L., Bengio, Y.: Object recognition with gradient-based learning. In: Shape, Contour and Grouping in Computer Vision. Springer (1999)
26. Leung, T., Malik, J.: Representing and recognizing the visual appearance of materials using three-dimensional textons. International Journal of Computer Vision (IJCV) (2001)
27. Lowe, D.G.: Object recognition from local scale-invariant features. In: Proceedings of the seventh IEEE international conference on computer vision. vol. 2, pp. 1150–1157. Ieee (1999)
28. Lowe, D.G.: Distinctive image features from scale-invariant keypoints. International journal of computer vision **60**(2), 91–110 (2004)
29. Luan, S., Chen, C., Zhang, B., Han, J., Liu, J.: Gabor convolutional networks. IEEE Transactions on Image Processing (2018)
30. Madry, A., Makelov, A., Schmidt, L., Tsipras, D., Vladu, A.: Towards deep learning models resistant to adversarial attacks. In: International Conference on Learning Representations (ICLR) (2018)
31. Marr, D.: Vision: A computational investigation into the human representation and processing of visual information (1982)
32. Montufar, G., Pascanu, R., Cho, K., Bengio, Y.: On the number of linear regions of deep neural networks. Advances in Neural Information Processing Systems (NeurIPS) (2014)
33. Moosavi-Dezfooli, S.M., Fawzi, A., Frossard, P.: Deepfool: A simple and accurate method to fool deep neural networks. IEEE Conference on Computer Vision and Pattern Recognition (CVPR) (2016)
34. Namuduri, K.R., Mehrotra, R., Ranganathan, N.: Edge detection models based on gabor filters. In: International Conference on Pattern Recognition. Conference C: Image, Speech and Signal Analysis, (1992)
35. Novak, C.L., Shafer, S.A., et al.: Anatomy of a color histogram.
36. Ouyang, W., Wang, X.: Joint deep learning for pedestrian detection. In: IEEE International Conference on Computer Vision (ICCV) (2013)
37. Papernot, N., McDaniel, P., Wu, X., Jha, S., Swami, A.: Distillation as a defense to adversarial perturbations against deep neural networks. In: IEEE Symposium on Security and Privacy (SP) (2016)
38. Poggi, M., Mattoccia, S.: A wearable mobility aid for the visually impaired based on embedded 3d vision and deep learning. In: 2016 IEEE Symposium on Computers and Communication (ISCC) (2016)

39. Rumelhart, D.E., Hinton, G.E., Williams, R.J., et al.: Learning representations by back-propagating errors. Cognitive modeling (1988)
40. Russakovsky, O., Deng, J., Su, H., Krause, J., Satheesh, S., Ma, S., Huang, Z., Karpathy, A., Khosla, A., Bernstein, M., et al.: Imagenet large scale visual recognition challenge. International Journal of Computer Vision (IJCV) (2015)
41. Sarwar, S.S., Panda, P., Roy, K.: Gabor filter assisted energy efficient fast learning convolutional neural networks. CoRR (2017)
42. Sedghi, H., Gupta, V., Long, P.M.: The singular values of convolutional layers. In: International Conference on Learning Representations (ICLR) (2019)
43. Shafahi, A., Najibi, M., Ghiasi, A., Xu, Z., Dickerson, J., Studer, C., Davis, L.S., Taylor, G., Goldstein, T.: Adversarial training for free! (2019)
44. Simonyan, K., Zisserman, A.: Very deep convolutional networks for large-scale image recognition. In: International Conference on Learning Representations (ICLR) (2015)
45. Su, Y.M., Wang, J.F.: A novel stroke extraction method for chinese characters using gabor filters. Pattern Recognition (2003)
46. Sun, J., Li, H., Xu, Z., et al.: Deep admm-net for compressive sensing mri. In: Advances in Neural Information Systems (NeurIPs) (2016)
47. Wong, E., Rice, L., Kolter, J.Z.: Fast is better than free: Revisiting adversarial training (2020)
48. Xu, K., Liu, S., Zhao, P., Chen, P., Zhang, H., Erdogmus, D., Wang, Y., Lin, X.: Structured adversarial attack: Towards general implementation and better interpretability. International Conference on Learning Representations (ICLR) (2019)
49. Yao, H., Chuyi, L., Dan, H., Weiyu, Y.: Gabor feature based convolutional neural network for object recognition in natural scene. In: International Conference on Information Science and Control Engineering (ICISCE) (2016)
50. Zagoruyko, S., Komodakis, N.: Wide residual networks. In: British Machine Vision Conference (BMVC) (2016)
51. Zhang, D., Zhang, T., Lu, Y., Zhu, Z., Dong, B.: You only propagate once: Accelerating adversarial training via maximal principle. In: Advances in Neural Information Processing Systems (NeuIPS) (2019)
52. Zhang, J., Ghanem, B.: Deep learning. IEEE Conference on Computer Vision and Patter Recognition (CVPR) (2017)
53. Zhong, Z., Jin, L., Xie, Z.: High performance offline handwritten chinese character recognition using googlenet and directional feature maps. In: International Conference on Document Analysis and Recognition (ICDAR) (2015)

# Supplementary Material

Next, we present the Supplementary Material for the paper "Gabor Layers Enhance Network Robustness". In Section 6 we report the details for the choice of the $p$ parameter of the Gabor layers for the various experiments we present in the paper; Section 7 shows the flip rates for all the experiments we reported; Section 8 presents the results for one experiment on ImageNet [40], both in terms of adversarial accuracy and flip rate; in Section 9 we present several comparisons of distributions of singular values of standard convolutional layers *vs.* Gabor layers; finally, Section 10 depicts visualizations of the filters learned in the Gabor layers when various architectures are trained on numerous datasets.

## 6   Number of families $p$ search space

We report the details for the search for the $p$ parameter in Tables 5 through 11.

Table 5: **Search Space for LeNet on MNIST** Test set accuracies and flip rates on MNIST with LeNet. Only the first convolutional layer was enhanced with $p$ families.

| $\epsilon$ | Standard | | $p = 2$ | | $p = 3$ | |
|---|---|---|---|---|---|---|
| | Accuracy | Flip Rate | Accuracy | Flip Rate | Accuracy | Flip Rate |
| 0 | 99.22 | - | 99.03 | - | 99.10 | - |
| 0.1 | 80.04 | 19.53 | 80.58 | 18.88 | 62.98 | 36.54 |
| 0.2 | 4.39 | 95.47 | 7.94 | 91.85 | 2.22 | 97.67 |
| 0.3 | 0.44 | 99.7 | 0.78 | 99.24 | 0.41 | 99.68 |

Table 6: **Search Space for AlexNet on CIFAR100.** Test set accuracies comparison with the standard AlexNet and his enhanced versions. Only the first convolutional layer was enhanced. S stands for standard.

| $\epsilon$ | S | **Accuracy** $p = 2$ | $p = 3$ | $p = 4$ | $p = 5$ | $p = 6$ | $p = 7$ | $p = 8$ | $p = 9$ |
|---|---|---|---|---|---|---|---|---|---|
| 0 | 46.48 | 42.14 | 42.81 | 41.99 | 43.19 | 42.32 | 45.15 | 42.69 | 43.36 |
| $2/255$ | 15.08 | 11.84 | 13.01 | 13.41 | 13.62 | 13.55 | 14.77 | 12.82 | 13.16 |
| $8/255$ | 4.80 | 6.82 | 7.40 | 6.74 | 7.81 | 7.01 | 7.71 | 7.86 | 7.73 |
| $16/255$ | 5.37 | 5.92 | 6.12 | 5.75 | 6.23 | 5.34 | 6.25 | 6.07 | 6.29 |

Table 7: **Search Space for VGG16 on CIFAR100 – First Layer.** Test set accuracies comparison with the standard VGG16 and his enhanced versions. The first Gabor layer number of families $p_1$ is explored. S stands for standard.

|  |  | Accuracy | | | | | |
|---|---|---|---|---|---|---|---|
| $\epsilon$ | S | $p_1=1$ | $p_1=2$ | $p_1=3$ | $p_1=4$ | $p_1=5$ | $p_1=6$ |
| 0 | 67.54 | 65.82 | 67.35 | 67.05 | 68.65 | 68.10 | 66.23 |
| $^2/_{255}$ | 27.22 | 28.35 | 27.89 | 27.31 | 27.83 | 27.32 | 24.83 |
| $^8/_{255}$ | 18.46 | 23.05 | 20.78 | 20.05 | 20.66 | 19.44 | 18.95 |
| $^{16}/_{255}$ | 10.49 | 13.92 | 12.18 | 11.30 | 11.02 | 10.64 | 11.29 |

Table 8: **Search Space for VGG16 on CIFAR100 – Second Layer.** Test set accuracies comparison with the standard VGG16 and his enhanced versions. The first convolutional layer is modified with $p_1 = 3$ families. The second Gabor layer number of families $p_2$ is explored. S stands for standard.

|  |  |  | Accuracy | | | | |
|---|---|---|---|---|---|---|---|
| $\epsilon$ | S | $p_1=3$ | $p_2=1$ | $p_2=2$ | $p_2=3$ | $p_2=4$ | $p_2=5$ |
| 0 | 67.54 | 67.05 | 64.10 | 67.68 | 66.67 | 65.06 | 64.81 |
| $^2/_{255}$ | 27.22 | 27.31 | 31.36 | 16.57 | 27.43 | 28.25 | 27.92 |
| $^8/_{255}$ | 18.46 | 20.05 | 26.09 | 10.41 | 21.15 | 22.50 | 22.00 |
| $^{16}/_{255}$ | 10.49 | 11.30 | 14.97 | 5.10 | 11.30 | 12.64 | 13.11 |

Table 9: **Search Space for VGG16 on CIFAR100 – Third Layer.** Test set accuracies comparison with the standard VGG16 and his enhanced versions. The first and second convolutional layers are modified with $p_1 = 3$ and $p_2 = 1$ families respectively. The third Gabor layer number of families $p_3$ is explored. S stands for standard.

|  |  |  | Accuracy | |
|---|---|---|---|---|
| $\epsilon$ | S | $p_1=3, p_2=1$ | $p_3=2$ | $p_3=3$ |
| 0 | 67.54 | 64.10 | 63.74 | 64.49 |
| $^2/_{255}$ | 27.22 | 31.36 | 17.13 | 31.12 |
| $^8/_{255}$ | 18.46 | 26.09 | 14.67 | 25.82 |
| $^{16}/_{255}$ | 10.49 | 14.97 | 9.20 | 15.40 |

Table 10: **Search Space for Wide-ResNet on SVHN – First Layer.** Test set accuracies comparison with the standard Wide-ResNet and his enhanced versions. The first Gabor layer number of families $p_1$ is explored. S stands for standard.

| | | **Accuracy** | | | | |
|---|---|---|---|---|---|---|
| $\epsilon$ | S | $p_1 = 2$ | $p_1 = 3$ | $p_1 = 4$ | $p_1 = 5$ | $p_1 = 6$ |
| 0 | 96.62 | 96.93 | 96.72 | 96.70 | 96.65 | 96.73 |
| $^2/_{255}$ | 40.27 | 45.84 | 41.37 | 49.35 | 43.56 | 45.66 |
| $^8/_{255}$ | 1.03 | 0.93 | 1.09 | 1.03 | 1.08 | 1.03 |
| $^{16}/_{255}$ | 1.32 | 1.11 | 1.32 | 1.42 | 1.38 | 1.28 |

Table 11: **Search Space for Wide-ResNet on SVHN – Second Layer.** Test set accuracies comparison with the standard Wide-ResNet and his enhanced versions. The first convolutional layer is modified with $p_1 = 4$ families. The second Gabor layer number of families $p_2$ is explored. S stands for standard.

| | | **Accuracy** | | | | |
|---|---|---|---|---|---|---|
| $\epsilon$ | S | $p_1 = 4$ | $p_2 = 1$ | $p_2 = 2$ | $p_2 = 3$ | $p_2 = 4$ |
| 0 | 96.62 | 96.70 | 96.67 | 96.71 | 96.60 | 96.71 |
| $^2/_{255}$ | 40.27 | 49.35 | 41.83 | 44.50 | 44.11 | 44.41 |
| $^8/_{255}$ | 1.03 | 1.03 | 1.01 | 0.99 | 0.98 | 1.03 |
| $^{16}/_{255}$ | 1.32 | 1.42 | 1.24 | 1.31 | 1.28 | 1.26 |

Table 12: **Flip rates comparison**. We compare Standard (S), Gabor-layered (G), and *regularized* Gabor-layered (G+r) architectures. For each attack strength ($\epsilon$), the lowest flip rate is in **bold**; second-lowest is underlined.

| $\epsilon$ | | $^2/_{255}$ | | | $^8/_{255}$ | | | $^{16}/_{255}$ | | |
|---|---|---|---|---|---|---|---|---|---|---|
| Dataset | Network | S | G | G+r | S | G | G+r | S | G | G+r |
| SVHN | WRN | 57.22 | 48.13 | **44.19** | 98.30 | **98.13** | 98.17 | 99.06 | 99.02 | **99.01** |
| SVHN | VGG16 | 39.53 | 34.11 | **32.62** | 92.82 | 83.83 | **82.33** | 97.40 | 91.56 | **90.78** |
| CIFAR10 | VGG16 | 60.03 | 56.42 | **55.71** | 74.51 | 67.83 | **67.46** | 86.47 | 80.84 | **80.18** |
| CIFAR100 | AN | 56.88 | 49.73 | **49.15** | 81.38 | **72.82** | 72.87 | 87.77 | **83.00** | 83.08 |
| CIFAR100 | WRN | 82.05 | 76.52 | **74.72** | 93.82 | 92.58 | **92.25** | 96.94 | 96.59 | **96.35** |
| CIFAR100 | VGG16 | 57.05 | **50.94** | 51.05 | 77.94 | 68.95 | **68.45** | 90.48 | 85.75 | **85.16** |

Table 13: **Flip rates comparison on MNIST**. We compare Standard (S), Gabor-layered (G), and *regularized* Gabor-layered (G+r) architectures on MNIST. For each attack strength ($\epsilon$), the lowest flip rate is in **bold**; second-lowest is underlined.

| $\epsilon$ | | 0.1 | | | 0.2 | | | 0.3 | | |
|---|---|---|---|---|---|---|---|---|---|---|
| Dataset | Network | S | G | G+r | S | G | G+r | S | G | G+r |
| MNIST | LeNet | 19.53 | 18.88 | **11.05** | 95.47 | 91.85 | **77.27** | 99.70 | **99.24** | 99.64 |

## 7   Flip rates

In our work, we assess the robustness of Deep Neural Networks (DNNs) through adversarial accuracies (reported in the main document) and flip rates. Next, we report the flip rates for the main experiments of the paper. Table 12 shows results on SVHN, CIFAR10 and CIFAR100, and Table 13 presents results on MNIST.

## 8   ImageNet Results

We conduct adversarial attacks for $\epsilon \in \{^8/_{255}, ^{16}/_{255}\}$. In Table 14 we report the adversarial accuracies and flip rates for VGG16 [44] trained on ImageNet [40].

## 9   Singular Values

We report the distribution of singular values of the filters of up to the first three convolutional layers of several Convolutional Neural Network (CNN) architectures (LeNet [25], AlexNet [21], WideResNet [50], and VGG16 [44]) trained in

Table 14: **Adversarial accuracy and flip rate comparison for VGG16 on ImageNet.** We compare Standard (S) and Gabor-layered (G) architectures. For each attack strength ($\epsilon$), the best performance is in **bold**.

| $\epsilon$ | Adv. Accuracy | | Flip Rate | |
|---|---|---|---|---|
| | S | G | S | G |
| 0 | **71.20** | 68.90 | - | - |
| $8/255$ | **2.95** | 2.24 | 95.07 | **94.46** |
| $16/255$ | **3.33** | 3.15 | 97.37 | **96.68** |

various datasets (MNIST [23], CIFAR10, CIFAR100 [20], and ImageNet [40]). The singular values of the layers are computed following [42,4]. The largest singular value of each layer's filter corresponds to the filters' Lipschitz constant [4].

In each histogram plot we show the distribution of singular values of (1) the standard architecture, in blue, and (2) the Gabor-layered architecture, in orange. For some experiments, we also show the distribution of singular values of the Gabor-layered architecture *with* regularization, in green.

For visualization purposes, we set the upper x-limit of each histogram plot to the 95th percentile of the distribution with the largest maximum value.

Next, we list the dataset-architecture pairs, and the Figures in which its distributions are shown:

- **MNIST-LeNet.** Figure 5.
- **CIFAR100-AlexNet.** Figure 6.
- **CIFAR10-VGG16.** Figures 7 - 9.
- **CIFAR100-VGG16.** Figures 10 - 12.
- **ImageNet-VGG16.** Figures 13 - 15.

In most cases, the distribution of singular values of the Gabor-layered version of the networks tends to be around smaller values, usually with high peaks between 0 and 0.5, than that of the standard network.

In terms of the Lipschitz constant of the layers, in most cases we observe that the Gabor-layered versions of the layers have lower Lipschitz constants, and that applying regularization results in even lower Lipschitz constants.

## 10    Filter visualizations

We report visualizations of the filters of the first convolutional layer of some of the architectures we experimented with. Figures 16 through 23 depict the filters. For the standard convolutional layers, we visualize each filter as an RGB image, where each "channel" of the filter is scaled to be between 0 and 1.

For the Gabor layers, each filter has 1 channel and, hence, we visualize it as is. We show the filters that share the Gabor function $G_\theta$ in the same row; while each column corresponds to one of the 8 $\alpha$-scaled rotations of the filter.

We report the filters for:

– **MNIST-LeNet.** Figures 16 - 18.
– **CIFAR100-AlexNet.** Figures 19 - 21.
– **ImageNet-VGG16.** Figures 22 - 23.

We note that, for LeNet on MNIST, the Gabor-layered version, without regularization, already has filters that strongly resemble a Dirac-delta function and that, as reported in the paper, this network already shows improvements in terms of robustness. Furthermore, when applying regularization, we observe that all filters in the layer become Dirac-delta functions (see Figure 18). Again, as reported in the paper, regularization showed large gains in robustness.

The filters from AlexNet are useful for visualizing the modeling capabilities of Gabor functions, as shown in Figure 20, where we observe blob-like patterns, and also oriented and scaled edges and bars. These patterns, while simpler than those of the standard convolutional layers (see Figure 19), provide on-pair accuracy with such layers, while also providing gains in robustness. For the case of AlexNet, however, we observe that regularization has virtually no impact in the form of the filters that are learnt (see Figure 21).

In Figure 22 we report the filters learnt by the standard version of VGG16 on ImageNet. Patterns are, however, not straightforward to visualize in these filters. By construction, the filters learnt by the Gabor-layered version of VGG16 when fine-tuned on ImageNet are more visually-appealing, as shown in Figure 23. Note, also, that some Gabor filters have strong similarities between one another.

Fig. 5: **Distribution of singular values for the first layer of LeNet trained on MNIST.** The legend shows the largest singular value, *i.e.* the Lipschitz constant of the layer.
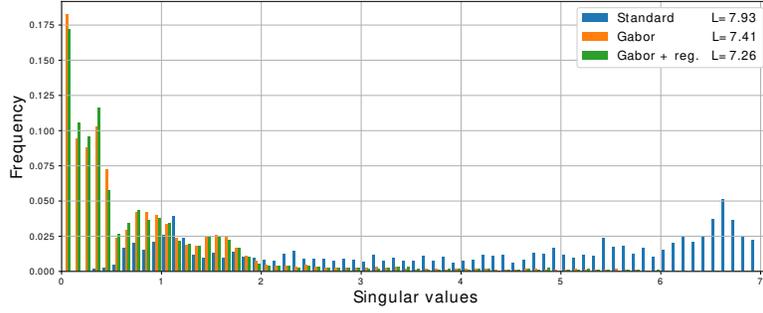


Fig. 6: **Distribution of singular values for the first layer of AlexNet trained on CIFAR100.** The legend shows the largest singular value, *i.e.* the Lipschitz constant of the layer.
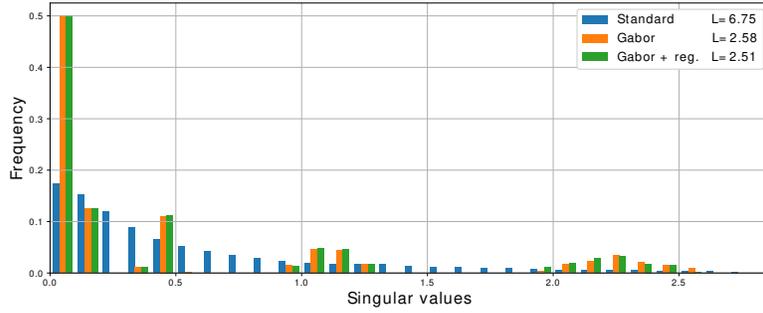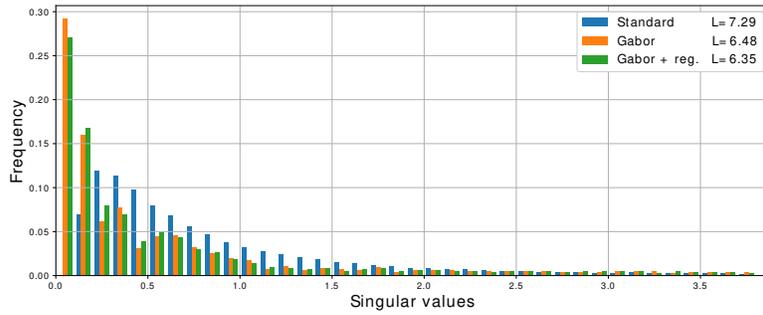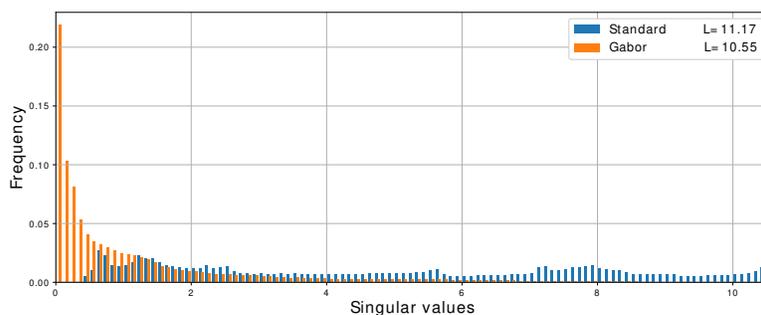
Fig. 7: **Distribution of singular values for the first layer of VGG16 trained on CIFAR10.** The legend shows the largest singular value, *i.e.* the Lipschitz constant of the layer.
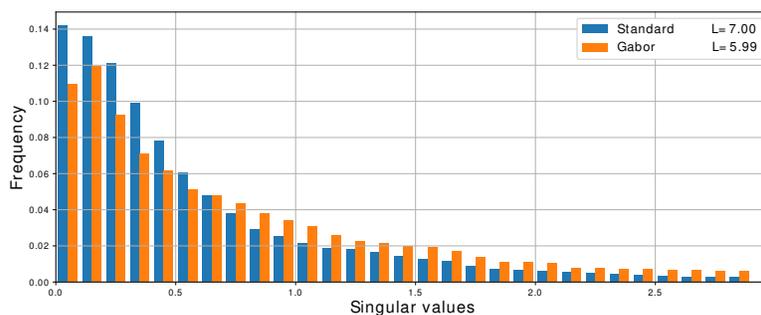


Fig. 8: **Distribution of singular values for the second layer of VGG16 trained on CIFAR10.** The legend shows the largest singular value, *i.e.* the Lipschitz constant of the layer.
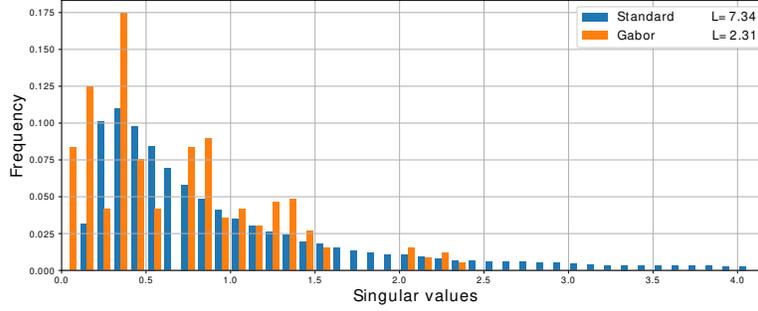


Fig. 9: **Distribution of singular values for the third layer of VGG16 trained on CIFAR10.** The legend shows the largest singular value, *i.e.* the Lipschitz constant of the layer.

Fig. 10: **Distribution of singular values for the first layer of VGG16 trained on CIFAR100.** The legend shows the largest singular value, *i.e.* the Lipschitz constant of the layer.



Fig. 11: **Distribution of singular values for the second layer of VGG16 trained on CIFAR100.** The legend shows the largest singular value, *i.e.* the Lipschitz constant of the layer.



Fig. 12: **Distribution of singular values for the third layer of VGG16 trained on CIFAR100.** The legend shows the largest singular value, *i.e.* the Lipschitz constant of the layer.

Fig. 13: **Distribution of singular values for the first layer of VGG16 trained on ImageNet.** The legend shows the largest singular value, *i.e.* the Lipschitz constant of the layer. Note that the Gabor-layered version was fine-tuned, starting from ImageNet-pretrained weights, due to computational constraints.
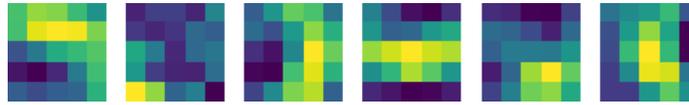


Fig. 14: **Distribution of singular values for the second layer of VGG16 trained on ImageNet.** The legend shows the largest singular value, *i.e.* the Lipschitz constant of the layer. Note that the Gabor-layered version was fine-tuned, starting from ImageNet-pretrained weights, due to computational constraints.

Fig. 15: **Distribution of singular values for the third layer of VGG16 trained on ImageNet.** The legend shows the largest singular value, *i.e.* the Lipschitz constant of the layer. Note that the Gabor-layered version was fine-tuned, starting from ImageNet-pretrained weights, due to computational constraints.



Fig. 16: **Standard LeNet-MNIST filters from the first convolutional layer.** The 6 grayscale filters are visualized.



Fig. 17: **Gabor-layered LeNet-MNIST filters from the first convolutional layer.** Each of the rows in the figure corresponds to each of the 7 families of filters of the layer. Each column is one of the 8 $\alpha$-scaled rotations of the filter. Note that the filters in the second row resemble a Dirac-delta function.

Fig. 18: **Gabor-layered LeNet-MNIST filters from the first convolutional layer with regularization.** Each of the rows in the figure corresponds to each of the 7 families of filters of the layer. Each column is one of the 8 $\alpha$-scaled rotations of the filter. Note that regularization enforces the filters to be Dirac-deltas and, as reported in the paper, this change results in large gains in robustness.
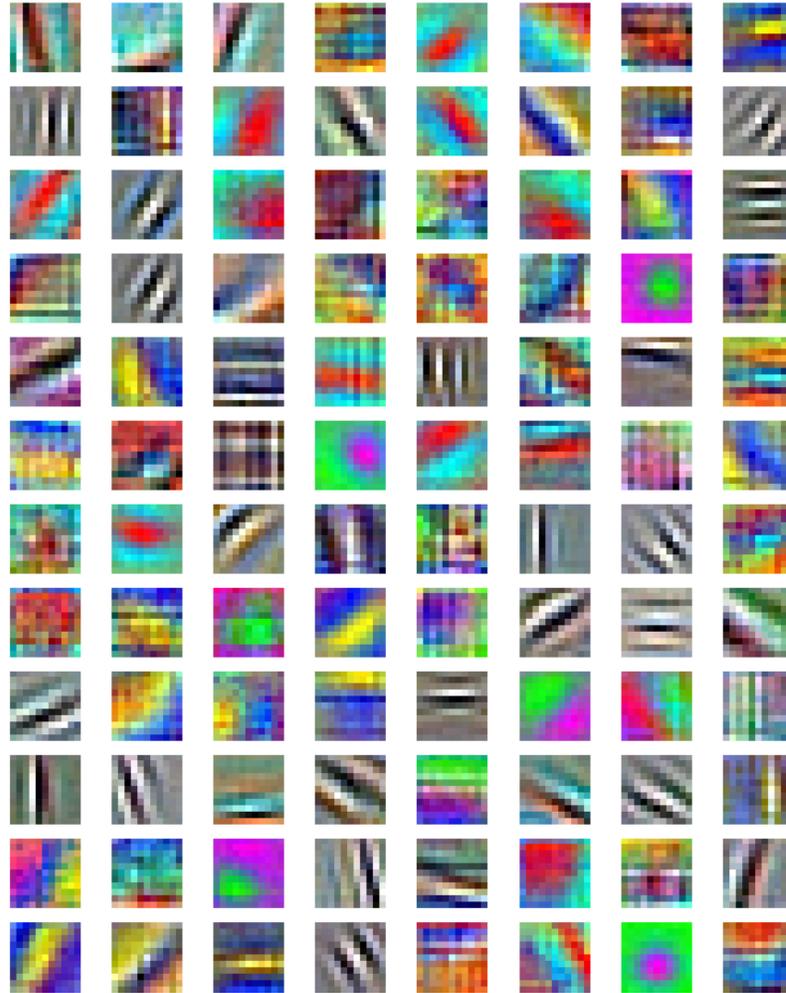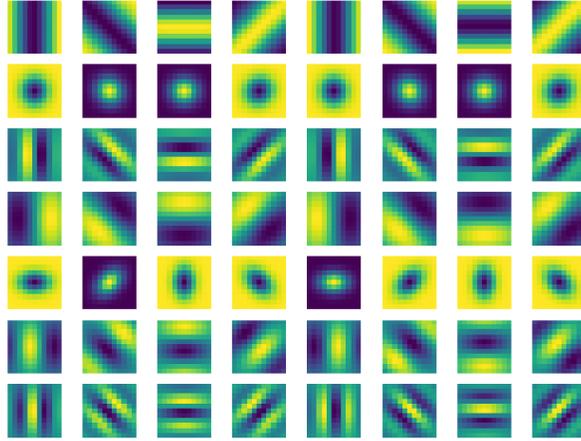
Fig. 19: **Standard AlexNet-CIFAR100 filters from the first convolutional layer.** The 96 filters are visualized by concatenating the RGB channels and mapping values from 0 to 1.

Fig. 20: **Gabor-layered AlexNet-CIFAR100 filters from the first convolutional layer.** Each of the rows in the figure corresponds to each of the 7 families of filters of the layer. Each column is one of the 8 $\alpha$-scaled rotations of the filter.
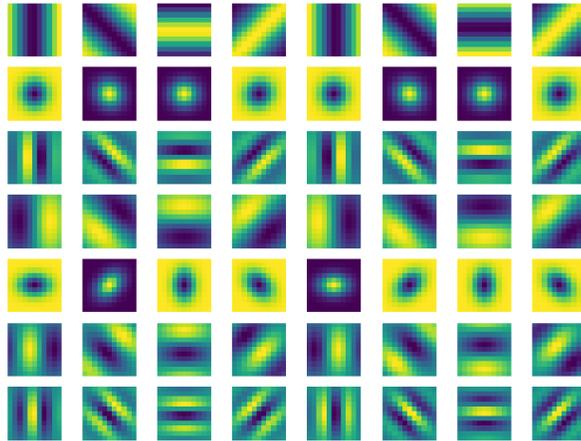


Fig. 21: **Gabor-layered AlexNet-CIFAR100 filters from the first convolutional layer with regularization.** Each of the rows in the figure corresponds to each of the 7 families of filters of the layer. Each column is one of the 8 $\alpha$-scaled rotations of the filter.
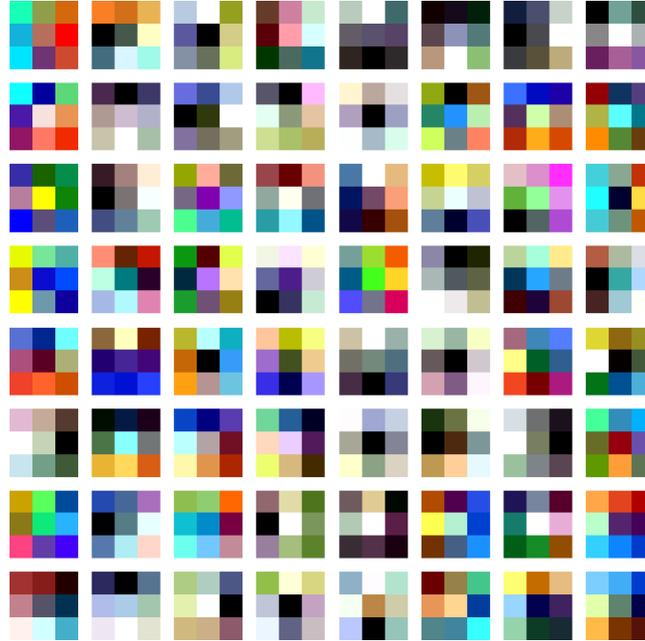
Fig. 22: **Standard VGG16-ImageNet filters from the first convolutional layer.** The 64 filters are visualized by concatenating the RGB channels and mapping values from 0 to 1.
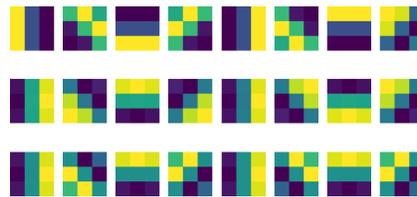


Fig. 23: **Gabor-layered VGG16-ImageNet filters from the first convolutional layer without regularization.** Each of the rows in the figure corresponds to each of the 3 families of filters of the layer. Each column is one of the 8 $\alpha$-scaled rotations of the filter.