

# CP-decomposition with Tensor Power Method for Convolutional Neural Networks Compression

Marcella Astrid

University of Science and Technology  
Daejeon, South Korea  
Email: marcella.astrid@ust.ac.kr

Seung-Ik Lee

Electronics and Telecommunications Research Institute  
University of Science and Technology  
Daejeon, South Korea  
Email: the\_silee@etri.re.kr

**Abstract**—Convolutional Neural Networks (CNNs) has shown a great success in many areas including complex image classification tasks. However, they need a lot of memory and computational cost, which hinders them from running in relatively low-end smart devices such as smart phones. We propose a CNN compression method based on CP-decomposition and Tensor Power Method. We also propose an iterative fine tuning, with which we fine-tune the whole network after decomposing each layer, but before decomposing the next layer. Significant reduction in memory and computation cost is achieved compared to state-of-the-art previous work with no more accuracy loss.

## I. INTRODUCTION

Convolutional neural networks (CNNs) have shown notable results in image recognition: VGG [10] and GoogleNet [11] achieved around 90% accuracy for top-5 classification in ImageNet2012 dataset; AlexNet [7] also achieved around 80% top-5 accuracy with the same dataset.

On the other hand, there is an emerging need for embedding or executing CNNs on smart devices such as mobile phones, robots, and embedded devices, in order to give them more intelligence. But the barrier here is that CNNs will require high amounts of memory and computational resources, which are unfortunately hard to be met by small-sized smart devices.

In order to tackle this problem, several approaches recently have been proposed based on tensor decomposition, including Tucker decomposition [5] and Canonical Polyadic (CP) decomposition [2], [8]. Kim et al. [5] successfully decompose all the layers by using Tucker decomposition. However, Tucker decomposition does not seem to compress as much as CP-decomposition due to the core tensors. Moreover, CP-decomposition [2], [8] has not been successful in compressing the whole convolution layers of a CNN because of the CP instability issue [8].

In this work, we further investigate low-rank CP-decomposition to compress the whole convolution layers of a CNN. Our method, called CP-TPM, is based on low-rank CP-decomposition with Tensor Power Method (TPM) for efficient optimization. We also propose iterative fine-tuning to overcome the CP-decomposition instability. We expect the followings with our method:

- **The whole convolution layer decomposition with CP:** To the best of our knowledge, the whole convolution layer decomposition has not been successful because of CP-decomposition's instability [8]. We expect that CP-TPM

can decompose the whole convolution layers in contrast to previous CP-decomposition approaches [2], [8].

- **Overcoming CP-decomposition instability by iterative fine-tuning:** The instability of CP-decomposition, in our view, is the cause of ill-training, such that the loss does not decrease or even amplified when all of the layers are decomposed by CP and fine-tuned once at the final stage. We believe that CP-decomposition on the whole convolution layers without fine-tuning in between causes the loss to be accumulated and become unrecoverable by the once-and-for-all fine-tuning after decomposition. We empirically prove in the experiment that this instability can be overcome by iterative fine-tuning.

Our CP-TPM achieves much reduction in memory and computational cost without a very small accuracy drop.

## II. METHOD

Our approach has two main steps: decomposition and fine-tuning. We apply the two steps layer-by-layer until all the layers of a CNN are decomposed and fine-tuned. All the layers except the fully connected layers are decomposed by CP decomposition, while the fully connected layers are decomposed by SVD. After each decomposition, the whole network is fine-tuned by back propagation.

Before starting the details, these are notations that we use in this paper. Tensors will be notated in calligraphy font capital letters, e.g.  $\mathcal{X}$ . Matrices will be notated in bold capital letters, e.g.  $\mathbf{X}$ . Vectors will be notated in bold small letters, e.g.  $\mathbf{x}$ . Scalars will be notated in regular font small letters, e.g.  $x$ . Regular font capital letters, e.g.  $X$ , will be used for dimension size.

### A. Kernel Tensor Decomposition

CP decomposes a tensor as a linear combination of rank one tensors (1). The number of components is the tensor rank  $R$ . Each component is an outer product of  $n$  vectors, where  $n$  corresponds to the number of ways of the target tensor. Rank  $R$  determines the amount of weight reduction, i.e., the smaller the  $R$  is, the more reduced the weights in a convolution layer.

$$\mathcal{X} = \sum_{r=1}^R \mathbf{a}_r \otimes \mathbf{b}_r \otimes \mathbf{c}_r \quad (1)$$

**Convolution kernel tensor:** In general, convolution layers in CNNs map a 3-way input tensor  $\mathcal{X}$  of size  $S \times W \times H$  into a 3-way output tensor  $\mathcal{Y}$  of size  $S \times W' \times H'$  using a 4-way kernel tensor  $\mathcal{K}$  of size  $T \times S \times D \times D$  with  $T$  corresponding to different output channels,  $S$  corresponding to different input channels, and the last two dimensions corresponding to the spatial dimension (for simplicity, we assume square shaped kernels and odd  $D$ )

$$\mathcal{Y}_{t,w',h'} = \sum_{s=1}^S \sum_{j=1}^D \sum_{i=1}^D \mathcal{K}_{t,s,j,i} \mathcal{X}_{s,w_j,h_i} \quad (2)$$

$$w_j = (w' - 1) \Delta + j - p \text{ and } h_i = (h' - 1) \Delta + i - p,$$

where  $\Delta$  is stride and  $p$  is zero-padding size.

**CP decomposition:** Now the problem is to approximate the kernel tensor  $\mathcal{K}$  with rank- $R$  CP-decomposition. This can be represented as in (3). Spatial dimensions are not decomposed as they are relatively small (e.g.,  $3 \times 3$  or  $5 \times 5$ ).

$$\mathcal{K}_{t,s,j,i} = \sum_{r=1}^R \mathbf{U}_{r,s}^{(1)} \mathcal{U}_{r,j,i}^{(2)} \mathbf{U}_{t,r}^{(3)} \quad (3)$$

where  $\mathbf{U}_{r,s}^{(1)}$ ,  $\mathcal{U}_{r,j,i}^{(2)}$ , and  $\mathbf{U}_{t,r}^{(3)}$  are the three components of sizes  $R \times S$ ,  $R \times D \times D$ , and  $T \times R$ , respectively.

Substituting (3) into (2) and performing simple manipulations gives (4) for the approximate evaluation of the convolution (2) from the input tensor  $\mathcal{X}$  into the output tensor  $\mathcal{Y}$ .

$$\mathcal{Y}_{t,w',h'} = \sum_{r=1}^R \mathbf{U}_{t,r}^{(3)} \left( \sum_{j=1}^D \sum_{i=1}^D \mathcal{U}_{r,j,i}^{(2)} \left( \sum_{s=1}^S \mathbf{U}_{r,s}^{(1)} \mathcal{X}_{s,w_j,h_i} \right) \right) \quad (4)$$

Equation (4) tells us that the output tensor  $\mathcal{Y}$  is computed by a sequence of three separate convolution operations from the input tensor  $\mathcal{X}$  with smaller kernels (Fig. 1(b)):

$$\mathcal{Z}_{r,w,h} = \sum_{s=1}^S \mathbf{U}_{r,s}^{(1)} \mathcal{X}_{s,w,h} \quad (5)$$

$$\mathcal{Z}'_{r,w',h'} = \sum_{j=1}^D \sum_{i=1}^D \mathcal{U}_{r,j,i}^{(2)} \mathcal{Z}_{t,w_j,h_i} \quad (6)$$

$$\mathcal{Y}_{t,w',h'} = \sum_{r=1}^R \mathbf{U}_{t,r}^{(3)} \mathcal{Z}'_{r,w',h'} \quad (7)$$

where  $\mathcal{Z}_{r,w,h}$  and  $\mathcal{Z}'_{r,w',h'}$  are intermediate tensors of sizes  $R \times W \times H$  and  $R \times W' \times H'$ , respectively.

**Fully Connected Layers (FC):** Fully connected layer calculation has the form in (8):

$$\mathbf{y}^T = \mathbf{x}^T \mathbf{W} \quad (8)$$

where  $\mathbf{y}^T$  is the transpose vector of the output of size  $M$ ,  $\mathbf{x}^T$  is the transpose vector of the input of size  $N$ , and  $\mathbf{W}$  is a weight matrix of size  $M \times N$ .

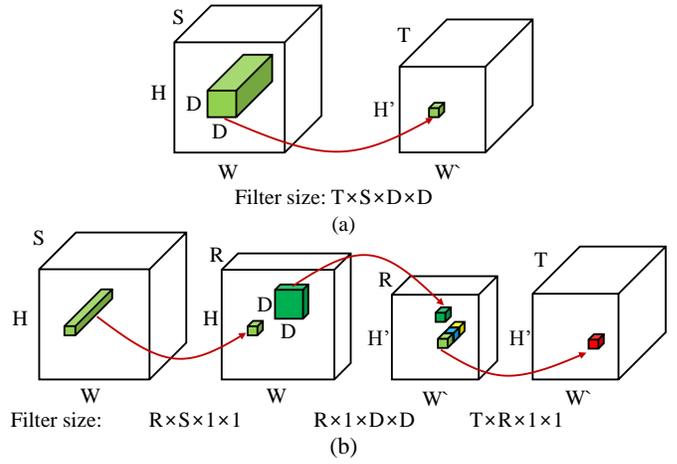


Fig. 1. Convolution layer. (a) Original convolution layer. (b) CP decomposed convolution layer in this paper.

**Decomposition of FC:** As the weights are in a matrix form, we apply Singular Value Decomposition (SVD) to decompose the weight matrix as in (9):

$$\mathbf{W} = \mathbf{U} \mathbf{D} \mathbf{V}^T = (\mathbf{U} \mathbf{D}) \mathbf{V}^T \quad (9)$$

where  $\mathbf{U}$  and  $\mathbf{V}^T$  are left and right singular matrices of sizes  $M \times R$  and  $R \times N$ , respectively, and  $\mathbf{D}$  is a  $R \times R$  diagonal singular-value matrix.

Substituting (9) into (8) and performing grouping gives (10) for approximate evaluation of a fully connected layer with smaller weight matrices.

$$\mathbf{y}^T = (\mathbf{x}^T (\mathbf{U} \mathbf{D})) \mathbf{V}^T \quad (10)$$

Therefore, one fully connected layer can be represented as two fully connected layers as in (11) and (12):

$$\mathbf{z}^T = \mathbf{x}^T (\mathbf{U} \mathbf{D}) \quad (11)$$

$$\mathbf{y}^T = \mathbf{z}^T \mathbf{V}^T \quad (12)$$

where  $\mathbf{z}^T$  is an intermediate layer of size  $R$ .

## B. Complexity Analysis

The initial convolution operation in (2) requires  $TSD^2$  parameters and  $TSD^2W'H'$  multiplication operation. With CP decomposition, the compression ratio  $E$  and speed-up ratio  $C$  are given by:

$$E = \frac{TSD^2}{RS + RD^2 + TR} \quad (13)$$

$$C = \frac{TSD^2W'H'}{RSWH + RD^2W'H' + TRW'H'} \quad (14)$$

The initial operations in FC of (8) is defined by  $MN$  parameters and requires the same number of *multiplication-addition* operations. Therefore, the compression ratio  $E$  and speed-up ratio  $C$  are the same and given by:

$$E = C = \frac{MN}{MR + RN} \quad (15)$$

### C. Rank Selection

Ranks play a key role in CP decomposition. If the rank is too high, compression would not be maximized, and if it is too low, the accuracy would drop too much to be recovered by fine-tuning. However, there is no straight algorithm to find the optimal tensor rank [6]. In fact, determining the rank is NP-hard [3].

Thus, we apply a primitive principle in determining the rank: the higher the accuracy loss caused by a layer, the higher rank the layer needs. Rank proportion is the proportion of rank of a layer among other layers. In order to figure out how sensitive a layer is to decomposition, we perform kind of prior decomposition of each layer with a very low, but constant rank (e.g., 5), and then fine-tune the whole network (one epoch).

### D. Computation of Tensor Decomposition

In general, tensor decomposition is an optimization problem, i.e., minimizing the difference between the decomposed tensor and the target tensor. We employ Tensor Power Method (TPM) [1]. TPM is known to explain the same variance with less rank compared to ALS [1] because the rank-1 tensors found in the early steps of the process explains most of the variances in the target tensor.

TPM approximates a target tensor  $\mathcal{W}$  by adding rank-1 tensors iteratively. First, TPM finds a rank-1 tensor,  $\mathcal{W}_{decomposed}$ , to approximate  $\mathcal{W}$  by minimizing  $\|\mathcal{W} - \mathcal{W}_{decomposed}\|_2$  in a coordinate descent manner. The main idea in the decomposition is that it utilizes the residual  $\mathcal{W}_{residual} = \mathcal{W} - \mathcal{W}_{decomposed}$ , so that the next iteration approximates the residual tensor  $\mathcal{W}_{residual}$  by minimizing  $\|\mathcal{W}_{residual} - \mathcal{W}_{decomposed}\|_2$ . This continues until the number of rank-1 tensors found is equal to  $R$ . More details can be found in [1].

### E. Fine-Tuning

As the accuracy will usually drop after decomposition because of the error in decomposition, fine-tuning is needed to recover the accuracy drop. However, as Lebedev et al. [8] pointed out, CP decomposition has not yet successfully applied to the whole convolution layers of a CNN with one-time fine-tuning because of its instability [5], [8].

To overcome the instability, we iteratively fine-tune the whole network after decomposing each layer in order to prevent the errors from getting too big to recover. In the iterative fine-tuning, no layer is frozen because freezing some layers makes the approach greedy, which usually tends to stuck in local minima. As experimented in [12], letting the layers unfrozen shows better results compared to freezing. In this way, all the layers including the already decomposed can adjust to the newly decomposed layer.

## III. EXPERIMENTS

In this section, we test our approach on AlexNet [7], one of the representative CNNs using Caffe framework [4]. Before describing the main experiments to all of layers, we first briefly introduce AlexNet.

### A. AlexNet Overview

AlexNet is one of famous object recognition architectures and its pre-trained model is available online in Caffe model zoo [4]. As a baseline, we evaluated the accuracy of the pre-trained model using 50,000 validation images from the ImageNet2012 [9] dataset for 1,000 class classification. Top-1 accuracy is 56.83% and top-5 accuracy is 79.95%. AlexNet has eight layers in total consisting of five convolution layers and three fully connected layers.

### B. Whole Network Decomposition

In this section, we explain the results of decomposing the whole network with CP-TPM and iterative fine tuning. As mentioned before, to the best of our knowledge, CP-based decomposition has not yet been successful to the whole convolution layer decomposition because of its instability causing the network to be ill-trained. To overcome the problem, we apply fine tuning iteratively so that the errors from each decomposition are not amplified, which usually is with normal one-shot fine tuning performed after the whole convolution layer decomposition.

Fig. 2 shows the steps of our method in AlexNet. Decomposition and fine tuning are performed for each layer. This process iterates from Conv1 to FC8 in sequence.

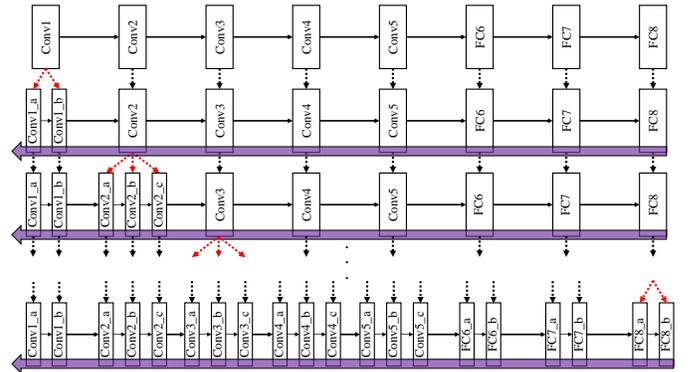


Fig. 2. Black solid arrow shows connection between layer. Red dotted line shows the decomposition process while black dotted line means that the weights are taken from the previous iteration. Purple block arrow means fine-tuning by backpropagation to all layers. First, Conv1 is decomposed into 2 layers (red dotted line) while the others remain the same (black dotted line). Then, fine-tuning is performed to all the layers. Afterward, the next layer, Conv2, is decomposed into 3 layers while the others are same. The process repeats until all the layers are decomposed and fine-tuned.

In order to decompose a layer, we first need to figure out the rank. It is desirable for a layer to have as low rank as possible assuming the same level of accuracy. That means each layer should have ranks proportional to its sensitivity, which is defined as the ratio of  $loss/total\_loss$ .

TABLE I

SETTINGS AND RESULTS COMPARISON WITH ORIGINAL NETWORK AND TUCKER METHOD OF KIM ET AL. [5]. LR: LEARNING RATE. STEPSIZE: FREQUENCY (IN NUMBER OF EPOCH) OF DECREASING LEARNING RATE BY FACTOR OF 10. OTH.: OTHER LAYERS. 1-2: CONV1 AND CONV2.

|           | Rank |     | Epoch              | Batch |      | Learning rate |       | Stepsize | Top-1 Acc        | Top-5 Acc        | Weights                    | Comp. Cost                |
|-----------|------|-----|--------------------|-------|------|---------------|-------|----------|------------------|------------------|----------------------------|---------------------------|
|           | Conv | FC  |                    | 1-2   | Oth. | 1-2           | Oth.  |          |                  |                  |                            |                           |
| Original  | -    | -   | -                  | -     | -    | -             | -     | -        | 56.83            | 79.95            | 61.0M                      | 724M                      |
| Tucker[5] | VBMF |     | 15<br>(one-shot)   | 128   |      | 0.001         |       | 5        | no data          | 78.33<br>(-1.62) | 11.2M<br>( $\times 5.46$ ) | 272M<br>( $\times 2.67$ ) |
| This work | 150  | 300 | 15<br>(each layer) | 128   | 32   | 0.001         | 0.002 | 5        | 54.98<br>(-1.84) | 78.53<br>(-1.42) | 8.7M<br>( $\times 6.98$ )  | 205M<br>( $\times 3.53$ ) |

In order to calculate sensitivity, we first give the same but very small rank (e.g., 5) to all layers and decompose the layers with that rank, which allows us to get the top-5 accuracy loss as in TABLE II. For example, the *total\_loss* for Conv layers is 57.99, and each convolutional layers is assigned with a rank proportional to its sensitivity value from a total of 750 ranks. The same method is also applied to the fully connected layers and in this case we use a total of 900 ranks, which is higher than for the convolutional layers because the FC layers will certainly have much more effects on accuracy.

TABLE II  
RANK CALCULATED FROM THE TOP-5 ACCURACY LOSS.

|                   | Layer        | Top-5 Acc | Top-5 Loss | Rank         |
|-------------------|--------------|-----------|------------|--------------|
| Conv<br>(Ave=150) | Conv1        | 74.57     | 5.38       | 69           |
|                   | Conv2        | 68.06     | 11.89      | 154          |
|                   | Conv3        | 68.09     | 11.86      | 153          |
|                   | Conv4        | 66.27     | 13.68      | 178          |
|                   | Conv5        | 64.78     | 15.17      | 196          |
|                   | <b>Total</b> |           |            | <b>57.99</b> |
| FC<br>(Ave=300)   | FC6          | 51.36     | 28.59      | 365          |
|                   | FC7          | 58.45     | 21.50      | 275          |
|                   | FC8          | 59.64     | 20.31      | 260          |
|                   | <b>Total</b> |           |            | <b>70.40</b> |

As seen in Table I, We achieved a better compression results in both the number of weights and computation cost compared with [5] while maintaining roughly the same level of accuracy. Specifically, we achieved 1.42% accuracy loss which is less than Tucker accuracy loss of 1.62%, with better compression rate. Our method achieved  $\times 6.98$  parameter reduction and  $\times 3.53$  speeding-up, which is better than Tucker-based method that shows  $\times 5.46$  and  $\times 2.67$  respectively.

#### IV. CONCLUSION

We have demonstrated that a TPM-based low-rank CP-decomposition combined with the iterative fine-tuning can achieve a whole network decomposition, which has not been tried before with CP. This approach outperforms the previous Tucker-based decomposition method of Kim et al. [5] for the whole network decomposition. Specifically, our method achieves  $\times 6.98$  parameter reduction and  $\times 3.53$  speeding-up in

Alexnet, while Tucker-based method shows  $\times 5.46$  and  $\times 2.67$  respectively.

There remains much work to be further researched. Firstly, there can be many variations in the iterative fine-tuning, for example, freezing layers after fine-tuning and unfreezing in the final fine-tuning; or starting fine-tuning from the last layer. Secondly, this work has focused more on convolution layers than fully connected layers, where we just have applied SVD. Thus, finding the best algorithms for fully connected layers still needs to be tried. Finally and the most importantly, figuring out the ranks in a systematic manner rather than arbitrary will be another key issue in the further research.

#### ACKNOWLEDGEMENT

This work was supported by the ICT R&D program of MSIP/IITP. [B0101-15-0551, Technology Development of Virtual Creatures with Digital Emotional DNA of Users].

#### REFERENCES

- [1] G. Allen. Sparse higher-order principal components analysis. In *AISTATS*, pages 27–36, 2012.
- [2] E. L. Denton, W. Zaremba, J. Bruna, Y. LeCun, and R. Fergus. Exploiting linear structure within convolutional networks for efficient evaluation. In *Advances in Neural Information Processing Systems*, pages 1269–1277, 2014.
- [3] C. J. Hillar and L.-H. Lim. Most tensor problems are np-hard. *Journal of the ACM (JACM)*, 60(6):45, 2013.
- [4] Y. Jia, E. Shelhamer, J. Donahue, S. Karayev, J. Long, R. Girshick, S. Guadarrama, and T. Darrell. Caffe: Convolutional architecture for fast feature embedding. *arXiv preprint arXiv:1408.5093*, 2014.
- [5] Y.-D. Kim, E. Park, S. Yoo, T. Choi, L. Yang, and D. Shin. Compression of deep convolutional neural networks for fast and low power mobile applications. *arXiv preprint arXiv:1511.06530*, 2015.
- [6] T. G. Kolda and B. W. Bader. Tensor decompositions and applications. *SIAM review*, 51(3):455–500, 2009.
- [7] A. Krizhevsky, I. Sutskever, and G. E. Hinton. Imagenet classification with deep convolutional neural networks. In *Advances in neural information processing systems*, pages 1097–1105, 2012.
- [8] V. Lebedev, Y. Ganin, M. Rakhuba, I. Oseledets, and V. Lempit-sky. Speeding-up convolutional neural networks using fine-tuned cp-decomposition. *arXiv preprint arXiv:1412.6553*, 2014.
- [9] O. Russakovsky, J. Deng, H. Su, J. Krause, S. Satheesh, S. Ma, Z. Huang, A. Karpathy, A. Khosla, M. Bernstein, A. C. Berg, and L. Fei-Fei. ImageNet Large Scale Visual Recognition Challenge. *International Journal of Computer Vision (IJCV)*, 115(3):211–252, 2015.
- [10] K. Simonyan and A. Zisserman. Very deep convolutional networks for large-scale image recognition. *arXiv preprint arXiv:1409.1556*, 2014.
- [11] C. Szegedy, W. Liu, Y. Jia, P. Sermanet, S. Reed, D. Anguelov, D. Erhan, V. Vanhoucke, and A. Rabinovich. Going deeper with convolutions. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, pages 1–9, 2015.
- [12] J. Yosinski, J. Clune, Y. Bengio, and H. Lipson. How transferable are features in deep neural networks? In *Advances in neural information processing systems*, pages 3320–3328, 2014.