# DEPTHWISE-STFT BASED SEPARABLE CONVOLUTIONAL NEURAL NETWORKS

*Sudhakar Kumawat and Shanmuganathan Raman*

Computer Science and Engineering, Indian Institute of Technology Gandhinagar, India
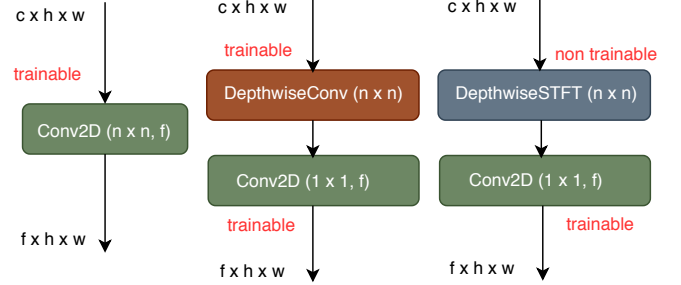{sudhakar.kumawat, shanmuga}@iitgn.ac.in

## ABSTRACT

In this paper, we propose a new convolutional layer called Depthwise-STFT Separable layer that can serve as an alternative to the standard depthwise separable convolutional layer. The construction of the proposed layer is inspired by the fact that the Fourier coefficients can accurately represent important features such as edges in an image. It utilizes the Fourier coefficients computed (channelwise) in the 2D local neighborhood (e.g., $3 \times 3$) of each position of the input map to obtain the feature maps. The Fourier coefficients are computed using 2D Short Term Fourier Transform (STFT) at multiple fixed low frequency points in the 2D local neighborhood at each position. These feature maps at different frequency points are then linearly combined using trainable pointwise $(1 \times 1)$ convolutions. We show that the proposed layer outperforms the standard depthwise separable layer based models on the CIFAR-10 and CIFAR-100 image classification datasets with reduced space-time complexity.

*Index Terms*— Convolutional neural networks, Short Term Fourier Transform, Separable convolutions

## 1. INTRODUCTION

Over the past few years, with the availability of large-scale datasets and computational power, deep learning has achieved impressive results on a wide range of applications in the field of computer vision. In general, the trend is to achieve higher performance by developing deep and complex models using large computational resources [1, 2, 3, 4, 5]. However, this progress is not necessarily making the networks more efficient with respect to memory and speed. In many real world and resource constraint applications such as robotics, satellites, and self-driving cars, the recognition tasks need to be carried out in a fast and computationally efficient manner. Therefore, there is a need to develop space-time efficient models for such applications.

In a standard convolutional layer, the convolution filters learn the spatial and channel correlations simultaneously. The depthwise separable convolutions factorize the above process

**Fig. 1**: **Various convolutional layer architectures.** Standard convolution (left), Depthwise separable convolution (center), Depthwise-STFT based separable convolution (right).

into two layers. In the first layer, a standard depthwise (channelwise) convolution is applied in order to learn the spatial correlations. In the second layer, pointwise convolutions ($1 \times 1$ convolutions) learn channel correlations by combining the outputs of the first layer. Fig. 1, compares various architectures. An important advantage of this factorization is that it significantly reduces the computation and model size. For example, for a filter of size $n = 3$, the depthwise separable convolution uses 8 to 9 times lesser parameters compared to the standard convolutional layer.

In this paper, we propose a new 2D convolutional layer named Depthwise-STFT separable layer. Similar to the standard depthwise separable convolution layer, our Depthwise-STFT separable layer has two sub-layers. The first layer named Depthwise-STFT captures the spatial correlations. For each channel in the input feature map, it computes the Fourier coefficients (at low frequency points) in a 2D local neighborhood ($n \times n$) at each position of the channel to obtain new feature maps. The Fourier coefficients are computed using 2D Short Term Fourier Transform (STFT) at multiple fixed low frequency points in the 2D local neighborhood at each position of the channel. The second layer named pointwise convolution uses $1 \times 1$ convolutions to learn channel correlations by combining feature maps obtained from the Depthwise-STFT layer. Note that unlike the case of standard depthwise separable layer, here only the second layer (pointwise convolutions) is trainable. Thus, the Depthwise-STFT separable layer has a lower space-time complexity when com-

pared to the depthwise separable convolutions. Furthermore, we show experimentally that the proposed layer achieves better performance compared to the many state-of-the-art depthwise separable based models such as MobileNet [6, 7] and ShuffleNet [8, 9].

## 2. RELATED WORKS

Recently, there has been a growing interest into developing space-time efficient neural networks for real time and resource restricted applications [10, 6, 7, 9, 8, 11, 12].

**Depthwise Separable Convolutions.** As discussed in Section 1, depthwise separable convolutions significantly reduce the space-time complexity of convolutional neural networks (CNNs) when compared to the standard convolutions by partitioning the steps of learning spatial and channel correlations. Recently, many depthwise separable convolutions based networks have been introduced such as MobileNet [6, 7], ShuffleNet [8, 9], and Xception [13]. Note that with the reduced complexity in the network architectures, these networks achieve trade-off between accuracy and space-time complexity.

**2D STFT based Convolutional Layers.** Recently, in [14], the authors introduced a 2D STFT based 2D convolutional layer named ReLPU for fundus image segmentation. The ReLPU layer when used inplace of the first convolutional layer (following the input layer) of the U-Net architecture [15] improved the performance of the baseline U-Net architecture. However, the ReLPU layer could not be used to replace all the convolutional layers in the network due to the extreme bottleneck used in it which reduced its learning capabilities. This work aims to solve this issue by introducing 2D STFT in depthwise separable convolutions.

## 3. METHOD

We will denote the feature map output by a layer in a 2D CNN network with the tensor $f(\mathbf{x}) \in \mathbb{R}^{c \times h \times w}$ where $h$, $w$, and $c$ are the height, width, and number of channels of the feature map, respectively. Fig. 1 presents the high-level architecture of our Depthwise-STFT based separable layer. Note that similar to the standard depthwise separable convolution layer, the Depthwise-STFT based separable layer has two sub-layers. In the first layer (named Depthwise-STFT), for each channel in the input feature map, we compute the Fourier coefficients in a 2D local neighborhood at each position of the channel to obtain the new feature maps. The Fourier coefficients are computed using 2D Short Term Fourier Transform (STFT) at multiple fixed low frequency points in the 2D local neighborhood at each position of the channel. The second layer (named pointwise convolutions) uses $1 \times 1$ convolutions to learn linear combinations of the feature maps obtained from the Depthwise-STFT layer. The detailed description of each

layer is as follows.

**Depthwise-STFT.** This layer takes a feature map $f(\mathbf{x}) \in \mathbb{R}^{c \times h \times w}$ as input from the previous layer. For simplicity, let us take $c = 1$. Hence, we can drop the channel dimension and rewrite the size of $f(\mathbf{x})$ as $h \times w$. Here, $\mathbf{x} \in \mathbb{Z}^2$ are the 2D coordinates of the elements in $f(\mathbf{x})$.

Next, each $\mathbf{x}$ in $f(\mathbf{x})$ has a $n \times n$ 2D neighborhood (denoted by $\mathcal{N}_{\mathbf{x}}$) which can be defined as shown in Equation 1.

$$\mathcal{N}_{\mathbf{x}} = \{\mathbf{y} \in \mathbb{Z}^2 \, ; \| \, (\mathbf{x} - \mathbf{y}) \, \|_{\infty} \leq r \, ; n = 2r + 1 ; r \in \mathbb{Z}_+ \} \quad (1)$$

Now, for all positions $\mathbf{x} = \{\mathbf{x}_1, \mathbf{x}_2, \ldots, \mathbf{x}_{h \times w}\}$ of the feature map $f(\mathbf{x})$, we use local 2D neighborhoods, $f(\mathbf{x} - \mathbf{y}), \forall \mathbf{y} \in \mathcal{N}_{\mathbf{x}}$ to derive the local frequency domain representation. For this, we use Short Term Fourier Transform (STFT) which is defined in Equation 2.

$$F(\mathbf{v}, \mathbf{x}) = \sum_{\mathbf{y}_i \in \mathcal{N}_{\mathbf{x}}} f(\mathbf{x} - \mathbf{y}_i) \exp^{-j2\pi \mathbf{v}^T \mathbf{y}_i} \quad (2)$$

Here $i = 1, \ldots, n^2$, $\mathbf{v} \in \mathbb{R}^2$ is a 2D frequency variable and $j = \sqrt{-1}$. Note that, due to the separability of the basis functions, 2D STFT can be efficiently computed using simple 1D convolutions for the rows and the columns, successively. Using vector notation, we can rewrite Equation 2 as shown in Equation 3.

$$F(\mathbf{v}, \mathbf{x}) = \mathbf{w}_{\mathbf{v}}^T \mathbf{f}_{\mathbf{x}} \quad (3)$$

Here, $\mathbf{w}_{\mathbf{v}}$ is a complex valued basis function (at frequency variable $\mathbf{v}$) of a linear transformation and is defined as shown in Equation 4.

$$\mathbf{w}_{\mathbf{v}}^T = [\exp^{-j2\pi \mathbf{v}^T \mathbf{y}_1}, \exp^{-j2\pi \mathbf{v}^T \mathbf{y}_2}, \ldots, \exp^{-j2\pi \mathbf{v}^T \mathbf{y}_{n^2}}] \quad (4)$$

$\mathbf{f}_{\mathbf{x}}$ is a vector containing all the elements from the neighborhood $\mathcal{N}_{\mathbf{x}}$ and is defined as shown in Equation 5.
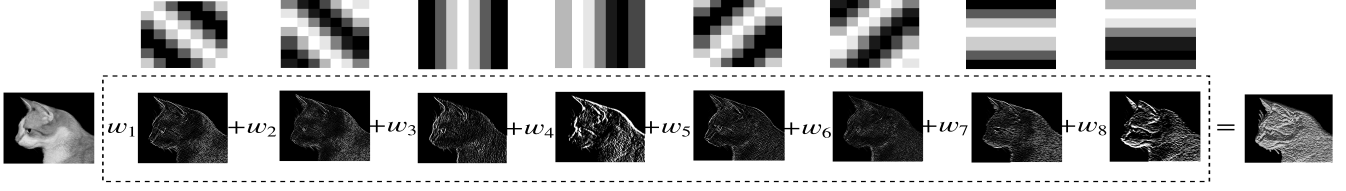
$$\mathbf{f}_{\mathbf{x}} = [f(\mathbf{x} - \mathbf{y}_1), f(\mathbf{x} - \mathbf{y}_2), \ldots, f(\mathbf{x} - \mathbf{y}_{n^2})]^T \quad (5)$$

In this work, we use four lowest non-zero frequency variables $\mathbf{v}_1 = [a, 0]^T, \mathbf{v}_2 = [0, a]^T, \mathbf{v}_3 = [a, a]^T$, and $\mathbf{v}_4 = [a, -a]^T$, where $a = 1/n$. Thus, from Equation 3, we can define the local frequency domain representation for the above four frequency variables as shown in Equation 6.

$$\mathbf{F}_{\mathbf{x}} = [F(\mathbf{v}_1, \mathbf{x}), F(\mathbf{v}_2, \mathbf{x}), F(\mathbf{v}_3, \mathbf{x}), F(\mathbf{v}_4, \mathbf{x})]^T \quad (6)$$

At each position $\mathbf{x}$, after separating the real and the imaginary parts of each component, we obtain a vector as shown in Equation 7.

$$\mathbf{F}_{\mathbf{x}} = [\Re\{F(\mathbf{v}_1, \mathbf{x})\}, \Im\{F(\mathbf{v}_1, \mathbf{x})\}, \ldots$$
$$\ldots, \Re\{F(\mathbf{v}_4, \mathbf{x})\}, \Im\{F(\mathbf{v}_4, \mathbf{x})\}]^T \quad (7)$$

**Fig. 2**: Visualization of the Depthwise-STFT Separable layer for $c = 1$ and $n = 7$. Here, weights $w_1, w_2, \ldots, w_8$ are learned by the network during training.

Here, $\Re\{\cdot\}$ and $\Im\{\cdot\}$ return the real and imaginary parts of a complex number, respectively. The corresponding $8 \times n^2$ transformation matrix can be written as shown in Equation 8.

$$\mathbf{W} = [\Re\{\mathbf{w}_{\mathbf{v}_1}\}, \Im\{\mathbf{w}_{\mathbf{v}_1}\}, \ldots, \Re\{\mathbf{w}_{\mathbf{v}_4}\}, \Im\{\mathbf{w}_{\mathbf{v}_4}\}]^T \quad (8)$$

Hence, from Equations 3 and 8, the vector form of STFT for all the four frequency points $\mathbf{v}_1, \mathbf{v}_2, \mathbf{v}_3$, and $\mathbf{v}_4$ can be written as shown in Equation 9.

$$\mathbf{F_x} = \mathbf{W}\mathbf{f_x} \quad (9)$$

Since $\mathbf{F_x}$ is computed for all positions $\mathbf{x}$ of the input $f(\mathbf{x})$, it results in an output feature map with size $8 \times h \times w$ corresponding to the four frequency variables. Remember that we took $c = 1$. Thus, for one channel, the Depthwise-STFT layer outputs a feature map of size $8 \times h \times w$ corresponding to the four frequency variables. Therefore, for $c$ channels, the Depthwise-STFT will output a feature map of size $8 \cdot c \times h \times w$.
**Pointwise convolutions.** This layer is the standard trainable $1 \times 1$ convolutional layer containing $f$ filters, each one of them has a depth equal to $8 \cdot c$ which takes as input a tensor of size $8 \cdot c \times h \times w$ and outputs a tensor of size $f \times h \times w$. Note that it is this layer that gets learned during the training phase of the CNN.

In Fig. 2, we present the visualization of the Depthwise-STFT based separable layer for $c = 1$ and $n = 7$. First the Fourier coefficients of the input feature map is extracted in a local $7 \times 7$ neighborhood of each position at four frequency points $[1/7, 0]^T, [0, 1/7]^T, [1/7, 1/7]^T$, and $[1/7, -1/7]$ to output the feature maps of size $8 \cdot c \times h \times w$ (after separating real and imaginary parts). Then, the output feature maps are linearly combined to output the final feature map.
**Parameter analysis.** Consider a standard 2D convolutional layer with $c$ input and $f$ output channels. Assume that spatial padding is done such that the spatial dimensions of the channels remain same. Let $n \times n$ be the size of the filters. In Table 1, we compare the number of trainable parameters in various convolutional layers. The number of trainable parameters in Depthwise-STFT separable layer is independent of the filter size $n$.
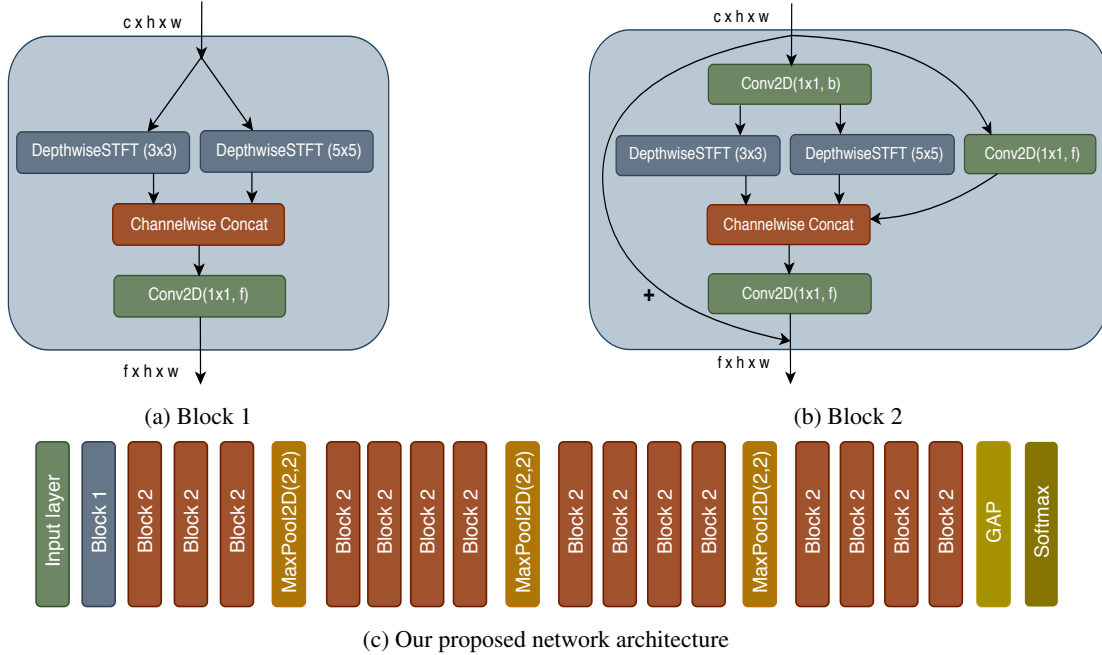
## 4. EXPERIMENTS

**Datasets.** We evaluate the Depthwise-STFT separable layer on the two popular CIFAR datasets− CIFAR-10 and CIFAR-

| Layer | # parameters |
|---|---|
| Standard Convolution | $c \cdot n^2 \cdot f$ |
| Depthwise Separable | $n^2 \cdot c + c \cdot f$ |
| Depthwise-STFT Separable | $8 \cdot c \cdot f$ |

**Table 1**: Comparison of the number of trainable parameters in various convolutional layers.

100 [16]. The CIFAR-10 and the CIFAR-100 datasets consist of 10 and 100 classes, respectively. Each dataset consists of natural RGB images of resolution $32 \times 32$ pixels with 50,000 images in the training set and 10,000 images in the testing set. We use the standard data augmentation scheme − horizontal flip/shifting/rotation ($\pm 20$) that is widely used for these two datasets. For preprocessing, we normalize the data using the channel means and standard deviations.
**Network Architecture.** We adopt a simple Inception-ResNet style bottleneck architecture [5]. Fig. 3 presents the building blocks of our network, *Block 1* (Fig. 3a) and *Block 2* (Fig. 3b). *Block 1* takes as input a feature map with $c$ channels and applies a bottleneck $1 \times 1$ convolution (trainable) to output a feature map of size $b$ such that $b < c$. This is followed by inception style non-trainable Depthwise-STFT layers with filter size 3 and 5. The intuition behind this design is to have a maximum number of possible pathways for information to flow through the network and let the network selectively choose the best frequency points/neighborhood sizes for computing local Fourier transform and to selectively give more weight to them during training. The outputs from the Depthwise-STFT layers are concatenated channel-wise and finally passed through an expansion $1 \times 1$ convolution (trainable) to output a feature map of size $f$ such that $f > b$. We use bottleneck architectures as they have been proved to be generalize better [7]. The architecture of *Block 2* is just a slightly augmented version of *Block 1* with skip connections (Fig. 3b). Each block is followed by a batch normalization layer which is followed by a LeakyReLU (with $\alpha = 0.3$) activation function [17]. For down-sampling, we use max pooling with pool size 2 and stride 2. Our final network as shown in Fig. 3c consists of 16 non-downsampling blocks (*Block 1* and *Block 2*), followed by a global average pooling layer connected to a final classification layer with softmax activation. We implemented the proposed network using Keras deep learning library [18] and

(a) Block 1

(b) Block 2

(c) Our proposed network architecture

**Fig. 3**: Various building blocks of the proposed network architecture and our proposed network

| Network | # params | # feat. | CIFAR-10 | CIFAR-100 |
|---|---|---|---|---|
| ShuffleNet [9] | 1.05M | 800 | 90.80 | 70.06 |
| ShuffleNet-V2 [8] | 1.35M | 1024 | 91.42 | 69.51 |
| MobileNet [6] | 3.31M | 1024 | 91.87 | 65.98 |
| MobileNet-V2 [7] | 2.36M | 1280 | 93.14 | 68.08 |
| ReLPU based [14] | 3.08M | 256 | 92.20 | 70.20 |
| Ours (b=8, f=128) | 0.90M | 128 | 93.16 | 70.19 |
| Ours (b=16, f=128) | 1.30M | 128 | 93.59 | 70.66 |
| Ours (b=32, f=128) | 2.21M | 128 | 93.72 | 71.08 |
| Ours (b=64, f=128) | 3.69M | 128 | 94.07 | 71.42 |
| Ours (b=64, f=256) | 8.21M | 256 | 94.25 | 73.01 |
| Ours (b=64, f=384) | 14.06M | 384 | **94.51** | **74.39** |

**Table 2**: Performance results of the Depthwise-STFT separable layer based architectures compared to the standard depthwise separable layer based architectures.

performed all the experiments on a system with Intel i7-8700 processor, 32 Gb RAM, and a single Nvidia Titan Xp GPU.

**Training.** For training out networks, we use Adam optimizer [19], categorical cross-entropy as loss function, and batch size of 64. All the trainable weights are initialized with orthogonal initializer [20]. We train our networks with a learning rate of 0.01 for 300 epochs. After that, we increase the batch size to 128 and further train the networks for 100 epochs. Note that the above training method is inspired from [21] which proposes that it is preferable to increase the batch size instead of decreasing the learning rate during training.

**Results and Analysis.** Table 2 reports the classification results of different resource efficient architectures on the CIFAR-10 and CIFAR-100 datasets. For fair comparison, we compare our networks with depthwise separable based architectures only such as MobileNet [6, 7] and ShuffleNet [8, 9]. As discussed in Section 2, we also compare with the ReLPU layer [14] based network which is form by replacing all the layers of the network of Fig. 3c with ReLPU layer. Our results show that the proposed layer outperforms the depthwise separable layer based and the ReLPU layer based architectures. Note that here our target is not to achieve state-of-the-art accuracy on the two datasets but to showcase the efficiency of our proposed layer when compared to the depthwise separable convolution based architectures. For analysis purpose, we ran two variants of the proposed network (Fig. 3c). In the first variant, the bottleneck parameter $b$ is increased while the expansion parameter $f$ is kept constant. In the second variant, we kept the bottleneck parameter $b$ constant while increasing the expansion parameter $f$. In both the settings, we observe improvement in the performance of the networks.

## 5. CONCLUSION

This paper proposes Depthwise-STFT separable layer, that can serve as an alternative to the standard depthwise separable layer. The proposed layer captures spatial correlations (channelwise) in the feature maps using STFT followed by pointwise convolutions to channel correlations. Our proposed layer outperforms the standard depthwise separable layer based models on the CIFAR-10 and CIFAR-100 datasets. Furthermore, it has lower space-time complexity when compared to standard depthwise separable layer.

# 6. REFERENCES

[1] Karen Simonyan and Andrew Zisserman, "Very deep convolutional networks for large-scale image recognition," *arXiv preprint arXiv:1409.1556*, 2014.

[2] Kaiming He, Xiangyu Zhang, Shaoqing Ren, and Jian Sun, "Deep residual learning for image recognition," in *CVPR*, 2016.

[3] Dongyoon Han, Jiwhan Kim, and Junmo Kim, "Deep pyramidal residual networks," in *CVPR*, 2017.

[4] Sergey Zagoruyko and Nikos Komodakis, "Wide residual networks," *BMVC*, 2016.

[5] Christian Szegedy, Sergey Ioffe, Vincent Vanhoucke, and Alexander A Alemi, "Inception-v4, inception-resnet and the impact of residual connections on learning," in *AAAI*, 2017.

[6] Andrew G Howard, Menglong Zhu, Bo Chen, Dmitry Kalenichenko, Weijun Wang, Tobias Weyand, Marco Andreetto, and Hartwig Adam, "Mobilenets: Efficient convolutional neural networks for mobile vision applications," *arXiv preprint arXiv:1704.04861*, 2017.

[7] Mark Sandler, Andrew Howard, Menglong Zhu, Andrey Zhmoginov, and Liang-Chieh Chen, "Mobilenetv2: Inverted residuals and linear bottlenecks," in *CVPR*, 2018.

[8] Ningning Ma, Xiangyu Zhang, Hai-Tao Zheng, and Jian Sun, "Shufflenet v2: Practical guidelines for efficient cnn architecture design," in *ECCV*, 2018.

[9] Xiangyu Zhang, Xinyu Zhou, Mengxiao Lin, and Jian Sun, "Shufflenet: An extremely efficient convolutional neural network for mobile devices," in *CVPR*, 2018.

[10] Sachin Mehta, Mohammad Rastegari, Anat Caspi, Linda Shapiro, and Hannaneh Hajishirzi, "Espnet: Efficient spatial pyramid of dilated convolutions for semantic segmentation," in *ECCV*, 2018.

[11] Sachin Mehta, Mohammad Rastegari, Linda Shapiro, and Hannaneh Hajishirzi, "Espnetv2: A light-weight, power efficient, and general purpose convolutional neural network," in *CVPR*, 2019.

[12] Sudhakar Kumawat and Shanmuganathan Raman, "Lp-3dcnn: Unveiling local phase in 3d convolutional neural networks," in *CVPR*, 2019.

[13] François Chollet, "Xception: Deep learning with depthwise separable convolutions," in *CVPR*, 2017.

[14] Sudhakar Kumawat and Shanmuganathan Raman, "Local phase u-net for fundus image segmentation," in *ICASSP*, 2019.

[15] Olaf Ronneberger, Philipp Fischer, and Thomas Brox, "U-net: Convolutional networks for biomedical image segmentation," in *MICCAI*, 2015.

[16] A Krizhevsky, "Learning multiple layers of features from tiny images," *Master's thesis, University of Tront*, 2009.

[17] Andrew L Maas, Awni Y Hannun, and Andrew Y Ng, "Rectifier nonlinearities improve neural network acoustic models," in *ICML Workshop on Deep Learning for Audio, Speech and Language Processing*, 2013.

[18] Antonio Gulli and Sujit Pal, *Deep Learning with Keras*, Packt Publishing Ltd, 2017.

[19] Diederik P Kingma and Jimmy Ba, "Adam: A method for stochastic optimization," *ICLR*, 2014.

[20] Kaiming He, Xiangyu Zhang, Shaoqing Ren, and Jian Sun, "Delving deep into rectifiers: Surpassing human-level performance on imagenet classification," in *ICCV*, 2015.

[21] Samuel L Smith, Pieter-Jan Kindermans, Chris Ying, and Quoc V Le, "Don't decay the learning rate, increase the batch size," in *ICLR*, 2017.