

UTILIZING EXCESS RESOURCES IN TRAINING NEURAL NETWORKS

Amit Henig and Raja Giryes

Electrical Engineering, Tel Aviv University

ABSTRACT

In this work, we suggest **Kernel Filtering Linear Overparameterization (KFLO)**, where a linear cascade of filtering layers is used during training to improve network performance in test time. We implement this cascade in a kernel filtering fashion, which prevents the trained architecture from becoming unnecessarily deeper. This also allows using our approach with almost any network architecture and let combining the filtering layers into a single layer in test time. Thus, our approach does not add computational complexity during inference. We demonstrate the advantage of KFLO on various network models and datasets in supervised learning.

Index Terms— linear overparameterization, structural reparameterization, kernel filtering / composition

1. INTRODUCTION

Deep neural networks have shown remarkable capabilities in computer vision and natural language processing. Research on improving the deployed networks' performance has been done in many fields and different problem frameworks. In this work, we focus on the common case where the deployed model architecture is chosen in advance of training, and resources are abundant during training relative to inference time. These resources can be used to improve performance. Specifically, we focus on using a larger network during training that generates a smaller one for inference.

One may treat this problem setting as a network compression or a knowledge transfer problem. A known strategy for this task is Knowledge Distillation (KD). There, knowledge is transferred from a pre-trained teacher network to a desired student network using signals from the teacher. Interestingly, the student can match or even outperform the teacher's performance while being notably less complex. KD was first introduced by [1] for model compression and generalized by Hinton et. al. [2]. Various KD extensions exist such as KD in generations [3], ensemble of students teaching each other [4], the teacher assistant [5], and many more.

A different approach suggests using structural reparameterization, term denoted by [6] for transforming one architecture's parameters to produce the parameters used by another architecture. Specifically, in the context of this paper, a collapsible overparameterization. Common networks are typi-

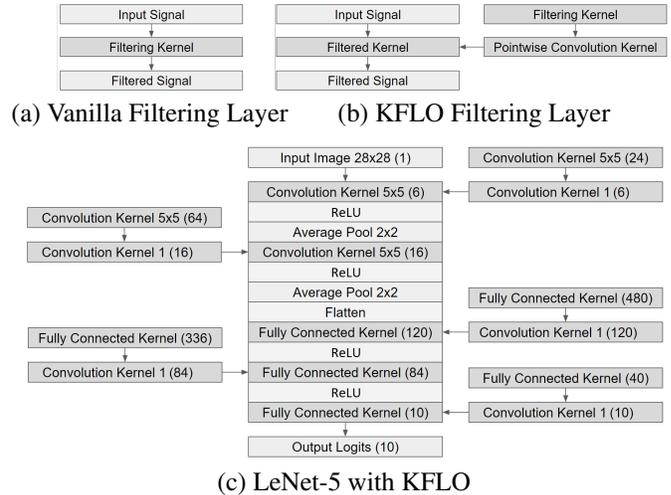


Fig. 1. A brief preview of our method. Layers containing filtering weights are marked with a darker shade. (a) A common vanilla filtering layer. (b) Filtering layer with KFLO: The kernel that filters the input signal is the result of a pointwise convolution applied on a different filtering kernel. The weights outside the input signal's path can be discarded after training, deploying only the Filtered Kernel. (c) LeNet-5 implemented with KFLO, for 28x28 grayscale images of 10 possible classes. Here too, the weights outside the signal's main path can be discarded after training, resulting in the original LeNet-5 architecture to be used at inference time.

cally overparameterized, yet numerous empirical results show that even deeper and/or wider networks can achieve better performance [7]. Recent works have suggested that this does not always result from stronger computational capabilities of the bigger network compared to the smaller one. The smaller one, which can have enough expressive power to memorize the dataset, is usually harder to train, that is, it is harder for the smaller network to achieve a good generalization as the bigger network; however, with the right training it might [5, 8]. Structural re-parameterization methods can improve performance by adding overparameterization to a desired architecture, while dropping the excess weights after training, so the deployed network at inference time is of the original desired architecture. We focus on linear overparameterization (LO) as most structural re-parameterization relies on linearity.

Linear networks and their convergence with gradient de-

cent have been studied for a long time. It has been shown that using a deep linear neural network can have advantages over the use of a single linear layer, despite having the same computational power, e.g., by accelerating convergence [9]. This work also notes the vanishing gradient problem that can occur in deep linear networks with near zero initialization.

A logical next step is to introduce LO to regular non linear networks by adding redundant (transformable / collapsible) linear layers. ExpandNets [10] used LO by expanding every 2D convolution to multiple consecutive convolutions, keeping the same receptive field size. While being promising on shallow networks, it is less effective and even harmful for deep networks. Also, even for shallow networks, too deep LO can harm performance due to vanishing gradients [10, 9].

Contribution. In this paper, we propose a novel framework, *Kernel Filtering Linear Overparameterization (KFLO)*, for improving the performance of a target network. Our approach attempts to better utilize linear overparameterization and handle the vanishing gradient problem, while also being compatible with every type of filtering layer and staying easy to implement. Conceptually redundant linear layers are added to the network in a kernel filtering fashion, which does not make the network deeper from the network’s input perspective. The overparameterized weights can be discarded after being used in training, thus making inference time light-weight. Our method enjoys a much smaller number of calculations needed during training compared to using a straightforward deeper LO. A brief preview of KFLO can be found in Figure 1. Code is available at <https://github.com/AmitHenig/KFLO>.

2. METHOD

Notations. An input/signal/feature-map tensor has the shape $[batch \times channels \times spatial]$ (the *spatial* dimensions of an image for example will be *Height* \times *Width*). The kernel of a filtering layer (such as convolution or fully connected) has the shape $[output_channels \times input_channels \times spatial']$. The kernel is regarded to have *output_channels* filters, each of shape $[1 \times input_channels \times spatial']$. From this point forward, we will examine a 2D convolutional layer as our filtering layer, though our solution can easily be extended to any spatial dimension (as we will see later), and fully connected layers can be treated as inputs and filters having spatial dimensions sizes of 1. Furthermore, for simplicity, we will omit the bias, as it can be easily added.

LO. Consider a linear cascade comprised of two 2D convolutional layers with corresponding kernels/weights $\{\mathbb{W}_i\}_{i=1}^2$, where the second layer is a pointwise convolution with a trivial sliding window. This cascade is a LO of some 2D convolutional layer with weights \mathbb{W}' and the same sliding window properties as \mathbb{W}_1 . Meaning, that for the same input, kernel \mathbb{W}' gives the same output as the cascade $\{\mathbb{W}_i\}_{i=1}^2$:

$$\mathbb{W}' * \mathbb{X}_1 = \mathbb{W}_2 * (\mathbb{W}_1 * \mathbb{X}_1) = \mathbb{W}_2 * \mathbb{X}_2 = \mathbb{X}_3$$

$$\begin{aligned} \mathbb{X}_1 &\in \mathbb{R}^{batch \times ch_1 \times (H \times W)} & ; & \quad \mathbb{W}_1 \in \mathbb{R}^{ch_2 \times ch_1 \times (M \times N)} \\ \mathbb{X}_2 &\in \mathbb{R}^{batch \times ch_2 \times (H' \times W')} & ; & \quad \mathbb{W}_2 \in \mathbb{R}^{ch_3 \times ch_2 \times (1 \times 1)} \\ \mathbb{X}_3 &\in \mathbb{R}^{batch \times ch_3 \times (H' \times W')} & ; & \quad \mathbb{W}' \in \mathbb{R}^{ch_3 \times ch_1 \times (M \times N)} \end{aligned}$$

Note that neither H' , W' nor the receptive field size are affected by the pointwise convolution, and so are solely dependent on the properties of \mathbb{W}_1 . This LO equality holds when:

$$\mathbb{W}'(\delta_3, \delta_1, m, n) = \sum_{\delta_2}^{ch_2} \mathbb{W}_2(\delta_3, \delta_2, 0, 0) \cdot \mathbb{W}_1(\delta_2, \delta_1, m, n) \quad (1)$$

This can be shown by rearranging the convolution:

$$\mathbb{X}_3(b, \delta_3, h', w') = \sum_{\delta_2}^{ch_2} \mathbb{W}_2(\delta_3, \delta_2, 0, 0) \cdot \mathbb{X}_2(b, \delta_2, h', w')$$

where $\mathbb{X}_2(b, \delta_2, h', w')$ is

$$\sum_{\delta_1}^{ch_1} \sum_m^M \sum_n^N \mathbb{W}_1(\delta_2, \delta_1, m, n) \cdot \mathbb{X}_1(b, \delta_1, \mu_{(h', w')}^{m, n}, \nu_{(h', w')}^{m, n})$$

and μ, ν are the kernel-to-input spatial mapping functions, defined by the sliding window properties (e.g., the mapping functions for a trivial sliding window are $\mu = h' + m$ and $\nu = w' + n$). Eq. 1 shows that each filter in \mathbb{W}' is a combination of \mathbb{W}_1 filters, weighted by a corresponding filter in \mathbb{W}_2 .

Kernel Filtering. We introduce our kernel filtering by transforming Eq. 1 to a convolutional representation. Given the nature of the pointwise convolution, kernel filtering for this case may be represented with a 1D convolution; making its implementation independent of the number of spatial dimensions. This is done by convolution with the reshaped kernels. For some spatial location $t = (M \cdot N) \cdot \delta_1 + N \cdot m + n$:

$$\begin{aligned} \mathbb{W}_1^R &\in \mathbb{R}^{1 \times ch_2 \times (ch_1 \cdot M \cdot N)} & ; & \quad \mathbb{W}_1^R(0, \delta_2, t) = \mathbb{W}_1(\delta_2, \delta_1, m, n) \\ \mathbb{W}_2^R &\in \mathbb{R}^{ch_3 \times ch_2 \times (1)} & ; & \quad \mathbb{W}_2^R(\delta_3, \delta_2, 0) = \mathbb{W}_2(\delta_3, \delta_2, 0, 0) \\ \mathbb{W}'^R &\in \mathbb{R}^{1 \times ch_3 \times (ch_1 \cdot M \cdot N)} & ; & \quad \mathbb{W}'^R(0, \delta_3, t) = \mathbb{W}'(\delta_3, \delta_1, m, n) \end{aligned}$$

$$\mathbb{W}'^R(0, \delta_3, t) = \sum_{\delta_2}^{ch_2} \mathbb{W}_2^R(\delta_3, \delta_2, 0) \cdot \mathbb{W}_1^R(0, \delta_2, t)$$

We get a 1D pointwise convolution \mathbb{W}_2^R for input \mathbb{W}_1^R with batch size of 1, resulting in the output \mathbb{W}'^R . Getting \mathbb{W}' from \mathbb{W}'^R is of course done by reshaping.

$$\mathbb{W}' * \mathbb{X}_1 = (\mathbb{W}'^R)^R * \mathbb{X}_1 = (\mathbb{W}_2^R * \mathbb{W}_1^R)^R * \mathbb{X}_1 = \mathbb{X}_3$$

KFLO. It can be shown that the above LO and kernel filtering equivalences hold for any number of pointwise kernels used in the linear cascade $\{\mathbb{W}_i \in \mathbb{R}^{ch_{i+1} \times ch_i \times (1 \times 1)}\}_{i=2}^B$, reshaping them in the same manner to 1D pointwise convolutions:

$$\begin{aligned} \mathbb{X}_{B+1} &= \mathbb{W}_B * (\mathbb{W}_{B-1} * \dots * (\mathbb{W}_2 * (\mathbb{W}_1 * \mathbb{X}_1)) \dots) & (2) \\ &= (\mathbb{W}_B^R * (\mathbb{W}_{B-1}^R * \dots * (\mathbb{W}_2^R * \mathbb{W}_1^R) \dots))^R * \mathbb{X}_1 & (3) \\ &= (\mathbb{W}'^R)^R * \mathbb{X}_1 = \mathbb{W}' * \mathbb{X}_1 & (4) \end{aligned}$$

For a deployed kernel $\mathbb{W}' \in \mathbb{R}^{ch_{out} \times ch_{in} \times (M \times N)}$, our method is implemented during training by creating weights $\{\mathbb{W}_i\}_{i=1}^B$ and calculating the effective convolution used in the network using Eqs. 3, 4. After training is complete, the LO weights $\{\mathbb{W}_i\}_{i=1}^B$ are discarded, only keeping \mathbb{W}' to be used at inference time (see Figure 1). Applying KFLO on a depthwise convolution layer is done exactly as for conventional convolution, keeping in mind that the number of \mathbb{W}_1 input channels is $\frac{\#input_channels}{\#groups}$. We define a scalar **width multiplier** ρ , and set all configurable LO widths $\{ch_i\}_{i=2}^B$ to $Round(\rho \cdot ch_{out})$. We also denote B as the **linear depth multiplier**. KFLO for a desired deployed architecture is implemented by defining $B > 1$ and $\rho > 0$ (mainly choosing values greater than 1 for LO), and applying our method as described above for every filtering layer in the network (example with $B=2$, $\rho=4$ in Figure 1.c).

We present two explanations for KFLO success. First, we know that overparameterization helps train networks in general. Second, from fully linear deep networks we know that this redundant overparameterization can improve convergence. Comparing to a LO with vanilla feature filtering, KFLO does not make the trained architecture deeper, which is a big advantage as a common problem of deep LO/networks is the vanishing gradient which can really harm performance. Also, our kernel filtering gives a great reduction in the number of operations that LO adds. While in feature filtering each LO layer is applied to the full spatial size of the layer’s input tensor, with our kernel filtering the LO layers filter a tensor with spatial size of a convolution layer, which is usually significantly smaller than that of the layer’s input signal.

3. RELATION TO OTHER LO WORKS

DO-Conv [11] has a similar approach as ours. They utilize a linear depthwise separable convolutional layer, showing that it can be collapsed to a single convolution kernel. Similarly to us, DO-Conv suggests a kind of kernel filtering (kernel composition), to replace the vanilla feature filtering during training. There are two key differences between our method and DO-Conv. First, the linear convolutional cascade of our method is different than the linear depthwise separable convolution in DO-Conv. Second, while we performs kernel filtering by passing forward \mathbb{W}_1^R through the 1D linear cascade, DO-Conv passes the pointwise kernel backwards by filtering it with a transposed depthwise convolution. The work in Orthogonal Over-Parameterized Training [12] uses a tailored LO filtering weights, with further restrictions to enforce an orthogonal transformation on static weights. ACNet [13] suggests Asymmetric Convolution Block (ACB), which can be applied on a square convolution kernel. During training the kernel is replaced with three parallel paths, each with a convolution kernel followed by batch normalization, and their outputs are summed. In two of the paths the kernels are of 1D, one horizontal and the other vertical. This method does not

make the network deeper during training in terms of filtering kernels and they show improvement in both small and big networks. Note that the commonly used pair of convolution followed by batch normalization can also act as a LO on its own. Diverse Branch Block (DBB) [14] suggests LO by replacing the convolution kernel with an Inception-like DBB instance during training. ACB (ACNet) can be viewed as a special case of DBB. LO in training was also utilized in [6, 15, 16]. The work in [17] uses LO during training alongside block-level KD on a student network initially obtained by compressing the pre-trained teacher. LO is applied in a vanilla feature filtering fashion by adding a pointwise convolution solely at the end of each defined block and is collapsed/merged after training. Compared to them, our kernel filtering is applied on every filtering layer with further wider LO (width multiplier $\rho > 1$).

4. EXPERIMENTS

In this section we demonstrate our approach on two popular datasets: CIFAR-10 and CIFAR-100 [18]. We evaluate KFLO and compare it to both the vanilla approach and other relevant training methods. We present experimental results for different architectures and datasets, showing the advantages of our method. We also compare the performance of KFLO with different linear depth multipliers, B , and different width multipliers, ρ . Finally, we show that KFLO can have a similar effect of transfer learning. Note that with infinite memory and time, many methods can be combined and applied simultaneously, multiple times. However, in the scope of this paper we will only refer to training with no more than one LO method applied at a time.

Networks. WRN-D-K will denote a Wide Residual Network [7], with the wide-dropout residual block, of depth D (with a 3 residual groups structure) and widening factor K . In addition, we test with Cifar-quick [19] and VGG16 [20] where the two fully connected layers are replaced with global average pooling followed by a fully connected layer that outputs 512 channels, as [13]. For both architectures, batch normalization were added after every convolution layer, again as [13].

Implementation Details. When using our training method, KFLO is applied to all the convolutional and fully connected kernels. KFLO with a linear depth multiplier B and width multiplier ρ will be denoted by **KFLO $B \times \rho$** . We apply a weak weight decay of 10^{-9} on the pointwise convolutions in the linear cascades $\{\mathbb{W}_i\}_{i=2}^B$, and initialize them with identity (dirac). We apply the vanilla training weight decay value on the generated weights \mathbb{W}' , and no weight decay on the filtered kernels \mathbb{W}_1 , thus not adding a hyper parameter to tune. DO-Conv is implemented with D_{mul} equals to the filter’s spatial size $M \times N$, as in [11]. Exponential moving average of the network weights when trained in a vanilla fashion will be referred to as **EMA**. The Experiments ran on NVIDIA GeForce RTX 2080 Ti. We followed the experimental settings of [13],

Method	Cifar-quick	VGG16	wResnet-16-8	wResnet-28-2	wResnet-28-5	wResnet-28-10
vanilla	86.26 ±0.22	93.85 ±0.18	95.44 ±0.12	94.49 ±0.12	95.50 ±0.06	95.85 ±0.10
EMA	86.26 ±0.21	93.87 ±0.17	95.44 ±0.12	94.52 ±0.10	95.49 ±0.05	95.84 ±0.09
ACNet	86.83 ±0.28	94.41 ±0.12	95.62 ±0.14	94.75 ±0.18	95.79 ±0.14	96.14 ±0.06
DO-Conv	86.85 ±0.14	94.00 ±0.21	95.50 ±0.14	94.72 ±0.09	95.61 ±0.15	96.01 ±0.12
KFLO 2x4	87.41 ±0.23	94.70 ±0.09	95.68 ±0.07	94.89 ±0.23	95.89 ±0.09	96.22 ±0.12

Table 1. CIFAR-10 accuracy results, methods comparison.

Method	Cifar-quick	VGG16	wResnet-16-8	wResnet-28-2	wResnet-28-5	wResnet-28-10
vanilla	57.80 ±0.30	73.97 ±0.28	78.58 ±0.22	75.37 ±0.26	78.89 ±0.39	80.18 ±0.25
EMA	57.80 ±0.26	73.97 ±0.29	78.60 ±0.21	75.44 ±0.27	78.88 ±0.39	80.20 ±0.25
ACNet	58.38 ±0.39	74.95 ±0.22	79.10 ±0.18	76.02 ±0.32	79.64 ±0.32	81.36 ±0.14
DO-Conv	58.86 ±0.20	74.38 ±0.17	79.01 ±0.28	75.73 ±0.29	79.10 ±0.19	80.59 ±0.29
KFLO 2x4	59.78 ±0.25	75.69 ±0.14	79.74 ±0.24	76.07 ±0.29	79.96 ±0.27	81.38 ±0.15

Table 2. CIFAR-100 accuracy results, methods comparison.

excluding WRN Experiments where we followed [4].

Transfer learning (TL) experiments were conducted on CIFAR-100, with only 40% of the training data available (uniformly sampled from each class). Networks pre-trained in the vanilla fashion on the full CIFAR-10 training set were used for weights initialization. In the classic strategy, a pre-trained network is used as the starting point and it is fine-tuned with a learning rate of an order of magnitude smaller.

With KFLO, more than one pre-trained model can be utilized in weight initialization. Setting the width multiplier ρ to the number of pre-trained networks results in KFLO filtering kernels (\mathbb{W}_1) wide enough to hold all relevant pre-trained kernel filters. Before Training starts, we stack the corresponding pre-trained kernels and use them as initialization weights for the relevant filtering kernels. Layers on which KFLO is not applied are initialized with the weights of the first pre-trained network. This with the Identity initialization of KFLO’s LO pointwise convolutions, means that the starting network gives identical outputs to the first pre-trained network.

4.1. Results

Comparison. We first present results comparing our method to vanilla training, EMA, ACNet [13] and DO-Conv [11]. Results for CIFAR-10 and CIFAR-100 can be found in Tables 1 and 2, respectively. We can see that KFLO achieves the highest accuracy with all featured architectures.

Ablation study. We compare the performance of our method with different linear depth multipliers B , and different width multipliers ρ . We show results for CIFAR-10 in Table 3 and for CIFAR-100 in Table 4. It seems that in most cases, the bigger the width multiplier the better the performance boost is. However, a larger depth multiplier harms performance.

Transfer learning. Table 5 shows results for the TL problem. It displays results for both networks trained from scratch, and networks trained with TL - denoted with the TL suffix. We can see the expected performance boost TL gives vanilla training. Yet, it would appear that KFLO does not benefit much from the pre-trained weights, and sometimes it can even harm performance. It seems that the straight forward method

Method	wResnet-16-8	wResnet-28-2	wResnet-28-5	wResnet-28-10
KFLO 2x1	95.50 ±0.17	94.59 ±0.10	95.83 ±0.09	96.16 ±0.13
KFLO 2x2	95.62 ±0.11	94.60 ±0.18	95.84 ±0.13	96.28 ±0.10
KFLO 2x3	95.59 ±0.10	94.88 ±0.08	95.79 ±0.15	96.26 ±0.15
KFLO 2x4	95.68 ±0.07	94.89 ±0.23	95.89 ±0.09	96.22 ±0.12
KFLO 3x2	94.65 ±0.15	91.98 ±0.14	94.68 ±0.17	95.65 ±0.11

Table 3. CIFAR-10 accuracy results, KFLO ablation study.

Method	wResnet-16-8	wResnet-28-2	wResnet-28-5	wResnet-28-10
KFLO 2x1	79.47 ±0.21	75.69 ±0.39	79.60 ±0.24	81.35 ±0.23
KFLO 2x2	79.68 ±0.31	75.91 ±0.31	79.76 ±0.22	81.36 ±0.15
KFLO 2x3	79.59 ±0.16	76.00 ±0.29	79.96 ±0.18	81.46 ±0.15
KFLO 2x4	79.74 ±0.24	76.07 ±0.29	79.96 ±0.27	81.33 ±0.20
KFLO 3x2	76.51 ±0.26	68.21 ±0.39	75.69 ±0.22	78.93 ±0.12

Table 4. CIFAR-100 accuracy results, KFLO ablation study.

Method	wResnet-16-8	wResnet-28-2	wResnet-28-5	wResnet-28-10
vanilla	69.81 ±0.19	66.69 ±0.24	69.56 ±0.19	70.79 ±0.41
vanilla TL	71.20 ±0.41	67.12 ±0.25	71.43 ±0.42	73.15 ±0.33
KFLO 2x1	71.50 ±0.23	67.06 ±0.53	71.04 ±0.27	73.11 ±0.26
KFLO 2x1 TL	71.62 ±0.37	67.26 ±0.41	71.33 ±0.34	72.19 ±0.35
KFLO 2x4	71.86 ±0.14	67.66 ±0.34	71.63 ±0.27	73.36 ±0.29
KFLO 2x4 TL	71.97 ±0.32	67.85 ±0.32	71.60 ±0.22	72.55 ±0.26

Table 5. Transfer Learning accuracy results. Only 40% of the CIFAR-100 training set available. The TL suffix denotes using networks pre-trained on the full CIFAR-10 training set.

to incorporate the pre-trained networks with KFLO is not optimal. Nonetheless, KFLO (without TL) outperforms vanilla TL, and it would seem that it achieves a similar effect although it does not have access to the extra data (CIFAR-10).

5. CONCLUSION

This work introduced KFLO, a novel approach to train a neural network in a structural re-parameterization fashion by modeling kernels as the result of linear convolutional kernel filtering. Interestingly, this over-parameterization, which can create a very large redundancy, improves the network training. This stands in line with the recent both practical and theoretical findings that overparameterization improves the generalization of neural networks [21, 22, 23, 24].

We have focused on demonstrating our approach in supervised learning and transfer learning settings, though we believe that our proposed framework has large potential in being applied to other problems such as domain adaptation, semi-supervised learning and incremental learning. Moreover, one may consider in our scheme many other topologies for the kernel filtering block. We hope that this new tool with its different variations will be used to improve many other tasks. Its simplicity allows using and extending it to new setups easily.

Acknowledgement. This research was supported by Wipro and ERC-StG grant no. 757497

6. REFERENCES

- [1] Cristian Bucila, Rich Caruana, and Alexandru Niculescu-Mizil, “Model compression,” in *KDD '06*, 2006.
- [2] Geoffrey Hinton, Oriol Vinyals, and Jeff Dean, “Distilling the knowledge in a neural network,” in *NIPS Deep Learning Workshop*, 2014.
- [3] Tommaso Furlanello, Zachary Chase Lipton, Michael Tschannen, Laurent Itti, and Anima Anandkumar, “Born again neural networks,” in *ICML*, 2018.
- [4] Y. Zhang, T. Xiang, T. M. Hospedales, and H. Lu, “Deep mutual learning,” in *Conference on Computer Vision and Pattern Recognition (CVPR)*, 2018, pp. 4320–4328.
- [5] Seyed-Iman Mirzadeh, Mehrdad Farajtabar, Ang Li, Nir Levine, Akihiro Matsukawa, and Hassan Ghasemzadeh, “Improved knowledge distillation via teacher assistant,” in *AAAI*, 2020, pp. 5191–5198.
- [6] Xiaohan Ding, X. Zhang, Ningning Ma, Jungong Han, Guiguang Ding, and Jian Sun, “Repvvg: Making vgg-style convnets great again,” in *CVPR*, 2021.
- [7] Sergey Zagoruyko and Nikos Komodakis, “Wide residual networks,” in *British Machine Vision Conference*, 2016.
- [8] C. Zhang, S. Bengio, Moritz Hardt, B. Recht, and Oriol Vinyals, “Understanding deep learning requires rethinking generalization,” in *ICLR*, 2017.
- [9] Sanjeev Arora, N. Cohen, and Elad Hazan, “On the optimization of deep networks: Implicit acceleration by overparameterization,” *ICML 2018*, 2018.
- [10] Shuxuan Guo, Jose M. Alvarez, and M. Salzmann, “Expandnets: Linear over-parameterization to train compact convolutional networks,” *NeurIPS*, 2020.
- [11] Jinming Cao, Yangyan Li, M. Sun, Ying Chen, D. Lischinski, D. Cohen-Or, B. Chen, and C. Tu, “Do-conv: Depthwise over-parameterized convolutional layer,” *ArXiv*, vol. abs/2006.12030, 2020.
- [12] Weiyang Liu, Rongmei Lin, Z. Liu, J. Rehg, Li Xiong, and L. Song, “Orthogonal over-parameterized training,” *ArXiv*, vol. abs/2004.04690, 2020.
- [13] Xiaohan Ding, Yuchen Guo, Guiguang Ding, and J. Han, “Acnet: Strengthening the kernel skeletons for powerful cnn via asymmetric convolution blocks,” *2019 IEEE/CVF International Conference on Computer Vision (ICCV)*, pp. 1911–1920, 2019.
- [14] Xiaohan Ding, Xiangyu Zhang, Jungong Han, and Guiguang Ding, “Diverse branch block: Building a convolution as an inception-like unit,” *ArXiv*, vol. abs/2103.13425, 2021.
- [15] K. Bhardwaj, M. Milosavljevic, A. Chalfin, Naveen Suda, Liam O’Neil, Dibakar Gope, Lingchuan Meng, Ramon Matas Navarro, and Danny Loh, “Collapsible linear blocks for super-efficient super resolution,” *ArXiv*, vol. abs/2103.09404, 2021.
- [16] Shoufa Chen, Y. Chen, S. Yan, and Jiashi Feng, “Efficient differentiable neural architecture search with meta kernels,” *ArXiv*, vol. abs/1912.04749, 2019.
- [17] Tianhong Li, Jianguo Li, Zhuang Liu, and Changshui Zhang, “Few sample knowledge distillation for efficient network compression,” *2020 IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, pp. 14627–14635, 2020.
- [18] Alex Krizhevsky, “Learning multiple layers of features from tiny images,” 2009.
- [19] Jasper Snoek, H. Larochelle, and Ryan P. Adams, “Practical bayesian optimization of machine learning algorithms,” in *NIPS*, 2012.
- [20] Karen Simonyan and Andrew Zisserman, “Very deep convolutional networks for large-scale image recognition,” *CoRR*, vol. abs/1409.1556, 2015.
- [21] Zeyuan Allen-Zhu, Yuanzhi Li, and Yingyu Liang, “Learning and generalization in overparameterized neural networks, going beyond two layers,” in *Advances in Neural Information Processing Systems 32*, H. Wallach, H. Larochelle, A. Beygelzimer, F. Alché-Buc, E. Fox, and R. Garnett, Eds., pp. 6158–6169. Curran Associates, Inc., 2019.
- [22] Behnam Neyshabur, Zhiyuan Li, Srinadh Bhojanapalli, Yann LeCun, and Nathan Srebro, “The role of overparameterization in generalization of neural networks,” in *International Conference on Learning Representations*, 2019.
- [23] Luca Venturi, Afonso S. Bandeira, and Joan Bruna, “Spurious valleys in one-hidden-layer neural network optimization landscapes,” *Journal of Machine Learning Research*, vol. 20, no. 133, pp. 1–34, 2019.
- [24] Alon Brutzkus and Amir Globerson, “Why do larger models generalize better? a theoretical perspective via the xor problem,” in *ICML*, 2019.