

TRP: Trained Rank Pruning for Efficient Deep Neural Networks

Yuhui Xu¹, Yuxi Li¹, Shuai Zhang², Wei Wen³, Botao Wang²,
Yingyong Qi², Yiran Chen³, Weiyao Lin¹, Hongkai Xiong¹

¹Shanghai Jiao Tong University ²Qualcomm AI Research ³Duke University

{yuhuixu, lyxok1, wylin, xionghongkai}@sjtu.edu.cn,
{shuazhan, botaow, yingyong}@qti.qualcomm.com, {wei.wen, yiran.chen}@duke.edu

Abstract

To enable DNNs on edge devices like mobile phones, low-rank approximation has been widely adopted because of its solid theoretical rationale and efficient implementations. Several previous works attempted to directly approximate a pre-trained model by low-rank decomposition; however, small approximation errors in parameters can ripple over a large prediction loss. As a result, performance usually drops significantly and a sophisticated effort on fine-tuning is required to recover accuracy. Apparently, it is not optimal to separate low-rank approximation from training. Unlike previous works, this paper integrates low rank approximation and regularization into the training process. We propose Trained Rank Pruning (TRP), which alternates between low rank approximation and training. TRP maintains the capacity of the original network while imposing low-rank constraints during training. A nuclear regularization optimized by stochastic sub-gradient descent is utilized to further promote low rank in TRP. The TRP trained network inherently has a low-rank structure, and is approximated with negligible performance loss, thus eliminating the fine-tuning process after low rank decomposition. The proposed method is comprehensively evaluated on CIFAR-10 and ImageNet, outperforming previous compression methods using low rank approximation.

1 Introduction

Deep Neural Networks (DNNs) have shown remarkable success in many computer vision tasks. Despite the high performance in server-based DNNs powered by cutting-edge parallel computing hardware, most state-of-the-art architectures are not yet ready to be deployed on mobile devices due to the limitations on computational capacity, memory and power.

This work was supported in part by the National Natural Science Foundation of China under Grants 61720106001, 61932022, and in part by the Program of Shanghai Academic Research Leader under Grant 17XD1401900.

To address this problem, many network compression and acceleration methods have been proposed. Pruning based methods [Han *et al.*, 2015b; He *et al.*, 2017; Liu *et al.*, 2017; Luo *et al.*, 2017] explore the sparsity in weights and filters. Quantization based methods [Han *et al.*, 2015b; Zhou *et al.*, 2017; Courbariaux and Bengio, 2016; Rastegari *et al.*, 2016; Xu *et al.*, 2018] reduce the bit-width of network parameters. Low-rank decomposition [Denton *et al.*, 2014; Jaderberg *et al.*, 2014; Guo *et al.*, 2018; Wen *et al.*, 2017; Alvarez and Salzmann, 2017] minimizes the channel-wise and spatial redundancy by decomposing the original network into a compact one with low-rank layers. In addition, efficient architectures [Sandler *et al.*, 2018; Ma *et al.*, 2018] are carefully designed to facilitate mobile deployment of deep neural networks. Different from precedent works, this paper proposes a novel approach to design low-rank networks.

Low-rank networks can be trained directly from scratch. However, it is difficult to obtain satisfactory results for several reasons. (1) *Low capacity*: compared with the original full rank network, the capacity of a low-rank network is limited, which causes difficulties in optimizing its performances. (2) *Deep structure*: low-rank decomposition typically doubles the number of layers in a network. The additional layers make numerical optimization much more vulnerable to gradients explosion and/or vanishing. (3) *Heuristic rank selection*: the rank of decomposed network is often chosen as a hyperparameter based on pre-trained networks; this may not be the optimal rank for the network trained from scratch.

Alternatively, several previous works [Zhang *et al.*, 2016; Guo *et al.*, 2018; Jaderberg *et al.*, 2014] attempted to decompose pre-trained models in order to get initial low-rank networks. However, the heuristically imposed low-rank could incur huge accuracy loss and network retraining is needed to recover the performance of the original network as much as possible. Some attempts were made to use sparsity regularization [Wen *et al.*, 2017; Chen *et al.*, 2015] to constrain the network into a low-rank space. Though sparsity regularization reduces the error incurred by decomposition to some extent, performance still degrades rapidly when compression rate increases.

This paper is an extension of [Xu *et al.*, 2019]. In this paper, we propose a new method, namely Trained Rank Pruning (TRP), for training low-rank networks. We embed the low-rank decomposition into the training process by gradu-

ally pushing the weight distribution of a well functioning network into a low-rank form, where all parameters of the original network are kept and optimized to maintain its capacity. We also propose a stochastic sub-gradient descent optimized nuclear regularization that further constrains the weights in a low-rank space to boost the TRP. The proposed solution is illustrated in Fig. 1.

Overall, our contributions are summarized below.

1. A new training method called the TRP is presented by explicitly embedding the low-rank decomposition into the network training;
2. A nuclear regularization is optimized by stochastic sub-gradient descent to boost the performance of the TRP;
3. Improving inference acceleration and reducing approximation accuracy loss in both channel-wise and spatial-wise decomposition methods.

2 Related Works

A lot of works have been proposed to accelerate the inference process of deep neural networks. Briefly, these works could be categorized into three main categories: quantization, pruning, and low-rank decomposition.

Quantization Weight quantization methods include training a quantized model from scratch [Chen *et al.*, 2015; Courbariaux and Bengio, 2016; Rastegari *et al.*, 2016] or converting a pre-trained model into quantized representation [Zhou *et al.*, 2017; Han *et al.*, 2015a; Xu *et al.*, 2018]. The quantized weight representation includes binary value [Rastegari *et al.*, 2016; Courbariaux and Bengio, 2016] or hash buckets [Chen *et al.*, 2015]. Note that our method is inspired by the scheme of combining quantization with training process, *i.e.* we embed the low-rank decomposition into training process to explicitly guide the parameter to a low-rank form.

Pruning Non-structured and structured sparsity are introduced by pruning. [Han *et al.*, 2015b] proposes to prune unimportant connections between neural units with small weights in a pre-trained CNN. [Wen *et al.*, 2016] utilizes group Lasso strategy to learn the structure sparsity of networks. [Liu *et al.*, 2017] adopts a similar strategy by explicitly imposing scaling factors on each channel to measure the importance of each connection and dropping those with small weights. In [He *et al.*, 2017], the pruning problem is formulated as a data recovery problem. Pre-trained filters are re-weighted by minimizing a data recovery objective function. Channels with smaller weight are pruned. [Luo *et al.*, 2017] heuristically selects filters using change of next layer’s output as a criterion.

Low-rank decomposition Original models are decomposed into compact ones with more lightweight layers. [Jaderberg *et al.*, 2014] considers both the spatial-wise and channel-wise redundancy and proposes decomposing a filter into two cascaded asymmetric filters. [Zhang *et al.*, 2016] further assumes the feature map lie in a low-rank subspace and decompose the convolution filter into $k \times k$ followed by 1×1 filters via SVD. [Guo *et al.*, 2018] exploits the low-rank assumption of convolution filters and decompose a regular convolution into several depth-wise and point-wise convolution structures. Although these works achieved notable performance in

network compression, all of them are based on the low-rank assumption. When such assumption is not completely satisfied, large prediction error may occur.

Alternatively, some other works [Wen *et al.*, 2017; Alvarez and Salzmann, 2017] implicitly utilize sparsity regularization to direct the neural network training process to learn a low-rank representation. Our work is similar to this low-rank regularization method. However, in addition to appending an implicit regularization during training, we impose an explicit sparsity constraint in our training process and prove that our approach can push the weight distribution into a low-rank form quite effectively.

3 Methodology

3.1 Preliminaries

Formally, the convolution filters in a layer can be denoted by a tensor $W \in \mathbb{R}^{n \times c \times k_w \times k_h}$, where n and c are the number of filters and input channels, k_h and k_w are the height and width of the filters. An input of the convolution layer $F_i \in \mathbb{R}^{c \times x \times y}$ generates an output as $F_o = W * F_i$. Channel-wise correlation [Zhang *et al.*, 2016] and spatial-wise correlation [Jaderberg *et al.*, 2014] are explored to approximate convolution filters in a low-rank space. In this paper, we focus on these two decomposition schemes. However, unlike the previous works, we propose a new training scheme TRP to obtain a low-rank network without re-training after decomposition.

3.2 Trained Rank Pruning

Trained Rank Pruning (TRP) is motivated by the strategies of training quantized nets. One of the gradient update schemes to train quantized networks from scratch [Li *et al.*, 2017] is

$$w^{t+1} = Q(w^t - \alpha \nabla f(w^t)) \quad (1)$$

where $Q(\cdot)$ is the quantization function, w^t denote the parameter in the t^{th} iteration. Parameters are quantized by $Q(\cdot)$ before updating the gradients.

In contrast, we propose a simple yet effective training scheme called Trained Rank Pruning (TRP) in a periodic fashion:

$$W^{t+1} = \begin{cases} W^t - \alpha \nabla f(W^t) & t \% m \neq 0 \\ T^z - \alpha \nabla f(T^z) & t \% m = 0 \end{cases} \quad (2)$$

$$T^z = \mathcal{D}(W^t), \quad z = t/m$$

where $\mathcal{D}(\cdot)$ is a low-rank tensor approximation operator, α is the learning rate, t indexes the iteration and z is the iteration of the operator \mathcal{D} , with m being the period for the low-rank approximation.

At first glance, this TRP looks very simple. An immediate concern arises: can the iterations guarantee the rank of the parameters converge, and more importantly would not increase when they are updated in this way? A positive answer (see Theorem 2) given in our theoretical analysis will certify the legitimacy of this algorithm.

For the network quantization, if the gradients are smaller than the quantization, the gradient information would be totally lost and become zero. However, it will not happen in

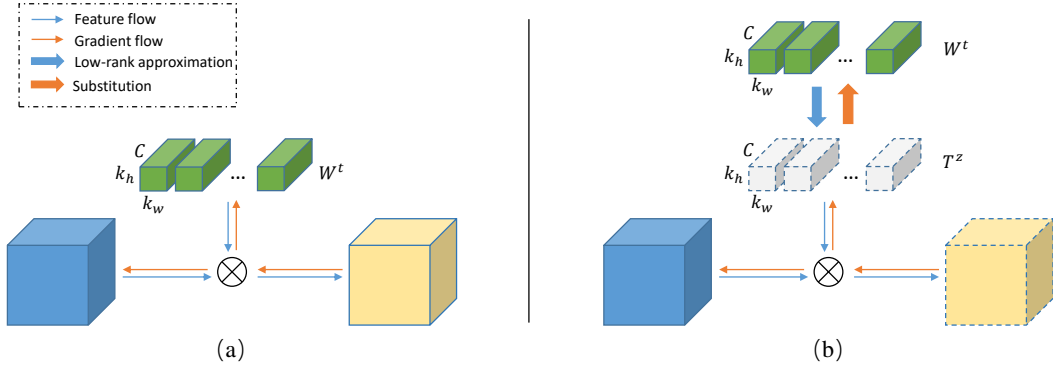


Figure 1: The training of TRP consists of two parts as illustrated in (a) and (b). (a) one normal iteration with forward-backward broadcast and weight update. (b) one training iteration inserted by TRP, where the low-rank approximation is first applied on filters before convolution. During backward propagation, the gradients are directly added on low-rank filters and the original weights are substituted by updated low-rank filters. (b) is applied once every m iterations (*i.e.* when gradient update iteration $t = zm, z = 0, 1, 2, \dots$), otherwise (a) is applied.

TRP because the low-rank operator is applied on the weight tensor. Furthermore, we apply low-rank approximation every m SGD iterations. This saves training time to a large extent. As illustrated in Fig. 1, for every m iterations, we perform low-rank approximation on the original filters, while gradients are updated on the resultant low-rank form. Otherwise, the network is updated via the normal SGD. Our training scheme could be combined with any low-rank operators. In the proposed work, we choose the low-rank techniques proposed in [Jaderberg *et al.*, 2014] and [Zhang *et al.*, 2016], both of which transform the 4-dimensional filters into 2D matrix and then apply the truncated singular value decomposition (TSVD). The SVD of matrix W^t can be written as:

$$W^t = \sum_{i=1}^{\text{rank}(W^t)} \sigma_i \cdot U_i \cdot (V_i)^T \quad (3)$$

where σ_i is the singular value of W^t with $\sigma_1 \geq \sigma_2 \geq \dots \geq \sigma_{\text{rank}(W^t)}$, and U_i and V_i are the singular vectors. The parameterized TSVD($W^t; e$) is to find the smallest integer k such that

$$\sum_{j=k+1}^{\text{rank}(W^t)} (\sigma_j)^2 \leq e \sum_{i=1}^{\text{rank}(W^t)} (\sigma_i)^2 \quad (4)$$

where e is a pre-defined hyper-parameter of the energy-pruning ratio, $e \in (0, 1)$.

After truncating the last $n - k$ singular values, we transform the low-rank 2D matrix back to 4D tensor. Compared with directly training low-rank structures from scratch, the proposed TRP has following advantages.

(1) Unlike updating the decomposed filters independently of the network training in literature [Zhang *et al.*, 2016; Jaderberg *et al.*, 2014], we update the network directly on the original 4D shape of the decomposed parameters, which enable jointly network decomposition and training by preserving its discriminative capacity as much as possible.

(2) Since the gradient update is performed based on the original network structure, there will be no exploding and vanishing gradients problems caused by additional layers.

(3) The rank of each layer is automatically selected during the training. We will prove a theorem certifying the rank of network weights convergence and would not increase in section 3.4.

3.3 Nuclear Norm Regularization

Nuclear norm is widely used in matrix completion problems. Recently, it is introduced to constrain the network into low-rank space during the training process [Alvarez and Salzmann, 2017].

$$\min \left\{ f(x; w) + \lambda \sum_{l=1}^L \|W_l\|_* \right\} \quad (5)$$

where $f(\cdot)$ is the objective loss function, nuclear norm $\|W_l\|_*$ is defined as $\|W_l\|_* = \sum_{i=1}^{\text{rank}(W_l)} \sigma_i^i$, with σ_i^i the singular values of W_l . λ is a hyper-parameter setting the influence of the nuclear norm. In [Alvarez and Salzmann, 2017] the proximity operator is applied in each layer independently to solve Eq. (5). However, the proximity operator is split from the training process and doesn't consider the influence within layers.

In this paper, we utilize stochastic sub-gradient descent [Avron *et al.*, 2012] to optimize nuclear norm regularization in the training process. Let $W = U\Sigma V^T$ be the SVD of W and let U_{tru}, V_{tru} be U, V truncated to the first $\text{rank}(W)$ columns or rows, then $U_{tru}V_{tru}^T$ is the sub-gradient of $\|W\|_*$ [Watson, 1992]. Thus, the sub-gradient of Eq. (5) in a layer is

$$\nabla f + \lambda U_{tru}V_{tru}^T \quad (6)$$

The nuclear norm and loss function are optimized simultaneously during the training of the networks and can further be combined with the proposed TRP.

3.4 Theoretic Analysis

In this section, we analyze the rank convergence of TRP from the perspective of matrix perturbation theory [Stewart, 1990]. We prove that rank in TRP is monotonously decreasing, *i.e.*, the model gradually converges to a more sparse model.

Let A be an $m \times n$ matrix, without loss of generality, $m \geq n$. $\Sigma = \text{diag}(\sigma_1, \dots, \sigma_n)$ and $\sigma_1 \geq \sigma_2 \geq \dots \geq \sigma_n$. $\tilde{\Sigma}$ is the diagonal matrix composed by all singular values of A . Let $\tilde{A} = A + E$ be a perturbation of A , and E is the noise matrix. $\tilde{\Sigma} = \text{diag}(\tilde{\sigma}_1, \dots, \tilde{\sigma}_n)$ and $\tilde{\sigma}_1 \geq \tilde{\sigma}_2 \geq \dots \geq \tilde{\sigma}_n$. $\tilde{\sigma}_i$ is the singular values of \tilde{A} . The basic perturbation bounds for the singular values of a matrix are given by

Theorem 1. *Mirsky's theorem [Mirsky, 1960]:*

$$\sqrt{\sum_i |\tilde{\sigma}_i - \sigma_i|^2} \leq \|E\|_F \quad (7)$$

where $\|\cdot\|_F$ is the Frobenius norm. Then the following corollary can be inferred from Theorem 1,

Corollary 1. *Let B be any $m \times n$ matrix of rank not greater than k , i.e. the singular values of B can be denoted by $\varphi_1 \geq \dots \geq \varphi_k \geq 0$ and $\varphi_{k+1} = \dots = \varphi_n = 0$. Then*

$$\|B - A\|_F \geq \sqrt{\sum_{i=1}^n |\varphi_i - \sigma_i|^2} \geq \sqrt{\sum_{j=k+1}^n \sigma_j^2} \quad (8)$$

Below, we will analyze the training procedure of the proposed TRP. Note that W below are all transformed into 2D matrix. In terms of Eq. (2), the training process between two successive TSVD operations can be rewritten as Eq. (9)

$$\begin{aligned} W^t &= T^z = \text{TSVD}(W^t; e) \\ W^{t+m} &= T^z - \alpha \sum_{i=0}^{m-1} \nabla f(W^{t+i}) \\ T^{z+1} &= \text{TSVD}(W^{t+m}; e) \end{aligned} \quad (9)$$

where W^t is the weight matrix in the t -th iteration. T^z is the weight matrix after applying TSVD over W^t . $\nabla f(W^{t+i})$ is the gradient back-propagated during the $(t+i)$ -th iteration. $e \in (0, 1)$ is the predefined energy threshold. Then we have following theorem.

Theorem 2. *Assume that $\|\alpha \nabla f\|_F$ has an upper bound G , if $G < \frac{\sqrt{e}}{m} \|W^{t+m}\|_F$, then $\text{rank}(T^z) \geq \text{rank}(T^{z+1})$.*

Proof. We denote σ_j^t and σ_j^{t+m} as the singular values of W^t and W^{t+m} respectively. Then at the t -th iteration, given the energy ratio threshold e , the TSVD operation tries to find the singular value index $k \in [0, n-1]$ such that :

$$\begin{aligned} \sum_{j=k+1}^n (\sigma_j^t)^2 &< e \|W^t\|_F^2 \\ \sum_{j=k}^n (\sigma_j^t)^2 &\geq e \|W^t\|_F^2 \end{aligned} \quad (10)$$

In terms of Eq. (10), T^z is a k rank matrix, i.e, the last $n-k$ singular values of T^z are equal to 0. According to Corollary 1, we can derive that:

$$\begin{aligned} \|W^{t+m} - T^z\|_F &= \left\| \alpha \sum_{i=0}^{m-1} \nabla f^{t+i} \right\|_F \\ &\geq \sqrt{\sum_{j=k+1}^n (\sigma_j^{t+m})^2} \end{aligned} \quad (11)$$

Given the assumption $G < \frac{\sqrt{e}}{m} \|W^{t+m}\|_F$, we can get:

$$\begin{aligned} \frac{\sqrt{\sum_{j=k+1}^n (\sigma_j^{t+m})^2}}{\|W^{t+m}\|_F} &\leq \frac{\|\alpha \sum_{i=0}^{m-1} \nabla f^{t+i}\|_F}{\|W^{t+m}\|_F} \\ &\leq \frac{\sum_{i=0}^{m-1} \|\alpha \nabla f^{t+i}\|_F}{\|W^{t+m}\|_F} \\ &\leq \frac{mG}{\|W^{t+m}\|_F} < \sqrt{e} \end{aligned} \quad (12)$$

Eq. (12) indicates that since the perturbations of singular values are bounded by the parameter gradients, if we properly select the TSVD energy ratio threshold e , we could guarantee that if $n-k$ singular values are pruned by previous TSVD iteration, then before the next TSVD, the energy for the last $n-k$ singular values is still less than the pre-defined energy threshold e . Thus TSVD should keep the number of pruned singular values or drop more to achieve the criterion in Eq. (10), consequently a weight matrix with lower or same rank is obtained, i.e. $\text{Rank}(T^z) \geq \text{Rank}(T^{z+1})$. We further confirm our analysis about the variation of rank distribution in Section 4. \square

Model	Top 1 (%)	Speed up
R-20 (baseline)	91.74	1.00×
R-20 (TRP1)	90.12	1.97×
R-20 (TRP1+Nu)	90.50	2.17 ×
R-20 ([Zhang <i>et al.</i> , 2016])	88.13	1.41×
R-20 (TRP2)	90.13	2.66×
R-20 (TRP2+Nu)	90.62	2.84 ×
R-20 ([Jaderberg <i>et al.</i> , 2014])	89.49	1.66×
R-56 (baseline)	93.14	1.00×
R-56 (TRP1)	92.77	2.31×
R-56 (TRP1+Nu)	91.85	4.48 ×
R-56 ([Zhang <i>et al.</i> , 2016])	91.56	2.10×
R-56 (TRP2)	92.63	2.43×
R-56 (TRP2+Nu)	91.62	4.51 ×
R-56 ([Jaderberg <i>et al.</i> , 2014])	91.59	2.10×
R-56 [He <i>et al.</i> , 2017]	91.80	2.00×
R-56 [Li <i>et al.</i> , 2016]	91.60	2.00×

Table 1: Experiment results on CIFAR-10. "R-" indicates ResNet-.

4 Experiments

4.1 Datasets and Baseline

We evaluate the performance of TRP scheme on two common datasets, CIFAR-10 [Krizhevsky and Hinton, 2009] and ImageNet [Deng *et al.*, 2009]. The CIFAR-10 dataset consists of colored natural images with 32×32 resolution and has totally 10 classes. The ImageNet dataset consists of 1000 classes of images for recognition task. For both of the datasets, we adopt ResNet [He *et al.*, 2016] as our baseline model since it is widely used in different vision tasks. We use ResNet-20, ResNet-56 for CIFAR-10 and ResNet-18, ResNet-50 for ImageNet. For evaluation metric, we adopt top-1 accuracy on

CIFAR-10 and top-1, top-5 accuracy on ImageNet. To measure the acceleration performance, we compute the FLOPs ratio between baseline and decomposed models to obtain the final speedup rate. Wall-clock CPU and GPU time is also compared. Apart from the basic decomposition methods, we compare the performance with other state-of-the-art acceleration algorithms [He *et al.*, 2017; Li *et al.*, 2016; Luo *et al.*, 2017; Zhou *et al.*, 2019].

Method	Top1(%)	Top5(%)	Speed up
Baseline	69.10	88.94	1.00×
TRP1	65.46	86.48	1.81×
TRP1+Nu	65.39	86.37	2.23 ×
[Zhang <i>et al.</i> , 2016] ¹	-	83.69	1.39×
[Zhang <i>et al.</i> , 2016]	63.10	84.44	1.41×
TRP2	65.51	86.74	2.60×
TRP2+Nu	65.34	86.61	3.18 ×
[Jaderberg <i>et al.</i> , 2014]	62.80	83.72	2.00×

Table 2: Results of ResNet-18 on ImageNet.

Method	Top1(%)	Top5(%)	Speed up
Baseline	75.90	92.70	1.00×
TRP1+Nu	72.69	91.41	2.30 ×
TRP1+Nu	74.06	92.07	1.80×
[Zhang <i>et al.</i> , 2016]	71.80	90.2	1.50×
[He <i>et al.</i> , 2017]	-	90.80	2.00
[Luo <i>et al.</i> , 2017]	72.04	90.67	1.58
[Luo <i>et al.</i> , 2018]	72.03	90.99	2.26
[Zhou <i>et al.</i> , 2019]	71.50	90.20	2.30

Table 3: Results of ResNet-50 on ImageNet.

4.2 Implementation Details

We implement our TRP scheme with NVIDIA 1080 Ti GPUs. For training on CIFAR-10, we start with base learning rate of 0.1 to train 164 epochs and degrade the value by a factor of 10 at the 82-th and 122-th epoch. For ImageNet, we directly finetune the model with TRP scheme from the pre-trained baseline with learning rate 0.0001 for 10 epochs. We adopt SGD solver to update weight and set the weight decay value as 10^{-4} and momentum value as 0.9. The accuracy improvement enabled by data dependent decomposition vanishes after fine-tuning. So we simply adopt the retrained data independent decomposition as our basic methods.

4.3 Results on CIFAR-10

Settings. Experiments on channel-wise decomposition (TRP1) and spatial-wise decomposition (TRP2) are both considered. The TSVD energy threshold in TRP and TRP+Nu is 0.02 and the nuclear norm weight λ is set as 0.0003. We decompose both the 1×1 and 3×3 layers in ResNet-56.

Results. As shown in Table 1, for both spatial-wise and channel-wise decomposition, the proposed TRP outperforms basic methods [Zhang *et al.*, 2016; Jaderberg *et al.*, 2014] on

¹the implementation of [Guo *et al.*, 2018]

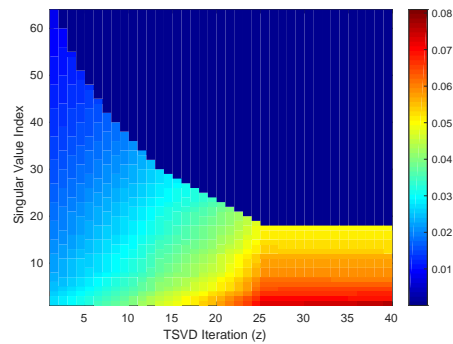


Figure 2: Visualization of rank selection, taken from the res3-1-2 convolution layer in ResNet-20 trained on CIFAR-10.

ResNet-20 and ResNet-56. Results become even better when nuclear regularization is used. For example, in the channel-wise decomposition (TRP2) of ResNet-56, results of TRP combined with nuclear regularization can even achieve $2 \times$ speed up rate than [Zhang *et al.*, 2016] with same accuracy drop. TRP also outperforms filter pruning [Li *et al.*, 2016] and channel pruning [He *et al.*, 2017]. The channel decomposed TRP trained ResNet-56 can achieve 92.77% accuracy with $2.31 \times$ acceleration, while [He *et al.*, 2017] is 91.80% and [Li *et al.*, 2016] is 91.60%. With nuclear regularization, our methods can approximately double the acceleration rate of [He *et al.*, 2017] and [Li *et al.*, 2016] with higher accuracy.

4.4 Results on ImageNet

Settings. We choose ResNet-18 and ResNet-50 as our baseline models. The TSVD energy threshold ϵ is set as 0.005. λ of nuclear norm regularization is 0.0003 for both ResNet-18 and ResNet-50. We decompose both the 3×3 and 1×1 Convolution layers in ResNet-50. TRP1 is the channel-wise decomposition and TRP2 is the spatial-wise decomposition.

Results. The results on ImageNet are shown in Table 2 and Table 3. For ResNet-18, our method outperforms the basic methods [Zhang *et al.*, 2016; Jaderberg *et al.*, 2014]. For example, in the channel-wise decomposition, TRP obtains $1.81 \times$ speed up rate with 86.48% Top5 accuracy on ImageNet which outperforms both the data-driven [Zhang *et al.*, 2016]¹ and data independent [Zhang *et al.*, 2016] methods by a large margin. Nuclear regularization can increase the speed up rates with the same accuracy.

For ResNet-50, to better validate the effectiveness of our method, we also compare TRP with pruning based methods. With $1.80 \times$ speed up, our decomposed ResNet-50 can obtain 74.06% Top1 and 92.07% Top5 accuracy which is much higher than [Luo *et al.*, 2017]. The TRP achieves $2.23 \times$ acceleration which is higher than [He *et al.*, 2017] with the same 1.4% Top5 degrade. Besides, with the same $2.30 \times$ acceleration rate, our performance is better than [Zhou *et al.*, 2019].

4.5 Rank Variation

To analyze the variation of rank distribution during training, we further conduct an experiment on the CIFAR-10 dataset with ResNet-20 and extract the weight from the *res3-1-2* convolution layer with channel-wise decomposition as our TRP

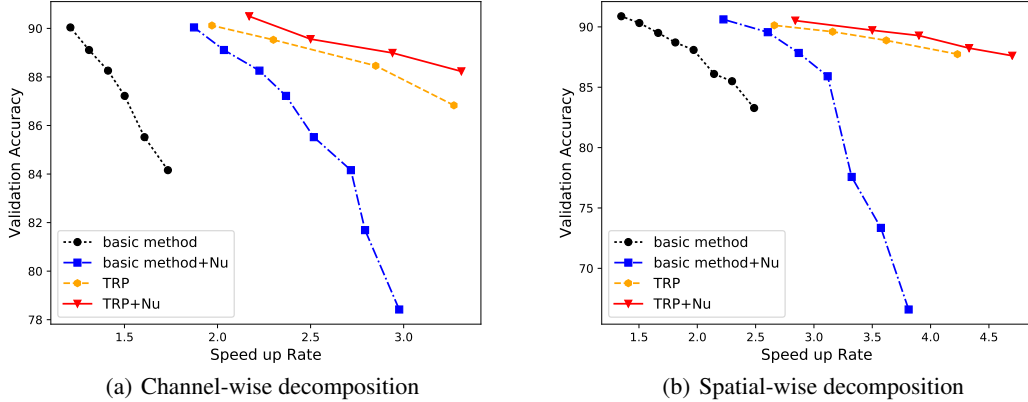


Figure 3: Ablation study on ResNet-20. Basic methods are data-independent decomposition methods (channel or spatial) with finetuning.

scheme. After each TSVD, we compute the normalized energy ratio $ER(i)$ for each singular value σ_i as Eq. (13).

$$ER(i) = \frac{\sigma_i^2}{\sum_{j=0}^{rank(T^z)} \sigma_j^2} \quad (13)$$

we record for totally 40 iterations of TSVD with period $m = 20$, which is equal to 800 training iterations, and our energy threshold e is pre-defined as 0.05. Then we visualize the variation of ER in Fig. 2. During our training, we observe that the theoretic bound value $\max_t \frac{mG}{\|W^t\|_F} \approx 0.092 < \sqrt{e} \approx 0.223$, which indicates that our basic assumption in theorem 2 always holds for the initial training stage.

And this phenomenon is also reflected in Fig. 2, at the beginning, the energy distribution is almost uniform w.r.t each singular value, and the number of dropped singular values increases after each TSVD iteration and the energy distribution becomes more dense among singular values with smaller index. Finally, the rank distribution converges to a certain point where the smallest energy ratio exactly reaches our threshold e and TSVD will not cut more singular values.

4.6 Ablation Study

In order to show the effectiveness of different components of our method, we compare four training schemes, basic methods [Zhang *et al.*, 2016; Jaderberg *et al.*, 2014], basic methods combined with nuclear norm regularization, TRP and TRP combined with nuclear norm regularization. The results are shown in Fig. 3. We can have following observations:

(1) *Nuclear norm regularization* After combining nuclear norm regularization, basic methods improve by a large margin. Since Nuclear norm regularization constrains the filters into low rank space, the loss caused by TSVD is smaller than the basic methods.

(2) *Trained rank pruning* As depicted in Fig. 3, when the speed up rate increases, the performance of basic methods and basic methods combined with nuclear norm regularization degrades sharply. However, the proposed TRP degrades very slowly. This indicates that by reusing the capacity of the network, TRP can learn a better low-rank feature representations than basic methods. The gain of nuclear norm regular-

ization on TRP is not as big as basic methods because TRP has already induced the parameters into low-rank space by embedding TSVD in training process.

Model	GPU time (ms)	CPU time (ms)
Baseline	0.45	118.02
TRP1+Nu (channel)	0.33	64.75
TRP2+Nu (spatial)	0.31	49.88

Table 4: Actual inference time per image on ResNet-18.

4.7 Runtime Speed up of Decomposed Networks

We further evaluate the actual runtime speed up of the compressed Network as shown in Table 4. Our experiment is conducted on a platform with one Nvidia 1080Ti GPU and Xeon E5-2630 CPU. The models we used are the original ResNet-18 and decomposed models by TRP1+Nu and TRP2+Nu. From the results, we observe that on CPU our TRP scheme achieves more salient acceleration performance. Overall the spatial decomposition combined with our TRP+Nu scheme has better performance. Because cuDNN is not friendly for 1×3 and 3×1 kernels, the actual speed up of spatial-wise decomposition is not as obvious as the reduction of FLOPs.

5 Conclusion

In this paper, we proposed a new scheme Trained Rank Pruning (TRP) for training low-rank networks. It leverages capacity and structure of the original network by embedding the low-rank approximation in the training process. Furthermore, we propose stochastic sub-gradient descent optimized nuclear norm regularization to boost the TRP. The proposed TRP can be incorporated with any low-rank decomposition method. On CIFAR-10 and ImageNet datasets, we have shown that our methods can outperform basic methods and other pruning based methods both in channel-wise decomposition and spatial-wise decomposition.

References

- [Alvarez and Salzmman, 2017] Jose M Alvarez and Mathieu Salzmman. Compression-aware training of deep networks. In *NIPS*, 2017.
- [Avron *et al.*, 2012] Haim Avron, Satyen Kale, Shiva Prasad Kasiviswanathan, and Vikas Sindhwani. Efficient and practical stochastic subgradient descent for nuclear norm regularization. In *ICML*, 2012.
- [Chen *et al.*, 2015] Wenlin Chen, James Wilson, Stephen Tyree, Kilian Weinberger, and Yixin Chen. Compressing neural networks with the hashing trick. In *ICML*, 2015.
- [Courbariaux and Bengio, 2016] Matthieu Courbariaux and Yoshua Bengio. Binarynet: Training deep neural networks with weights and activations constrained to +1 or -1. *arXiv preprint arXiv:1602.02830*, 2016.
- [Deng *et al.*, 2009] Jia Deng, Wei Dong, Richard Socher, Li-Jia Li, Kai Li, and Li Fei-Fei. Imagenet: A large-scale hierarchical image database. *CVPR*, 2009.
- [Denton *et al.*, 2014] Emily Denton, Wojciech Zaremba, Joan Bruna, Yann Lecun, and Rob Fergus. Exploiting linear structure within convolutional networks for efficient evaluation. In *NIPS*, 2014.
- [Guo *et al.*, 2018] Jianbo Guo, Yuxi Li, Weiyao Lin, Yurong Chen, and Jianguo Li. Network decoupling: From regular to depthwise separable convolutions. In *BMVC*, 2018.
- [Han *et al.*, 2015a] Song Han, Huizi Mao, and William J. Dally. Deep compression: Compressing deep neural network with pruning, trained quantization and huffman coding. *CoRR*, abs/1510.00149, 2015.
- [Han *et al.*, 2015b] Song Han, Jeff Pool, John Tran, and William Dally. Learning both weights and connections for efficient neural network. In *NIPS*, 2015.
- [He *et al.*, 2016] Kaiming He, Xiangyu Zhang, Shaoqing Ren, and Jian Sun. Deep residual learning for image recognition. 2016.
- [He *et al.*, 2017] Yihui He, Xiangyu Zhang, and Jian Sun. Channel pruning for accelerating very deep neural networks. In *ICCV*, 2017.
- [Jaderberg *et al.*, 2014] Max Jaderberg, Andrea Vedaldi, and Andrew Zisserman. Speeding up convolutional neural networks with low rank expansions. *arXiv preprint arXiv:1405.3866*, 2014.
- [Krizhevsky and Hinton, 2009] Alex Krizhevsky and Geoffrey Hinton. Learning multiple layers of features from tiny images. *Computer Science*, 2009.
- [Li *et al.*, 2016] Hao Li, Asim Kadav, Igor Durdanovic, Hanan Samet, and Hans Peter Graf. Pruning filters for efficient convnets. *arXiv preprint arXiv:1608.08710*, 2016.
- [Li *et al.*, 2017] Hao Li, Soham De, Zheng Xu, Christoph Studer, Hanan Samet, and Tom Goldstein. Training quantized nets: A deeper understanding. In *NIPS*, 2017.
- [Liu *et al.*, 2017] Zhuang Liu, Jianguo Li, Zhiqiang Shen, Gao Huang, Shoumeng Yan, and Changshui Zhang. Learning efficient convolutional networks through network slimming. In *ICCV*, 2017.
- [Luo *et al.*, 2017] Jian-Hao Luo, Jianxin Wu, and Weiyao Lin. Thinet: A filter level pruning method for deep neural network compression. *ICCV*, 2017.
- [Luo *et al.*, 2018] Jian-Hao Luo, Hao Zhang, Hong-Yu Zhou, Chen-Wei Xie, Jianxin Wu, and Weiyao Lin. Thinet: pruning cnn filters for a thinner net. *TPAMI*, 2018.
- [Ma *et al.*, 2018] Ningning Ma, Xiangyu Zhang, Hai-Tao Zheng, and Jian Sun. Shufflenet v2: Practical guidelines for efficient cnn architecture design. *arXiv preprint arXiv:1807.11164*, 2018.
- [Mirsky, 1960] Leon Mirsky. Symmetric gauge functions and unitarily invariant norms. *The quarterly journal of mathematics*, 11(1):50–59, 1960.
- [Rastegari *et al.*, 2016] Mohammad Rastegari, Vicente Ordonez, Joseph Redmon, and Ali Farhadi. Xnor-net: Imagenet classification using binary convolutional neural networks. In *ECCV*, 2016.
- [Sandler *et al.*, 2018] Mark Sandler, Andrew Howard, Menglong Zhu, Andrey Zhmoginov, and Liang-Chieh Chen. Mobilenetv2: Inverted residuals and linear bottlenecks. In *CVPR*, June 2018.
- [Stewart, 1990] Gilbert W Stewart. Matrix perturbation theory. 1990.
- [Watson, 1992] G Alistair Watson. Characterization of the subdifferential of some matrix norms. *Linear algebra and its applications*, 170:33–45, 1992.
- [Wen *et al.*, 2016] Wei Wen, Chunpeng Wu, Yandan Wang, Yiran Chen, and Hai Li. Learning structured sparsity in deep neural networks. In *NIPS*, 2016.
- [Wen *et al.*, 2017] Wei Wen, Cong Xu, Chunpeng Wu, Yandan Wang, Yiran Chen, and Hai Li. Coordinating filters for faster deep neural networks. In *ICCV*, 2017.
- [Xu *et al.*, 2018] Yuhui Xu, Yongzhuang Wang, Aojun Zhou, Weiyao Lin, and Hongkai Xiong. Deep neural network compression with single and multiple level quantization. *CoRR*, abs/1803.03289, 2018.
- [Xu *et al.*, 2019] Yuhui Xu, Yuxi Li, Shuai Zhang, Wei Wen, Botao Wang, Yingyong Qi, Yiran Chen, Weiyao Lin, and Hongkai Xiong. Trained rank pruning for efficient deep neural networks. In *NIPS EMC2 workshop*, 2019.
- [Zhang *et al.*, 2016] Xiangyu Zhang, Jianhua Zou, Kaiming He, and Jian Sun. Accelerating very deep convolutional networks for classification and detection. *TPAMI*, 38(10):1943–1955, 2016.
- [Zhou *et al.*, 2017] Aojun Zhou, Anbang Yao, Yiwen Guo, Lin Xu, and Yurong Chen. Incremental network quantization: Towards lossless cnns with low-precision weights. *arXiv preprint arXiv:1702.03044*, 2017.
- [Zhou *et al.*, 2019] Yuefu Zhou, Ya Zhang, Yanfeng Wang, and Qi Tian. Accelerate cnn via recursive bayesian pruning. In *ICCV*, 2019.