

# Multikernel activation functions: formulation and a case study<sup>\*</sup>

Simone Scardapane<sup>1</sup>[0000-0003-0881-8344], Elena Nieddu<sup>2</sup>, Donatella Firmani<sup>2</sup>[0000-0003-0358-3208], and Paolo Meriardo<sup>2</sup>[0000-0002-3852-8092]

<sup>1</sup> DIET Department, Sapienza University of Rome, Rome, Italy  
`{simone.scardapane}@uniroma1.it`

<sup>2</sup> Department of Engineering, Roma Tre University, Rome, Italy  
`firstname.lastname@uniroma3.com`

**Abstract.** The design of activation functions is a growing research area in the field of neural networks. In particular, instead of using fixed point-wise functions (e.g., the rectified linear unit), several authors have proposed ways of learning these functions directly from the data in a non-parametric fashion. In this paper we focus on the kernel activation function (KAF), a recently proposed framework wherein each function is modeled as a one-dimensional kernel model, whose weights are adapted through standard backpropagation-based optimization. One drawback of KAFs is the need to select a single kernel function and its eventual hyperparameters. To partially overcome this problem, we motivate an extension of the KAF model, in which multiple kernels are linearly combined at every neuron, inspired by the literature on multiple kernel learning. We provide an application of the resulting *multi-KAF* on a realistic use case, specifically handwritten Latin OCR, on a large dataset collected in the context of the ‘In Codice Ratio’ project. Results show that multi-KAFs can improve the accuracy of the convolutional networks previously developed for the task, with faster convergence, even with a smaller number of overall parameters.

**Keywords:** Activation function · multikernel · OCR · Latin.

## 1 Introduction

The recent successes in deep learning owe much to concurrent advancements in the design of activation functions, especially the introduction of the rectified linear unit (ReLU) [7], and its several variants [11,2,13]. Albeit the vast majority of applications consider fixed activation functions, a growing trend of research lately has focused on designing *flexible* ones, that are able to adapt their shape from the data through the use of additional trainable parameters. Common

---

<sup>\*</sup> The work of S. Scardapane was supported in part by Italian MIUR, “*Progetti di Ricerca di Rilevante Interesse Nazionale*”, GAUCHO project, under Grant 2015YYPXH4W\_004.

examples of this are parametric versions of ReLU [11], or the more recent beta-Swish function [15].

In the more extreme case, multiple authors have advocated for *non-parametric* formulations, in which the overall flexibility and number of parameters can be chosen freely by the user on the basis of one or more hyper-parameters. As a result, the trained functions can potentially approximate a much larger family of shapes. Different proposals, however, differ on the way in which each function is modeled, resulting in vastly different characteristics in terms of approximation, optimization, and simplicity of implementation. Examples of non-parametric activation functions are the maxout neuron [9,20], defined as the maximum over a fixed number of affine functions of its input, the Fourier activation function [3], defined as a linear combination of a predetermined trigonometric basis expansion, or the Hermitian-based expansion [18]. We refer to [16] for a fuller overview on the topic.

In this paper we focus on the recently proposed kernel activation function (KAF) [16], in which each (scalar) function is modeled as a one-dimensional kernel expansion, with the linear mixing coefficients adapted together with all the other parameters of the network during optimization. In [16] it was shown that KAFs can greatly simplify the design of neural networks, allowing to reach higher accuracies, sometimes with a smaller number of hidden layers. Linking neural networks with kernel methods also allows to leverage a large body of literature on the learning of kernel functions (e.g., kernel filters [14]), particularly with respect to their approximation capabilities. At the same time, compared to ReLUs, KAFs introduce a number of additional design choices, most notably the selection of which kernel function to use, and its eventual hyper-parameters (e.g., the bandwidth of the Gaussian kernel). Although in [16] and successive works we mostly focused on the Gaussian kernel, it is not guaranteed to be the optimal one in all applications.

**Contribution of the paper** To solve the kernel selection problem of KAFs, in this paper we propose an extension inspired to the theory of multiple kernel learning [8,1]. In the proposed *multi-KAF*, different kernels are linearly combined for every neuron through an additional set of mixing coefficients, adapted during training. In this way, the optimal kernel for each neuron (or a specific mixture of them) can be learned in a principled way during the optimization process. In addition, since in our KAF implementation the points where the kernels are evaluated are fixed, a large amount of computation can be shared between the different kernels, leading to a very small computational overhead overall, as we show in the following sections.

**Case study: In Codice Ratio** To show the usefulness of the proposed activation functions, we provide a realistic use case by applying them to the data from the ‘In Codice Ratio’ (ICR) project [4], whose aim is the automatic transcription of a large part of the Vatican Secret Archive.<sup>3</sup> A key component of

<sup>3</sup> <http://www.archiviosegreto.vaticano.va/>

the project is an OCR tool applied to characters from a Latin handwritten text (see Section 4 for additional details). In [5] we presented a convolutional neural network (CNN) for this task, which we applied to a dataset of 23 different Latin characters extracted from a sample selection of pages from the Vatican Register. In this paper we show that, using multi-KAFs, we can increase the accuracy of the CNN even while reducing the number of filters per layer.

**Organization of the paper** In Section 2 and 3 we describe standard activation functions for neural networks, KAFs [16], and the proposed multi-KAFs. The dataset from the ICR project is described in Section 4. We then perform a set of experiments in Section 5, before concluding in Section 6.

## 2 Preliminaries

### 2.1 Feedforward neural networks

Consider a generic feedforward NN layer, taking as input a vector  $\mathbf{x} \in \mathbb{R}^d$  and producing in output a vector  $\mathbf{y} \in \mathbb{R}^c$ :

$$\mathbf{y} = g(\mathbf{W}\mathbf{x} + \mathbf{b}) , \quad (1)$$

where  $\{\mathbf{W}, \mathbf{b}\}$  are adaptable weight matrices, and  $g(\cdot)$  is an element-wise activation function. Multiple layers can be stacked to obtain a complete NN. In the following we focus especially on the choice of  $g(\cdot)$ , but we note that everything extends immediately to more complex types of layer, including convolutional layers (wherein the matrix product is replaced by a convolutional operator), or recurrent layers [10].

Generally speaking, the activation functions  $g(\cdot)$  for the hidden (not last) layers are chosen as simple operations, such as the rectified linear unit (ReLU), originally introduced in [7]:

$$g(s) = \max\{0, s\} , \quad (2)$$

where we use the letter  $s$  to denote a generic scalar input to the function, i.e., a single *activation* value. An NN is a generic composition of  $L$  such layers, denoted as  $f(\mathbf{x})$ , that is trained with a dataset of  $N$  training samples  $\{\mathbf{x}^i, \mathbf{d}^i\}_{i=1}^N$ . In the experimental section, in particular, we deal with multi-class classification with  $C$  classes, where the desired output  $\mathbf{d}^i$  represents a one-hot encoding of the target class. We train it by minimizing a regularized cross-entropy cost:

$$\min \left\{ - \sum_{i=1}^N \sum_{c=1}^C d_c^i \log (f_c(\mathbf{x}^i)) + \lambda \cdot \|\mathbf{w}\|^2 \right\} , \quad (3)$$

where  $\mathbf{w}$  is a weight vector collecting all the adaptable weights of the network,  $\lambda$  is a positive scalar, and we use a subscript to denote the  $c$ -th element of a vector.

## 2.2 Kernel activation functions

Differently from (2), a KAF can be adapted from the data. In particular, each activation function is modeled in terms of  $D$  expansions of the activation  $s$  with a kernel function  $\kappa$ :

$$g(s) = \sum_{i=1}^D \alpha_i \kappa(s, d_i), \quad (4)$$

where the scalars  $d_i$  form the so-called dictionary, while the scalars  $\alpha_i$  are called the mixing coefficients. To make the training problem simpler, the dictionary is fixed beforehand (and not adapted) by sampling  $D$  values from the real line, uniformly around zero, and it is shared across the network, while a different set of mixing coefficients is adapted for every neuron. This makes implementation extremely efficient. The integer  $D$  is the key hyper-parameter of the model: a higher value of  $D$  increases the overall flexibility of each function, at the expense of adding additional mixing coefficients to be adapted.

Any kernel function from the literature can be used in (4), provided it respects the semi-definiteness property:

$$\sum_{i=1}^D \sum_{j=1}^D \alpha_i \alpha_j \kappa(d_i, d_j) \geq 0, \quad (5)$$

for any choice of the mixing coefficients and the dictionary. In practice, [16] and all subsequent papers only used the one-dimensional Gaussian kernel defined as:

$$\kappa(s, d_i) = \exp \left\{ -\gamma (s - d_i)^2 \right\}, \quad (6)$$

where  $\gamma > 0$  is a parameter of the kernel. The value of  $\gamma$  influences the ‘locality’ of each  $\alpha_i$  with respect to the dictionary. In [16] we proposed the following rule-of-thumb, found empirically:

$$\gamma = \frac{1}{6\Delta^2}, \quad (7)$$

where  $\Delta$  is the distance between any two dictionary elements. However, note that neither the Gaussian kernel in (6) nor the rule-of-thumb in (7) are optimal in general. For example, simple smooth shapes (like slowly varying polynomials) could be more easily modeled via different types of kernels or much larger values of  $\gamma$  with a possibly smaller  $D$ . These are common problems also in the kernel literature. Leveraging it, in the next section we propose an extension of KAF to mitigate both problems.

## 3 Proposed multiple kernel activation functions

In order to mitigate the problems mentioned in the previous section, assume to have available a set of  $M$  *candidate* kernel functions  $\kappa_1, \dots, \kappa_M$ . These can

be entirely different functions or the same kernel with different choices of its parameters. There has been a vast research on how to successfully combine different kernels to obtain a new one, going under the name of multiple kernel learning (MKL) [1]. For the purpose of this paper we adopt a simple approach, in order to evaluate its feasibility. In particular, we build each KAF with a new kernel given by a linearly weighted sum of the constituents (base) kernels:

$$g(s) = \sum_{i=1}^D \alpha_i \left[ \sum_{m=1}^M \mu_m \kappa_m(s, d_i) \right] = \sum_{i=1}^D \alpha_i \tilde{\kappa}(s, d_i), \quad (8)$$

where  $\{\mu_m\}_{m=1}^M$  are an additional set of mixing coefficients. From the properties of reproducing kernel Hilbert spaces, it is straightforward to show that  $\tilde{\kappa}$  is a valid kernel if its constituents are also valid kernel functions. We call the resulting activation functions *multi-KAFs*. Note that such an approach only introduces  $M - 1$  additional parameters for each neuron, where  $M$  is generally small. In addition, since in our implementation the dictionary is fixed, all kernels are evaluated on the same points, which can greatly simplify the implementation and allows to share a large part of the computation.

More in detail, for our experiments we consider an implementation with  $M = 3$ , where  $\kappa_1$  is the Gaussian kernel in (6),  $\kappa_2$  is chosen as the (isotropic) rational quadratic [6]:

$$\kappa_2(s, d_i) = 1 + \frac{(s - d_i)^2}{(s - d_i)^2 + c} \quad (9)$$

with  $c$  being a parameter, and  $\kappa_3$  is chosen as the polynomial kernel of order 2:

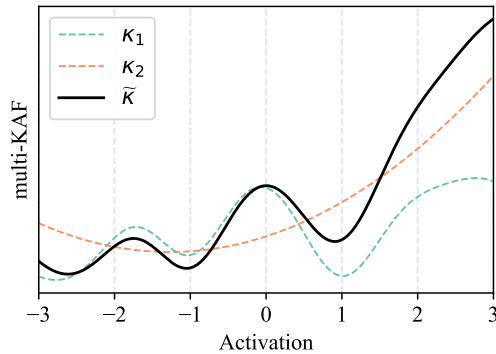
$$\kappa_3(s, d_i) = (1 + sd_i)^2. \quad (10)$$

The rational quadratic is similar to the Gaussian, but it is sometimes preferred in practical applications [6], while the polynomial kernel allows to introduce a smoothly varying global trend to the function. We show a simple example of the mix of  $\kappa_1$  and  $\kappa_3$  in Fig. 1.

In order to simplify optimization, especially for deeper architectures, we apply the kernel ridge regression initialization procedure described in [16] to initialize all multi-KAFs to a known activation function, i.e., the exponential linear unit (ELU) [2]. For this purpose, denote by  $\mathbf{t}$  the vector of ELU values computed on our dictionary points. We initialize all  $\mu_j$  to  $\frac{1}{3}$ , and initialize the vector of mixing coefficients  $\boldsymbol{\alpha}$  as:

$$\boldsymbol{\alpha} = \left( \tilde{\mathbf{K}} + \varepsilon \mathbf{I} \right)^{-1} \mathbf{t}, \quad (11)$$

where  $\tilde{\mathbf{K}} \in \mathbb{R}^{D \times D}$  is the kernel matrix computed between  $\mathbf{t}$  and  $\mathbf{d}$  using  $\tilde{\kappa}$ , and  $\varepsilon = 10^{-4}$ . As a final remark, note that in (8) we considered an unrestricted linear combination of the constituting kernels. We can easily obtain more restricted formulations (which are sometimes found in the MKL literature [8]) by applying some nonlinear transformation to the mixing coefficients  $\mu_m$ , e.g., a softmax function to obtain convex combinations. We leave such comparisons to a future work.

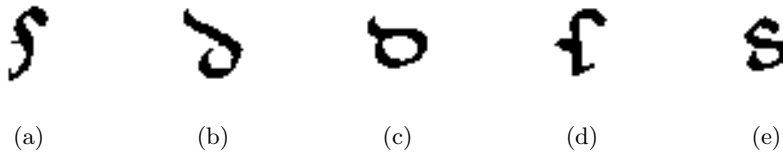


**Fig. 1.** An example of multi-KAF obtained from the mixture of a Gaussian kernel (green dashed line) and a polynomial kernel (red dashed line). In the example, we have  $D = 15$  elements in the dictionary equispaced in  $[-3.0, 3.0]$ , mixing coefficients  $\alpha_i$  sampled from the normal distribution, and  $\mu_1 = 0.5$ ,  $\mu_2 = 0.01$ .

#### 4 Case Study: In Codice Ratio

As stated in the introduction, we apply the proposed multi-KAF on a realistic case study taken from the ICR project. Apart from the details described below, we refer the reader to [5,4] for a fuller description of the project.

The overall goal of ICR is the transcription of a large portion of the Vatican Secret Archives, one of the largest existing historical libraries. In the first phase of the project, we collected and manually annotated a set of 23000 images representing 23 different types of characters (one of which is a special ‘non-character’ class). All characters were extracted from a sample of 30 (handwritten) pages of private correspondence of Pope Honorii III from the XIII century. Each character was then annotated using a crowdsourcing platform and 120 volunteer students.<sup>4</sup> A few examples taken from the dataset are shown in Fig. 2.



**Fig. 2.** Examples taken from the Latin OCR dataset. (a) and (d) are examples of the character ‘s’; (b) and (c) are examples of character ‘d’; (e) is an example of a different version for the character ‘s’ (considered as a separate class in the dataset).

<sup>4</sup> The dataset is available on the web at <http://www.dia.uniroma3.it/db/icr/>.

In [5] we described the design and evaluation of a CNN for tackling the problem of automatically assigning a class to each character, represented as a  $56 \times 56$  black-and-white image. The final CNN had the following architecture: (i) a convolutive block with 42 filters of size  $5 \times 5$ ; (ii) max-pooling with size  $2 \times 2$ ; (iii) two additional series of convolutive blocks and max-pooling, this time with 28 filters per layer; (iv) a fully connected layer with 100 neurons, followed by (v) an output layer of 23 neurons. In the original implementation, all linear operations were preceded by dropout [19] (with probability 50%), and followed by ReLU nonlinearities. We will use this dataset and architecture as baseline for our experiments. Note that this design was heavily fine-tuned, and it was not possible to increase the testing accuracy by simply adding more neurons / layers to the resulting CNN model.

## 5 Experimental results

### 5.1 Experimental setup

For our comparisons, we use the same CNN architecture described in Section 4, but we replace the ReLU functions by either KAFs or the proposed multi-KAFs, with the hyper-parameters described in the previous sections. In order to make the networks comparable in terms of parameters, for KAF and multi-KAF we decrease the number of filters and neurons in the linear layers by 10%. To stabilize training, we also replace dropout with a batch normalization step [12] before applying KAF-based nonlinearities.

We train the networks following a similar procedure as [5]. We use the Adam optimization algorithm on random mini-batches of 32 elements, with a small regularization factor  $\lambda = 0.001$ . After every 10 iterations of the optimization algorithm we evaluate the accuracy on a randomly held-out set of 2500 samples, taken from the original training set. Training is stopped whenever the validation accuracy has not improved for at least 250 iterations. The networks are then evaluated on a second independent held-out set of 2300 examples. All networks are implemented in PyTorch, and experiments are run on a CUDA backend using the Google Colaboratory platform.

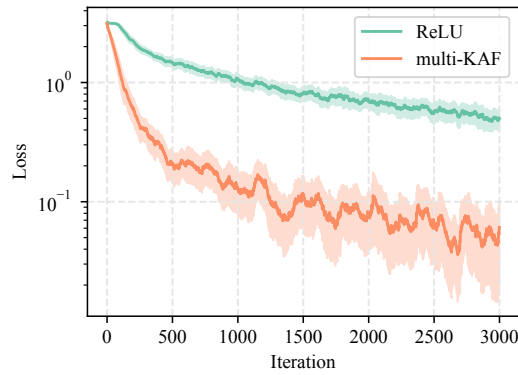
### 5.2 Experimental results

The results of the experiments, averaged over 5 different repetitions, are shown in Table 1, together with the number of trainable parameters of the different architectures. It can be seen that, while KAF fails to provide a meaningful improvement in this case, the multi-KAF architecture obtains a significant gain in testing accuracy, stably throughout the repetitions. Most notably, this gain is obtained with a strong decrease in the number of trainable parameters, which is around  $16 \cdot 10^4$  for the KAF-based architectures, compared to  $\approx 19 \cdot 10^4$  parameters for the baseline one.

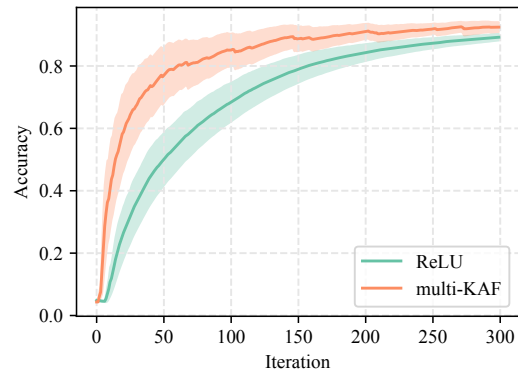
This gain in accuracy is not only obtained with a smaller number of overall parameters, but also with a much faster rate of convergence. To see this, we

**Table 1.** Results of comparing the different activation functions on the OCR dataset described in Section 4. See the text for details on the experimental procedure.

Activation function	Testing accuracy [%]	Trainable parameters
ReLU	94.99 ( $\pm 0.1$ )	191871
KAF	94.00 ( $\pm 0.3$ )	158840
Proposed multi-KAF	96.31 ( $\pm 0.2$ )	159552



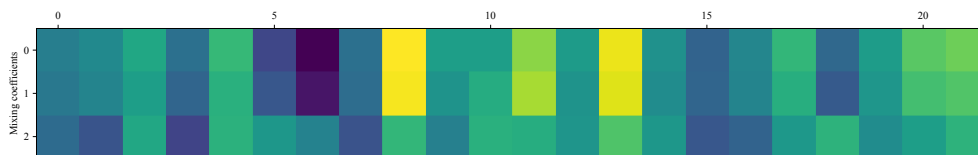
(a) Training loss



(b) Validation accuracy

**Fig. 3.** Loss and validation accuracy evolution for the baseline network (using ReLUs) and the proposed multi-KAF. Standard deviation for all curves is shown with a lighter color.





**Fig. 4.** Mixing coefficients of the multiple kernel after training, for the first convolutive layer and 25 randomly sampled filters. Light colors indicate stronger mixing coefficients.

plot in Fig. 3 the evolution of the loss function in (3) and the evolution of the accuracy on the validation portion of the dataset.

Finally, we also show the final mixing coefficients  $\mu_m$  for 25 randomly sampled neurons of the first convolutive layer of the multi-KAF network, after training, in Fig. 4. Interestingly, different neurons require vastly different combinations of base kernels, reflected by light/dark colors in Fig. 4.

## 6 Conclusions

In this paper we investigated a new non-parametric activation function for neural networks, which extends the recently proposed kernel activation function, by incorporating ideas from the field of multiple kernel learning to simplify the choice of the kernel function and further increase the expressiveness.

We evaluated the resulting multi-KAF on a benchmark dataset of Latin handwritten characters recognition, in the context of an ongoing real-world project. While in a sense these are only preliminary results, they point to the greater flexibility of such activation functions, coming with a faster rate of convergence and an overall smaller number of trainable parameters for the full architecture. We are currently in the process of collecting a larger dataset, considering a bigger amount of possible classes for the characters, in order to further evaluate the proposed architecture.

Additional research directions will consider the application of multi-KAFs on different types of benchmarks, going beyond standard CNNs, particularly with respect to recurrent models, complex-valued kernels [17], and generative adversarial networks. Furthermore, we plan to test more extensively additional types of kernels, such as the periodic ones [6], in order to evaluate the scalability of multi-KAFs in the presence of a larger number of constituent kernels. Generalization bounds for the architecture are also a promising research direction.

## References

1. Aioli, F., Donini, M.: Easymkl: a scalable multiple kernel learning algorithm. *Neurocomputing* **169**, 215–224 (2015)
2. Clevert, D.A., Unterthiner, T., Hochreiter, S.: Fast and accurate deep network learning by exponential linear units (ELUs). In: *Proc. 2016 International Conference on Learning Representations (ICLR)* (2016)

3. Eisenach, C., Wang, Z., Liu, H.: Nonparametrically learning activation functions in deep neural nets. In: 5th International Conference for Learning Representations (Workshop Track) (2017)
4. Firmani, D., Maiorino, M., Merialdo, P., Nieddu, E.: Towards knowledge discovery from the vatican secret archives. in codice ratio-episode 1: Machine transcription of the manuscripts. arXiv preprint arXiv:1803.03200 (2018)
5. Firmani, D., Merialdo, P., Nieddu, E., Scardapane, S.: In codice ratio: Ocr of hand-written latin documents using deep convolutional networks. In: 11th International Workshop on Artificial Intelligence for Cultural Heritage (AI\*CH 2017). pp. 9–16 (2017)
6. Genton, M.G.: Classes of kernels for machine learning: a statistics perspective. *Journal of machine learning research* **2**(Dec), 299–312 (2001)
7. Glorot, X., Bordes, A., Bengio, Y.: Deep sparse rectifier neural networks. In: Proc. 14th International Conference on Artificial Intelligence and Statistics (AISTATS). p. 275 (2011)
8. Gönen, M., Alpaydm, E.: Multiple kernel learning algorithms. *Journal of machine learning research* **12**(Jul), 2211–2268 (2011)
9. Goodfellow, I.J., Warde-Farley, D., Mirza, M., Courville, A., Bengio, Y.: Maxout networks. In: Proc. 30th International Conference on Machine Learning (ICML) (2013)
10. Goodfellow, I., Bengio, Y., Courville, A.: *Deep learning*. MIT Press (2016)
11. He, K., Zhang, X., Ren, S., Sun, J.: Delving deep into rectifiers: Surpassing human-level performance on imagenet classification. In: Proc. IEEE International Conference on Computer Vision (ICCV). pp. 1026–1034 (2015)
12. Ioffe, S., Szegedy, C.: Batch normalization: Accelerating deep network training by reducing internal covariate shift. In: Proc. 32nd International Conference on Machine Learning (ICML). pp. 448–456 (2015)
13. Jin, X., Xu, C., Feng, J., Wei, Y., Xiong, J., Yan, S.: Deep learning with S-shaped rectified linear activation units. In: Proc. Thirtieth AAAI Conference on Artificial Intelligence (2016)
14. Liu, W., Principe, J.C., Haykin, S.: *Kernel adaptive filtering: a comprehensive introduction*. John Wiley & Sons (2011)
15. Ramachandran, P., Zoph, B., Le, Q.V.: Searching for activation functions. arXiv preprint arXiv:1710.05941 (2017)
16. Scardapane, S., Van Vaerenbergh, S., Totaro, S., Uncini, A.: Kafnets: kernel-based non-parametric activation functions for neural networks. *Neural Networks* (2018), in press.
17. Scardapane, S., Van Vaerenbergh, S., Hussain, A., Uncini, A.: Complex-valued neural networks with non-parametric activation functions. *IEEE Transactions on Emerging Topics in Computational Intelligence* (2018), in press.
18. Siniscalchi, S.M., Salerno, V.M.: Adaptation to new microphones using artificial neural networks with trainable activation functions. *IEEE Transactions on Neural Networks and Learning Systems* **28**(8), 1959–1965 (2017)
19. Srivastava, N., Hinton, G.E., Krizhevsky, A., Sutskever, I., Salakhutdinov, R.: Dropout: a simple way to prevent neural networks from overfitting. *Journal of Machine Learning Research* **15**(1), 1929–1958 (2014)
20. Zhang, X., Trmal, J., Povey, D., Khudanpur, S.: Improving deep neural network acoustic models using generalized maxout networks. In: 2014 IEEE International Conference on Acoustics, Speech and Signal Processing (ICASSP). pp. 215–219. IEEE (2014)