

---

# Highly Parallel Deep Ensemble Learning

---

Anonymous Author(s)

Affiliation

Address

email

## Abstract

1 In this paper, we propose a novel highly parallel deep ensemble learning, which  
2 leads to highly compact and parallel deep neural networks. The main idea is to first  
3 represent the data in tensor form, apply a linear transform along certain dimension  
4 and split the transformed data into different independent spectral data sets; then  
5 the matrix product in conventional neural networks is replaced by tensor product,  
6 which in effect imposes certain transformed-induced structure on the original  
7 weight matrices, e.g., a block-circulant structure. The key feature of the proposed  
8 spectral tensor network is that it consists of parallel branches with each branch being  
9 an independent neural network trained using one spectral subset of the training  
10 data. Besides, the joint data/model parallel amiable for GPU implementation.  
11 The outputs of the parallel branches, which are trained on different independent  
12 spectral, are combined for ensemble learning to produce an overall network with  
13 substantially stronger generalization capability than that of those parallel branches.  
14 Moreover, benefiting from the reducing size of inputs, the proposed spectral tensor  
15 network exhibits an inherent network compression, and as a result, reduction  
16 in computation complexity, which leads to the acceleration of training process.  
17 The high parallelism from the massive independent operations of the parallel  
18 spectral subnetworks enable a further acceleration in training and inference process.  
19 We evaluate the proposed spectral tensor networks on the MNIST, CIFAR-10  
20 and ImageNet data sets, to highlight that they simultaneously achieve network  
21 compression, reduction in computation and parallel speedup.

## 22 1 Introduction

23 Deep neural networks (DNNs) [1] have made impressive successes in many applications, such  
24 as computer vision [2][3][4], online game [5][6][7], natural language processing [8][9][10], au-  
25 tonomous driving [11][12][13], and robotics [14][15][16]. However, DNNs are memory-intensive  
26 and computation-intensive, which are two major challenges for wider adoption, e.g., in Internet of  
27 Things (IoT) applications [17]. Modern DNNs may have billions of parameters that consume exces-  
28 sive amount of memory and usually require long training time. For example, AlexNet [18] consists  
29 of three fully-connected layers and five convolutional layers, containing 60 million parameters and  
30 consuming about 250 MB of memory and about 40 hours for training on ImageNet data set [19].

31 To mitigate the impact of memory and computation intensive, there are two types of techniques  
32 most investigated by researchers, namely model compression and parallel training. The model  
33 compression methods, which represents the conventional neural layers in DNNs by compact layers,  
34 such as the strucred linear layers, tensor layers, *etc.*, are one of the most efficient methods to reduce  
35 memory consumption of DNNs. Besides, the proposed parallel training methods can be generally  
36 categorized into data parallelism and model parallelism methods. It should be noted that tensor  
37 layer provides potential parallelism for efficient GPU computing as well as the reduction of the

Table 1: Different tensor layers. An  $N^2 \times N^2$  weight matrix is organized into a  $d$ -th order tensor with the size of each dimension  $n$ , i.e.,  $N^4 = n^d$ , and  $r$  is the tensor rank.

Methods	Model size	Inference time
Fully connected layer	$O(N^4)$	$O(N^4)$
Low-rank matrix [29]	$O(N^2r)$	$O(N^2r)$
CP tensor layer [31]	$O(nr)$	$O(nr)$
Tucker tensor layer [32]	$O(nr + r^3)$	$O(nr)$
HT tensor layer [32]	$O(dn^2r + r^d)$	$O(dnr^2N^2)$
Tensor-train layer [29]	$O(dnr^2)$	$O(dr^3N^2)$

38 memory consumption. Thus, it naturally motivates us to employ tensor layer to achieve both the  
 39 model compression and parallel computing.

40 In this paper, we propose a unified approach that simultaneously achieves both model compression  
 41 and parallel learning *without* communication overhead. The key technique is a novel spectral tensor  
 42 layer that enables a joint *data/model-parallel* implementation of a DNN as follows: 1) The training  
 43 data set is split into multiple orthogonal spectral sets; 2) The neural network is split into parallel  
 44 branches with each branch being a conventional neural network, that are trained asynchronously and  
 45 independently on the corresponding spectral sets; 3) The outputs of the parallel branches are finally  
 46 combined to yield an overall neural network with substantially stronger generalization capability than  
 47 that of those parallel branches.

48 The remainder of this paper is organized as follows. Section 2 presents an brief overview of the  
 49 study on the model compression and parallel computing. Section 3 presents multiple spectral tensor  
 50 layers, including fully connected layers, convolution layers and recurrent layers. Section 4 presents  
 51 the experimental results and we conclude this paper in Section 5.

## 52 2 Related Works

53 **Network compression:** Consider a linear layer that is a central building block of modern DNNs,  
 54 where an input  $\mathbf{x} \in \mathbb{R}^{N^2}$  (e.g., an  $N \times N$  image) is transformed by a weight matrix  $\mathbf{W} \in \mathbb{R}^{N^2 \times N^2}$   
 55 [1], i.e.,  $\mathbf{y} = \mathbf{W}\mathbf{x} \in \mathbb{R}^{N^2}$ . The memory size for storing  $\mathbf{W}$  is  $O(N^4)$  and the computational  
 56 complexity of the matrix-vector product is also  $O(N^4)$ , both are very demanding for smartphones,  
 57 robots, and embedded devices.

58 Many works [20][21][22][23] showed that over 95% of the parameters are redundant, thus the  
 59 network can be greatly compressed. The *structured linear layer* imposes certain structures on  $\mathbf{W}$ ,  
 60 including circulant [21][22][24], circulant-block [25][26], and Teoplitz-like structures [21][27][23].  
 61 For example, for a circulant weight matrix  $\mathbf{W}$  [21][22][27][24], the memory size becomes  $O(N^2)$   
 62 and the computational complexity becomes  $O(N^2 \log N)$ , at the expense of a slight drop in the  
 63 inference performance.

64 On the other hand, the *tensor layer* exploits different types of low-rank tensor representations of  
 65 the weight matrix  $\mathbf{W}$  [28][29][30][31][32]. For example, if the weight matrix  $\mathbf{W} \in \mathbb{R}^{N^2 \times N^2}$  has  
 66 rank  $r \ll N^2$ , such that  $\mathbf{W} = \mathbf{A}\mathbf{B}$ ,  $\mathbf{A} \in \mathbb{R}^{N^2 \times r}$ ,  $\mathbf{B} \in \mathbb{R}^{r \times N^2}$ , then the memory size becomes  
 67  $O(2N^2r)$ , and the computational complexity becomes  $O(2N^2r)$ . Table 1 summarizes the model size  
 68 and computational complexity for various tensor representation schemes.

69 **Parallel training:** Existing works on parallel machine learning take advantage of either data-  
 70 parallelism [33][34][35] or model-parallelism [36][37][38] or both [39]. In the data-parallel approach,  
 71 the data are distributed among multiple processors that apply the same model [33]. The processors  
 72 periodically exchange the outputs and gradients. In the model-parallel approach, exact copies of the  
 73 whole data set are processed by multiple processors that operate in parallel on different parts of the  
 74 same model [36]. There are frequent aggregation of model parts and distribution of gradients [40].

75 From the above discussion, we see that the computation associated with a DNN can be speed up via  
 76 two separate ways: one is through model compression so that reduced number of weight parameters  
 77 leads to reduced number of multiplications and additions; and the other is through parallel computing,

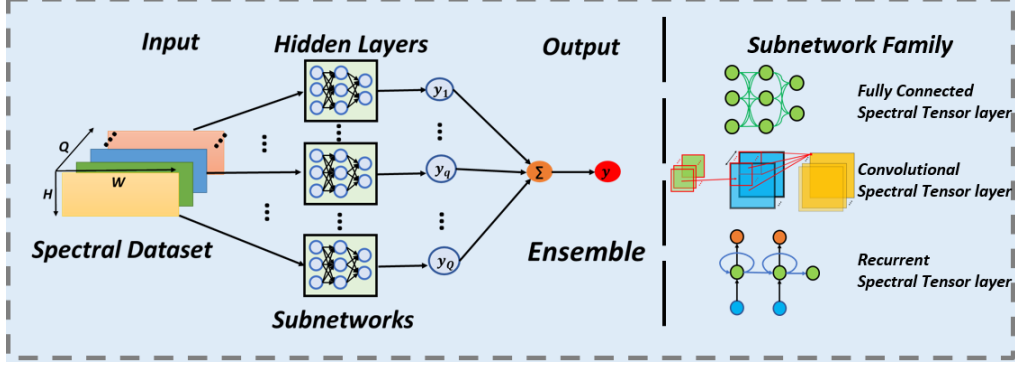


Figure 1: The framework of highly parallel deep ensemble learning.

78 where the speedup is due to distributing the computation among parallel processors, at the expense of  
 79 communication overhead.

### 80 3 Highly Parallel Spectral Tensor Networks

#### 81 3.1 Overview

82 As shown in Fig.1, a batch of data is firstly pre-processed to the spectral dataset and splitted into  
 83  $Q$  independent subsets. Then,  $Q$  different subnetwork, namely spectral tensor layer, works on  $Q$   
 84 independent joint spectral dataset to output own result. Last, the final result is the ensemble result  
 85 from  $Q$  subnetworks. The subnetwork family includes the fully connected spectral tensor layer,  
 86 convolutional spectral tensor layer and recurrent spectral tensor layer. The feature of joint data and  
 87 model parallel makes it suitable for computing on GPUs. The operations in the forward pass, such as  
 88 t-product further provides potential high parallelism for acceleration on GPUs.

#### 89 3.2 Notations and Basic Tensor Operations

90 Scalars, vectors, matrices and tensors are denoted by lowercase, boldface lowercase, boldface capital,  
 91 and calligraphic letters, e.g.,  $a \in \mathbb{R}$ ,  $\mathbf{a} \in \mathbb{R}^n$ ,  $\mathbf{A} \in \mathbb{R}^{n_1 \times n_2}$ ,  $\mathcal{A} \in \mathbb{R}^{n_1 \times n_2 \times n_3}$ , respectively. We use  
 92  $\mathcal{A}(:, :, k)$ ,  $\mathcal{A}(:, j, :)$ ,  $\mathcal{A}(i, :, :)$  to denote the frontal, lateral, and horizontal slices.

93 Given an invertible discrete linear transform  $\mathcal{L} : \mathbb{R}^{n_3} \rightarrow \mathbb{C}^{n_3}$ , let  $\mathcal{L}$  and its inverse  $\mathcal{L}^{-1}$  be taken along  
 94 the third-dimension of third-order tensors. That is, for  $\mathcal{A} \in \mathbb{R}^{n_1 \times n_2 \times n_3}$ ,  $\tilde{\mathcal{A}} = \mathcal{L}(\mathcal{A}) \in \mathbb{C}^{n_1 \times n_2 \times n_3}$ ,  
 95 with  $\tilde{\mathcal{A}}(i, j, :) = \mathcal{L}(\mathcal{A}(i, j, :))$ ,  $i = 1, \dots, n_1, j = 1, \dots, n_2$ . And for  $\tilde{\mathcal{A}} \in \mathbb{C}^{n_1 \times n_2 \times n_3}$ ,  $\mathcal{A} = \mathcal{L}^{-1}(\tilde{\mathcal{A}})$ ,  
 96 with  $\mathcal{A}(i, j, :) = \mathcal{L}^{-1}(\tilde{\mathcal{A}}(i, j, :))$ ,  $i = 1, \dots, n_1, j = 1, \dots, n_2$ . The general spectral tensor product  
 97 [41][42][43] is defined as

$$\mathcal{C} = \mathcal{A} \bullet \mathcal{B} = \mathcal{L}^{-1}(\mathcal{L}(\mathcal{A}) \triangle \mathcal{L}(\mathcal{B})), \quad (1)$$

98 where  $\triangle$  denotes the frontal-slice wise multiplication, i.e., for  $\tilde{\mathcal{A}} \in \mathbb{C}^{n_1 \times n' \times n_3}$ ,  $\tilde{\mathcal{B}} \in \mathbb{C}^{n' \times n_2 \times n_3}$ , if  
 99  $\tilde{\mathcal{C}} = \tilde{\mathcal{A}} \triangle \tilde{\mathcal{B}}$ , then  $\tilde{\mathcal{C}}(:, :, k) = \tilde{\mathcal{A}}(:, :, k) \tilde{\mathcal{B}}(:, :, k)$ ,  $k = 1, \dots, n_3$ . The t-product [44] is a special case  
 100 of (1) where the transform  $\mathcal{L}$  is the DFT [45].

101 For a tensor  $\mathcal{A} \in \mathbb{R}^{n_1 \times n_2 \times n_3}$ , we define the following operations.  $\text{bcirc}(\mathcal{A}) \in \mathbb{R}^{n_1 n_3 \times n_2 n_3}$  organizes  
 102 its  $n_3$  frontal slices into a *block-circulant* matrix

$$\text{bcirc}(\mathcal{A}) = \begin{bmatrix} \mathcal{A}(:, :, 1) & \mathcal{A}(:, :, n_3) & \cdots & \mathcal{A}(:, :, 2) \\ \mathcal{A}(:, :, 2) & \mathcal{A}(:, :, 1) & \cdots & \mathcal{A}(:, :, 3) \\ \vdots & \vdots & \ddots & \vdots \\ \mathcal{A}(:, :, n_3) & \mathcal{A}(:, :, n_3 - 1) & \cdots & \mathcal{A}(:, :, 1) \end{bmatrix}. \quad (2)$$

103 Further,  $\text{unfold}(\cdot)$  is defined as

$$\text{unfold}(\mathcal{A}) = [\mathcal{A}(:, :, 1)^T, \dots, \mathcal{A}(:, :, n_3)^T]^T \in \mathbb{R}^{n_1 n_3 \times n_2},$$

104 and  $\text{fold}(\cdot)$  organizes it back to  $\mathcal{A}$ , such that

$$\text{fold}(\text{unfold}(\mathcal{A})) = \mathcal{A}. \quad (3)$$

105 Given  $\mathcal{A} \in \mathbb{R}^{n_1 \times n' \times n_3}$  and  $\mathcal{B} \in \mathbb{R}^{n' \times n_2 \times n_3}$ , the t-product [45] can be expressed as follows

$$\mathcal{A} *_t \mathcal{B} = \text{fold}(\text{bcirc}(\mathcal{A}) \cdot \text{unfold}(\mathcal{B})) \in \mathbb{R}^{n_1 \times n_2 \times n_3}. \quad (4)$$

106 The  $\text{vec}(\cdot)$  operation maps a matrix in  $\mathbb{R}^{n_1 \times n_2}$  into a vector in  $\mathbb{R}^{n_1 n_2}$ , while  $\text{vec}^{-1}(\cdot)$  is the inverse  
107 mapping.

### 108 3.3 Fully Connected Spectral Tensor Layer

109 A conventional  $N$ -layer fully connected network takes  $m$  input vectors each of size  $\ell'_0 \times 1$  and  
110 represents them as a matrix  $\mathbf{X}^0 \in \mathbb{R}^{\ell'_0 \times m}$ . For example, if each input vector represents a color  
111 image of size  $n \times n \times 3$ , then  $\ell'_0 = 3n^2$ . Each input vector will be classified to one of the  $L$  classes.

112 The network parameters at the  $j$ -th layer include a weight matrix  $\mathbf{W}^j \in \mathbb{R}^{\ell'_j \times \ell'_{j-1}}$  and an offset  
113  $\mathbf{B}^j = \underbrace{[\mathbf{b}^j, \dots, \mathbf{b}^j]}_m \in \mathbb{R}^{\ell'_j \times m}$ , and the forward pass can be represented as [1]

$$\mathbf{X}^j = \sigma(\mathbf{W}^j \cdot \mathbf{X}^{j-1} + \mathbf{B}^j), \quad j = 1, \dots, N-1, \quad (5)$$

114 where  $\mathbf{X}^j \in \mathbb{R}^{\ell'_j \times m}$ , and  $\sigma(\cdot)$  is an element-wise activation function, e.g., linear, sigmoid, ReLU,  
115 and softmax. The last, i.e.,  $N$ -th, layer produces the output  $\mathbf{Y} \in \mathbb{R}^{L \times m}$  corresponding to the  $m$  input  
116 vectors, where

$$\mathbf{X}^N = \mathbf{W}^N \cdot \mathbf{X}^{N-1}, \quad (6)$$

$$\mathbf{Y} = f(\mathbf{X}^N), \quad (7)$$

117 and the output function  $f(\mathbf{X})$  operates on the columns of  $\mathbf{X}$ , i.e.,  $f(\mathbf{X}(:, s))$  maps  $\mathbf{X}(:, s)$  to an  
118 output score vector  $\mathbf{Y}(:, s) \in \mathbb{R}^L$  representing the probabilities that the  $s$ -th input data vector  $\mathbf{X}^0(:, s)$   
119 belongs to different classes. For example  $f(\cdot)$  can be a softmax.

120 For our proposed fully connected spectral tensor network, the  $m$  input data vectors are organized as a  
121 tensor  $\mathcal{X}^0 \in \mathbb{R}^{\ell_0 \times m \times Q}$ . For the  $n \times n \times 3$  color image example, we can set  $Q = 3n$  and  $\ell_0 = n$ ,  
122 namely each image is a lateral slice of  $\mathcal{X}^0$ . Using the weight tensor  $\mathcal{W}^j \in \mathbb{R}^{\ell_j \times \ell_{j-1} \times Q}$  and offset  
123 tensor  $\mathcal{B}^j \in \mathbb{R}^{\ell_j \times m \times Q}$ , a fully connected tensor layer corresponding to (5) and (6) becomes

$$\mathcal{X}^j = \varrho(\mathcal{W}^j \bullet \mathcal{X}^{j-1} + \mathcal{B}^j), \quad j = 1, \dots, N-1, \quad (8)$$

$$\mathcal{X}^N = \mathcal{W}^N \bullet \mathcal{X}^{N-1}, \quad (9)$$

124 where  $\mathcal{X}^j \in \mathbb{R}^{\ell_j \times m \times Q}$ , the spectral tensor product  $\bullet$  is given in (1), and the tensor-activation function  
125  $\varrho(\cdot)$  under transform  $\mathcal{L}$  is defined by applying the conventional element-wise activation function  $\sigma(\cdot)$   
126 in the spectral domain, i.e.,

$$\varrho(\mathcal{X}) = \mathcal{L}^{-1}(\sigma(\mathcal{L}(\mathcal{X}))). \quad (10)$$

127 A salient feature of the proposed spectral tensor network is the *fully parallel implementation*. Specifi-  
128 cally, for  $\mathcal{X} \in \mathbb{R}^{\ell \times m \times Q}$ , denote  $\tilde{\mathcal{X}} \in \mathbb{C}^{\ell \times m \times Q}$  as the transform of  $\mathcal{X}$  along the third dimension, i.e.,  
129  $\tilde{\mathcal{X}}(i, s, :) = \mathcal{L}(\mathcal{X}(i, s, :))$ ,  $i = 1, \dots, \ell$ ,  $s = 1, \dots, m$ . Denote further  $\tilde{\mathcal{X}}_q = \tilde{\mathcal{X}}(:, :, q)$ . Then according  
130 to (1), (8)-(9) can be split into  $Q$  branches of matrix computations

$$\tilde{\mathcal{X}}_q^j = \sigma\left(\tilde{\mathcal{W}}_q^j \cdot \tilde{\mathcal{X}}_q^{j-1} + \tilde{\mathcal{B}}_q^j\right), \quad (11)$$

$$\tilde{\mathcal{X}}_q^N = \tilde{\mathcal{W}}_q^N \cdot \tilde{\mathcal{X}}_q^{N-1}, \quad q = 1, \dots, Q, \quad (12)$$

131 where  $\tilde{\mathcal{X}}_q^j \in \mathbb{C}^{\ell_j \times m}$ ,  $\tilde{\mathcal{B}}_q^j = \underbrace{[\tilde{\mathbf{b}}_q^j, \dots, \tilde{\mathbf{b}}_q^j]}_m \in \mathbb{C}^{\ell_j \times m}$ , and  $\tilde{\mathcal{W}}_q^j \in \mathbb{C}^{\ell_j \times \ell_{j-1}}$ .

132 We further assume that the weight tensor  $\mathcal{W}^j$  in (8)-(9) has low tubal-rank [45][46] such that  $\mathcal{W}^j =$   
133  $\mathcal{C}^j \bullet \mathcal{D}^j$ , where  $\mathcal{C}^j \in \mathbb{C}^{\ell_j \times r \times Q}$ ,  $\mathcal{D}^j \in \mathbb{C}^{r \times \ell_{j-1} \times Q}$ , and  $r \ll \min\{\ell_0, \dots, \ell_N\}$ . Correspondingly, the  
134 weight matrix of each branch has low-rank structure, i.e.,

$$\tilde{\mathcal{W}}_q^j = \tilde{\mathcal{C}}_q^j \cdot \tilde{\mathcal{D}}_q^j, \quad q = 1, \dots, Q, \quad (13)$$

135 where  $\widetilde{\mathbf{C}}_q^j \in \mathbb{C}^{\ell_j \times r}$  and  $\widetilde{\mathbf{D}}_q^j \in \mathbb{C}^{r \times \ell_{j-1}}$ . Then, (11)-(12) become

$$\widetilde{\mathbf{Z}}_q^j = \widetilde{\mathbf{D}}_q^j \cdot \widetilde{\mathbf{X}}_q^{j-1}, \quad j = 1, \dots, N, \quad (14)$$

$$\widetilde{\mathbf{X}}_q^j = \sigma \left( \widetilde{\mathbf{C}}_q^j \cdot \widetilde{\mathbf{Z}}_q^j + \widetilde{\mathbf{B}}_q^j \right), \quad j = 1, \dots, N-1, \quad (15)$$

$$\widetilde{\mathbf{X}}_q^N = \widetilde{\mathbf{C}}_q^N \cdot \widetilde{\mathbf{Z}}_q^N, \quad (16)$$

136 where  $\widetilde{\mathbf{Z}}_q^j \in \mathbb{C}^{r \times m}$ ,  $q = 1, \dots, Q$ .

137 Therefore, an  $N$ -layer fully connected spectral tensor network in (8)-(9) is split into a  $2N$ -layer  
 138 network, such that each layer in (11) is now implemented by two sub-layers, namely a linear layer  
 139 (14) and a nonlinear layer (15), while the  $N$ -th layer in (12) is implemented by two linear sub-layers,  
 140 namely (14) and (16). There are multiple parallel matrix multiplications with the same size and along  
 141 the same dimension in  $Q$  subnetworks. Therefore, we employ the batch matrix multiplication using  
 142 GPUs to accelerate the computations.

143 Finally, we specify the network output. At each branch  $q$ , the output function  $f(\cdot)$  is applied to the  
 144 last layer output  $\widetilde{\mathbf{X}}_q^N$ , as in (7), i.e.,

$$\mathbf{Y}_q = f(\widetilde{\mathbf{X}}_q^N), \quad q = 1, \dots, Q. \quad (17)$$

145 Finally the network output is the weighted sum of the outputs of the  $Q$  branches, i.e.,

$$\mathbf{Y} = \sum_{q=1}^Q \omega_q \mathbf{Y}_q, \quad \text{s.t. } \omega_q \geq 0, \quad \sum_{q=1}^Q \omega_q = 1. \quad (18)$$

146 The loss function can be a cross-entropy function as follows:

$$\text{Loss} = - \sum_{s=1}^m \sum_{c=1}^L \mathbb{1}(\mathbf{y}_s(c) = 1) \cdot \ln(\mathbf{Y}(c, s)). \quad (19)$$

147 Once the spectral tensor network in Fig. 1 is trained, for inference, given a new data sample  $\mathbf{x} \in \mathbb{R}^{\ell_0 Q}$ ,  
 148 we first matricize it into  $\mathbf{X} \in \mathbb{R}^{\ell_0 \times Q}$  and then take transform along each row to obtain  $\widetilde{\mathbf{X}}$ . We input  
 149 the  $q$ -th column of  $\widetilde{\mathbf{X}}$ , i.e.,  $\widetilde{\mathbf{X}}(:, q)$ , to the  $q$ -th sub-network and obtain the output  $\mathbf{y}_q$ ,  $q = 1, \dots, Q$ .  
 150 The final output is then  $\mathbf{y} = \sum_{q=1}^Q \omega_q \mathbf{y}_q$ .

### 151 3.4 Convolutional Spectral Tensor Layer

152 A convolutional neural network [1][47] takes  $m$  input images each of dimension  $H_0 \times W_0 \times C_0$ , i.e.,  
 153 each image has a size  $H_0 \times W_0$  and the number of channels is  $C_0$ , and represents them as a fourth-order  
 154 tensor  $\mathbf{X}^0 \in \mathbb{R}^{H_0 \times W_0 \times C_0 \times m}$ . The input to the  $j$ -th layer is  $\mathbf{X}^{j-1} \in \mathbb{R}^{H_{j-1} \times W_{j-1} \times C_{j-1} \times m}$ , which is  
 155 processed by a convolutional kernel  $\mathbf{W}^j \in \mathbb{R}^{H_j' \times W_j' \times C_{j-1} \times C_j}$  and an offset  $\mathbf{B}^j \in \mathbb{R}^{H_j'' \times W_j'' \times C_j \times m}$ ,  
 156 to yield  $\mathbf{Y}^j \in \mathbb{R}^{H_j'' \times W_j'' \times C_j \times m}$ , with  $H_j'' = H_{j-1} - H_j' + 1$ ,  $W_j'' = W_{j-1} - W_j' + 1$ , where

$$\mathbf{Y}^j(h, w, c, :) = \sum_{d=1}^{C_{j-1}} \sum_{\ell=0}^{W_j'-1} \sum_{i=0}^{H_j'-1} \mathbf{X}^{j-1}(h+i, w+\ell, d, :) \cdot \mathbf{W}^j(i, \ell, d, c) + \mathbf{B}^j(h, w, c, :), \quad (20)$$

$$h = 1, \dots, H_j'', \quad w = 1, \dots, W_j'', \quad c = 1, \dots, C_j.$$

157 To introduce the convolutional spectral tensor network, we represent (20) as a matrix product form  
 158 that is similar to (5)

$$\mathbf{Y}^j = \mathbf{W}^j \cdot \mathbf{X}^{j-1} + \mathbf{B}^j, \quad (21)$$

159 where  $\mathbf{Y}^j \in \mathbb{R}^{H_j'' W_j'' C_j \times m}$ ,  $\mathbf{W}^j \in \mathbb{R}^{H_j'' W_j'' C_j \times H_{j-1} W_{j-1} C_{j-1}}$ , and  $\mathbf{X}^{j-1} \in \mathbb{R}^{H_{j-1} W_{j-1} C_{j-1} \times m}$   
 160 are formed from  $\mathbf{Y}^j$ ,  $\mathbf{W}^j$  and  $\mathbf{X}^{j-1}$ . In particular,

$$\mathbf{Y}^j(:, s) = \text{vec}(\text{unfold}(\mathbf{Y}^j(:, :, :, s))), \quad (22)$$

$$\mathbf{X}^{j-1}(:, s) = \text{vec}(\text{unfold}(\mathbf{X}^{j-1}(:, :, :, s))), \quad s = 1, \dots, m.$$

161 Assume that  $C_j = nB$  for  $j = 0, \dots, N$ , then similar to (8)-(9), (21) leads to a convolutional tensor  
 162 layer

$$\mathcal{Y}^j = \mathcal{W}^j \bullet \mathcal{X}^{j-1} + \mathcal{B}^j, \quad (23)$$

163 where  $\mathcal{Y}^j \in \mathbb{R}^{H_j'' W_j'' n \times m \times B}$ ,  $\mathcal{W}^j \in \mathbb{R}^{H_j'' W_j'' n \times H_{j-1} W_{j-1} n \times B}$ ,  $\mathcal{X}^j \in \mathbb{R}^{H_{j-1} W_{j-1} n \times m \times B}$ , and  
 164  $\mathcal{B}^j \in \mathbb{R}^{H_j'' W_j'' n \times m \times B}$ .

165 Consider the case when  $\mathcal{L}$  is DFT, according to (4), (23) can be written as (21) where  
 166  $\mathbf{Y}^j = \text{unfold}(\mathcal{Y}^j)$ ,  $\mathbf{X}^{j-1} = \text{unfold}(\mathcal{X}^{j-1})$ ,  $\mathbf{B}^j = \text{unfold}(\mathcal{B}^j)$ , and  $\mathbf{W}^j = \text{bcirc}(\mathcal{W}^j) \in$   
 167  $\mathbb{R}^{H_j'' W_j'' n B \times H_{j-1} W_{j-1} n B}$  has a block-circulant structure, namely  $B \times B$  blocks organized in a  
 168 circulant form and each block has size  $H_j'' W_j'' n \times H_{j-1} W_{j-1} n$ . Recall that  $\mathbf{W}^j$  in (21) is derived  
 169 from the convolutional kernel  $\mathbf{W}^j \in \mathbb{R}^{H_j'' \times W_j'' \times n B \times n B}$  in (20), following a linear mapping that is  
 170 consistent with (22). Therefore, the block-circulant structure of  $\mathbf{W}^j$  in (21) implies a block-circulant  
 171 structure of each matrix  $\mathbf{W}^j(i, \ell, :, :)$  in (20).

172 Similar to the fully connected case in Section 3.3, the proposed convolutional spectral tensor network  
 173 also features a *fully parallel implementation*. Specifically, for  $\mathcal{X} \in \mathbb{R}^{H W n \times m \times B}$ , denote  $\tilde{\mathcal{X}} =$   
 174  $\mathcal{L}(\mathcal{X}) \in \mathbb{C}^{H W n \times m \times B}$  as the transform of  $\mathcal{X}$  along the third dimension. Denote  $\tilde{\mathbf{X}}_b = \tilde{\mathcal{X}}(:, :, b)$ .  
 175 Then, (23) can be split into  $B$  parallel branches of matrix computations as follows

$$\tilde{\mathbf{Y}}_b^j = \tilde{\mathbf{W}}_b^j \cdot \tilde{\mathbf{X}}_b^{j-1} + \tilde{\mathbf{B}}_b^j, \quad b = 1, \dots, B, \quad (24)$$

176 where  $\tilde{\mathbf{Y}}_b^j \in \mathbb{C}^{H_j'' W_j'' n \times m}$ ,  $\tilde{\mathbf{W}}_b^j \in \mathbb{C}^{H_j'' W_j'' n \times H_{j-1} W_{j-1} n}$ , and  $\tilde{\mathbf{X}}_b^{j-1} \in \mathbb{C}^{H_{j-1} W_{j-1} n \times m}$ . To fully  
 177 utilize the parallelism of  $B$  branches, we employ batch matrix multiplication on GPUs to accelerate  
 178 (24).

179 We convert (24) back to the convolutional form in (20), using an inverse mapping of (22) as follows

$$\begin{aligned} \tilde{\mathbf{Y}}_b^j(:, :, :, s) &= \text{fold}(\text{vec}^{-1}(\tilde{\mathbf{Y}}_b^j(:, s))), \\ \tilde{\mathbf{X}}_b^{j-1}(:, :, :, s) &= \text{fold}(\text{vec}^{-1}(\tilde{\mathbf{X}}_b^{j-1}(:, s))), \quad s = 1, \dots, m. \end{aligned} \quad (25)$$

180 For the  $b$ -th branch in (24), the input is  $\tilde{\mathbf{X}}_b^{j-1} \in \mathbb{C}^{H_{j-1} \times W_{j-1} \times n \times m}$ , the output feature map is  
 181  $\tilde{\mathbf{Y}}_b^j \in \mathbb{C}^{H_j'' \times W_j'' \times n \times m}$ , and the kernel weight is  $\tilde{\mathbf{W}}_b^j \in \mathbb{C}^{H_j'' \times W_j'' \times n \times n}$ . Then for  $b = 1, \dots, B$ , (24)  
 182 can be rewritten as

$$\begin{aligned} \tilde{\mathbf{Y}}_b^j(h, w, c, :) &= \sum_{d=1}^n \sum_{\ell=0}^{W_j''-1} \sum_{i=0}^{H_j''-1} \tilde{\mathbf{X}}_b^{j-1}(h+i, w+\ell, d, :) \cdot \tilde{\mathbf{W}}_b^j(i, \ell, d, c) + \tilde{\mathbf{B}}_b^j(h, w, c, :), \\ h &= 1, \dots, H_j'', \quad w = 1, \dots, W_j'', \quad c = 1, \dots, n. \end{aligned} \quad (26)$$

183 Similar to the fully connected tensor networks in Section 3.3, we apply the activation function in the  
 184 spectral domain as follows

$$\tilde{\mathbf{Z}}_b^j = \sigma(\tilde{\mathbf{Y}}_b^j) \in \mathbb{C}^{H_j'' \times W_j'' \times n \times m}. \quad (27)$$

185 Then, the pooling operation is performed at the  $j$ -th layer of each branch, resulting in the output  
 186  $\tilde{\mathbf{X}}_b^j \in \mathbb{C}^{H_j \times W_j \times n \times m}$ .

187 At the last layer, the output function  $f(\cdot)$  is applied to  $\tilde{\mathbf{X}}_b^N$  as in (22).

### 188 3.5 Recurrent Spectral Tensor Layer

189 To present the recurrent spectral tensor layer, We take the same case in Section 3.3. Different  
 190 from the fully connected layer, the network parameters of the recurrent layer at the  $j$ -th layer  
 191 include a weight matrix  $\mathbf{W}_X^j \in \mathbb{R}^{\ell_j' \times \ell_{j-1}'}$ ,  $\mathbf{W}_H^j \in \mathbb{R}^{\ell_{j-1}' \times \ell_{j-1}'}$ ,  $\mathbf{W}_Y^j \in \mathbb{R}^{\ell_{j-1}' \times \ell_{j-1}'}$  and an offset  
 192  $\mathbf{B}^j = \underbrace{[\mathbf{b}^j, \dots, \mathbf{b}^j]}_m \in \mathbb{R}^{\ell_j' \times m}$ , and the forward pass can be represented as [1]

$$\begin{aligned} \mathbf{H}_t^j &= \sigma_1(\mathbf{W}_X^j \cdot \mathbf{X}_t^{j-1} + \mathbf{W}_H^j \cdot \mathbf{H}_t^{j-1} + \mathbf{B}^j), \\ \mathbf{X}_t^j &= \sigma_2(\mathbf{W}_Y^j \cdot \mathbf{H}_t^j), \quad j = 1, \dots, N, \quad t = 1, 2, \dots, T, \end{aligned} \quad (28)$$

193 where  $\mathbf{X}_t^j \in \mathbb{R}^{\ell_j^j \times m}$  is the input matrix at the  $j$ -t time step, and  $\mathbf{Y}_t^N = \mathbf{X}_t^N$ .

194 For the introduced recurrent spectral tensor layer, the forward can be computed as

$$\begin{aligned} \mathcal{H}_t^j &= \varrho_1(\mathcal{W}_X^j \bullet \mathcal{X}_t^{j-1} + \mathcal{W}_H^j \bullet \mathcal{H}_t^{j-1} + \mathcal{B}^j), \\ \mathcal{X}_t^j &= \varrho_2(\mathcal{W}_Y^j \bullet \mathcal{H}_t^j), \quad j = 1, \dots, N, t = 1, 2, \dots, T, \end{aligned} \quad (29)$$

195 where  $\mathcal{X}_t^j \in \mathbb{R}^{\ell_j \times m \times Q}$ .

196 Furthermore, the implementation of  $Q$  branches for computations can be written as

$$\begin{aligned} \tilde{H}_{q,t}^j &= \varrho_1(\tilde{W}_X^j \bullet \tilde{X}_{q,t}^{j-1} + \tilde{W}_H^j \bullet \tilde{H}_{q,t}^{j-1} + \tilde{B}^j), \\ \tilde{X}_{q,t}^j &= \varrho_2(\tilde{W}_Y^j \bullet \tilde{H}_{q,t}^j), \quad j = 1, \dots, N, t = 1, 2, \dots, T, \end{aligned} \quad (30)$$

197 where  $\tilde{X}_{q,t}^j \in \mathbb{C}^{\ell_j \times m}$  is the input for the  $q$ -th subnetwork at the  $t$ -th time step, likewise for  $\tilde{H}_{q,t}^j$ .  
 198 There are  $2Q$  parallel matrix multiplications with the same size and along the same dimension, thus  
 199 we use the batch matrix multiplication on GPUs to accelerate (30).

## 200 4 Performance Evaluation

201 We first describe the experimental settings, then present the results on the MNIST, CIFAR-10 and  
 202 ImageNet data sets.

### 203 4.1 Data Sets and Performance Metrics

204 We verify the performance of the proposed spectral tensor networks on the following three widely  
 205 used data sets: 1) MNIST [48] contains grayscale images of handwritten digits. Each image has  
 206  $28 \times 28$  pixels. The training set has 60,000 images and the testing set has 10,000 images. 2)  
 207 CIFAR-10 [49] contains 60,000 color images in 10 classes, where each image has size  $32 \times 32 \times 3.3$   
 208 ImageNet-1K [19]: It contains 12,000,000 training images and 50,000 testing images with size of  
 209  $224 \times 224 \times 3$ , labeled with the presence or absence of 1000 object categories that do not overlap  
 210 with each other.

211 We are interested in the following performance metrics: 1) *Compression ratio*: the ratio of the  
 212 conventional network size to the spectral tensor network size, which is the also the total reduction in  
 213 computation due to the reduced number of non-zero network weights; 2) *Parallel speedup*: the ratio  
 214 of the training time of a conventional network to that of the spectral tensor network, due to the fully  
 215 parallel training of all sub-networks; 3) *Convergence*: the loss value versus the training iterations; 4)  
 216 *Accuracy*: the percentage of correctly estimated labels. Both the training and testing processes are  
 217 executed on a DGX-2 server [50] that has two 64 core AMD CPUs, 8 NVIDIA A100 GPUs and 2  
 218 TB of memory. The operating system is Ubuntu 20.04 with CUDA 10.1. We use PyTorch [51] to  
 219 implement neural networks.

220 We summarize the compression ratio, the reduction of computation, and the parallel speedup in  
 221 Table 2. They are theoretical upper bounds, while their actual values depend on data sets and  
 222 implementations. For the compression ratio and reduction in computation, each fully connected  
 223 network / convolutional network has two columns: the right one corresponds to select-the-best  
 224 weighting, and the left one to other weighting methods.

### 225 4.2 Evaluation of Fully Connected Spectral Tensor Networks

226 For comparison, we consider a conventional fully connected network (FC) [1], the tNN [44], and the  
 227 fully connected spectral tensor network (FC-tensor) in Section 3. All three methods use the ReLU  
 228 activation function as  $\sigma(\cdot)$  in the hidden layers, the softmax function as the output function  $f(\cdot)$  in  
 229 the last layer, and the cross-entropy loss function in (19). We use  $N = 8$  layers in each method and  
 230 the DCT transform in tNN and the proposed FC-tensor method. The learning rate was set to be 0.01,  
 231 the batch size was set to be 64, and we used the Adam optimizer [52].

232 For the MNIST data set, the conventional FC method has  $n = 28$ ,  $\ell_0 = \dots = \ell_7 = 784$ , and  $L = 10$ .  
 233 Both the tNN method and our FC-tensor method have  $n = 28$ ,  $Q = 28$ ,  $\ell_0 = \dots = \ell_7 = 28$ , and  
 234  $L = 10$ . Our FC-tensor method has  $r = 8$ . For the CIFAR-10 data set, the following parameters are

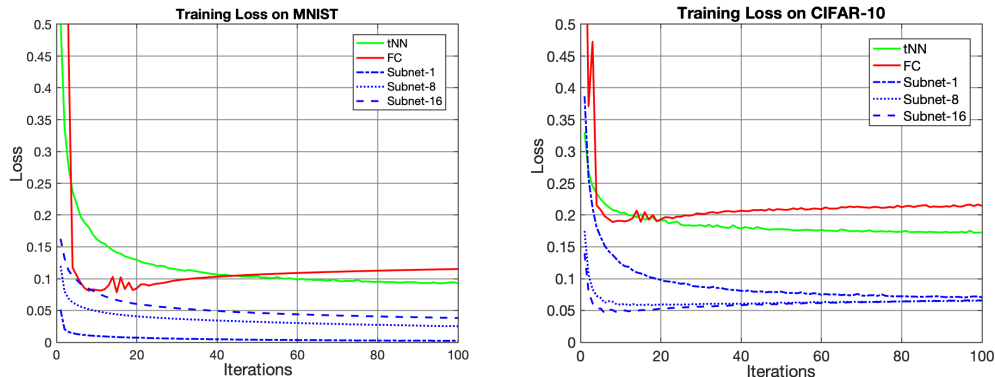


Figure 2: Training loss of fully connected networks on the MNIST data set (left) and CIFAR-10 data set (right).

Table 2: Upper bounds for model compression and parallel speedup. For fully connected networks,  $n$  denotes the input size of each sub-network,  $r$  is a rank value, and there are  $Q$  branches. For convolutional networks, there are  $B$  branches.

-	Fully Connected	Convolutional (1D)
Compression ratio	$O(nQ/2r), O(nQ^2/2r)$	$O(B), O(B^2)$
Reduction in computation	$O(nQ/2r), O(nQ^2/2r)$	$O(B), O(B^2)$
Parallel speedup	$O(Q)$	$O(B)$

235 different:  $n = 32, Q = 32, \ell'_0 = \dots = \ell'_7 = 1,024$ , and  $\ell_0 = \dots = \ell_7 = 32$ . Therefore, our methods  
 236 achieve a compression ratio of  $49\times$  and  $64\times$  for the two data sets, respectively.

237 The training loss over iterations is shown in Fig. 2, with the left one for the MNIST data set and  
 238 the right one for the CIFAR-10 data set. Our scheme converges faster than tNN and FC, while the  
 239 training process is more stable than FC. The possible reason is that the FC-tensor has much less  
 240 parameters so that a more stable model can be learned from the same amount of data samples<sup>1</sup>. The  
 241 loss values of our sub-networks are lower than both tNN and FC.

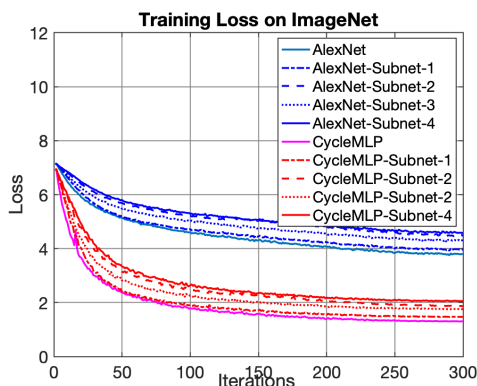


Figure 3: Training loss on ImageNet-1K data set.

Table 3: MNIST and CIFAR-10 data sets.

Methods	MNIST	CIFAR-10
FC [1]	98.71%	<b>59.19%</b>
tNN [44]	97.59%	44.50%
FC-tensor ( <b>average</b> )	97.43%	47.24%
FC-tensor ( <b>weighted sum</b> )	98.02%	48.13%
FC-tensor ( <b>geometric</b> )	<b>99.01%</b>	48.33%

242 In Table 3, we report accuracy results on both MNIST and CIFAR-10 data sets. Among the four  
 243 schemes for weighting the sub-networks, the geometric weighting gives the best performance. For  
 244 the MNIST data set, all three methods achieve a relative high accuracy, i.e., over 97%, while our  
 245 FC-tensor method reaches 98.36%. For the CIFAR-10 data set, all three methods achieve a relative  
 246 low accuracy, i.e., below 60%. This is consistent with the known fact that fully connected layers are

<sup>1</sup>Note that we use the same number of layers and the same batch size.



Table 3: Results on the ImageNet-1K data set.

Methods	Accuracy	Size	Training Time
AlexNet [18]	<b>63.44%</b>	244 MB	40.8 h
AlexNet-spectral ( <b>average</b> )	61.26%	130 MB	31.9 h
AlexNet-spectral ( <b>weighted sum</b> )	58.01%	130 MB	31.9 h
AlexNet-spectral ( <b>geometric</b> )	62.26%	130 MB	31.9 h
AlexNet-spectral ( <b>select-the-best</b> )	56.45%	32.5 MB	31.9 h
CycleMLP [53]	<b>83.23%</b>	103 MB	93.6 h
CycleMLP-spectral ( <b>average</b> )	78.80%	76 MB	60.4 h
CycleMLP-spectral ( <b>weighted sum</b> )	77.54%	76 MB	60.4 h
CycleMLP-spectral ( <b>geometric</b> )	83.20%	76 MB	60.4 h
CycleMLP-spectral ( <b>select-the-best</b> )	72.45%	19 MB	60.4 h

247 not enough for the classification task on CIFAR-10. Note that both tNN and FC-tensor achieve lower  
248 accuracy than the FC method.

### 249 4.3 Evaluation on ImageNet Data set

250 The ImageNet data set [19] is split into  $B = 4$  spectral subsets, where each image is organized into  
251 a tensor of size  $56 \times 56 \times 3 \times 16$  and then processed into a spectrum tensor using DCT transform.  
252 Note that the three RGB channels are processed independently.

253 Our proposed spectral tensor methods have the same structure in Fig. 1, where each branch is replaced  
254 by either AlexNet [18] or CycleMLP. We use the DCT transform in our spectral methods. We set  
255 the learning rate 0.01 and the batch size 128. We follow the standard practice in the community by  
256 reporting the top-1 accuracy on the testing set.

257 For the ImageNet data set, the training loss over training iterations is shown in Fig. 3. Our spectral  
258 sub-networks have similar loss curve to their original networks. In Table 3, we report the accuracy,  
259 model size, and training time. For the AlexNet structure, our spectral network achieves  $1.88\times$  model  
260 compression and  $1.28\times$  speedup in training time, at the cost of an accuracy drop of 1.18%. For the  
261 CycleMLP structure, our spectral network achieves  $1.36\times$  model compression and  $1.55\times$  speedup in  
262 training time, at the cost of an accuracy drop of 0.03%.

## 263 5 Conclusions

264 In this paper, we have proposed a spectral tensor form of deep neural networks that is inherently com-  
265 pressive and allows communication-free parallel/distributed implementations. The data is organized  
266 into tensors and a linear transform is applied along certain dimension, resulting in different spectral  
267 subsets. The overall network consists of parallel branches of networks, each independently performs  
268 training and inference on a spectral data subset. We tested the proposed spectral networks, including  
269 fully connected, convolutional, AlexNet, and CycleMLP, on the MNIST, CIFAR-10 and ImageNet  
270 data sets, and results show that they can achieve relatively high accuracy with substantial network  
271 compression, computation reduction, and parallel speedup, compared with conventional networks.  
272 Limited by the pages, We do not provide experiments using the recurrent spectral tensor layer.

273 For future works, we would like to explore an ensemble-style approach *model soup* [54] that takes  
274 average over multiple trained models and achieves state-of-the-art performance on the ImageNet data  
275 set.

## 276 References

- 277 [1] I. Goodfellow, Y. Bengio, and A. Courville, *Deep learning*, MIT Press, 2016.
- 278 [2] Yang Liu, Peng Sun, Nickolas Wergeles, and Yi Shang, “A survey and performance evaluation  
279 of deep learning methods for small object detection,” *Expert Systems with Applications*, vol.  
280 172, pp. 114602, 2021.

- 281 [3] Xiyang Dai, Yinpeng Chen, Bin Xiao, Dongdong Chen, Mengchen Liu, Lu Yuan, and Lei  
282 Zhang, “Dynamic head: Unifying object detection heads with attentions,” in *Proceedings of the*  
283 *IEEE/CVF Conference on Computer Vision and Pattern Recognition*, 2021, pp. 7373–7382.
- 284 [4] Peize Sun, Rufeng Zhang, Yi Jiang, Tao Kong, Chenfeng Xu, Wei Zhan, Masayoshi Tomizuka,  
285 Lei Li, Zehuan Yuan, Changhu Wang, et al., “Sparse r-cnn: End-to-end object detection with  
286 learnable proposals,” in *Proceedings of the IEEE/CVF Conference on Computer Vision and*  
287 *Pattern Recognition*, 2021, pp. 14454–14463.
- 288 [5] Xiangjun Wang, Junxiao Song, Penghui Qi, Peng Peng, Zhenkun Tang, Wei Zhang, Weimin Li,  
289 Xiongjun Pi, Jujie He, Chao Gao, et al., “Scc: an efficient deep reinforcement learning agent  
290 mastering the game of starcraft ii,” in *International Conference on Machine Learning*. PMLR,  
291 2021, pp. 10905–10915.
- 292 [6] Tianhao Zhang, Yueheng Li, Chen Wang, Guangming Xie, and Zongqing Lu, “Fop: Factorizing  
293 optimal joint policy of maximum-entropy multi-agent reinforcement learning,” in *International*  
294 *Conference on Machine Learning*. PMLR, 2021, pp. 12491–12500.
- 295 [7] Bo Liu, Qiang Liu, Peter Stone, Animesh Garg, Yuke Zhu, and Anima Anandkumar, “Coach-  
296 player multi-agent reinforcement learning for dynamic team composition,” in *International*  
297 *Conference on Machine Learning*. PMLR, 2021, pp. 6860–6870.
- 298 [8] Zhiqi Huang, Fenglin Liu, Xian Wu, Shen Ge, Helin Wang, Wei Fan, and Yuexian Zou, “Audio-  
299 oriented multimodal machine comprehension via dynamic inter-and intra-modality attention,” in  
300 *Proceedings of the AAAI Conference on Artificial Intelligence*, 2021, vol. 35, pp. 13098–13106.
- 301 [9] Tianyang Zhao, Zhao Yan, Yunbo Cao, and Zhoujun Li, “Asking effective and diverse questions:  
302 a machine reading comprehension based framework for joint entity-relation extraction,” in  
303 *Proceedings of the Twenty-Ninth International Conference on International Joint Conferences*  
304 *on Artificial Intelligence*, 2021, pp. 3948–3954.
- 305 [10] Cong Sun, Zhihao Yang, Lei Wang, Yin Zhang, Hongfei Lin, and Jian Wang, “Biomedical  
306 named entity recognition using bert in the machine reading comprehension framework,” *Journal*  
307 *of Biomedical Informatics*, vol. 118, pp. 103799, 2021.
- 308 [11] Aditya Prakash, Kashyap Chitta, and Andreas Geiger, “Multi-modal fusion transformer for  
309 end-to-end autonomous driving,” in *Proceedings of the IEEE/CVF Conference on Computer*  
310 *Vision and Pattern Recognition*, 2021, pp. 7077–7087.
- 311 [12] Yingfeng Cai, Tianyu Luan, Hongbo Gao, Hai Wang, Long Chen, Yicheng Li, Miguel Angel  
312 Sotelo, and Zhixiong Li, “Yolov4-5d: An effective and efficient object detector for autonomous  
313 driving,” *IEEE Transactions on Instrumentation and Measurement*, vol. 70, pp. 1–13, 2021.
- 314 [13] Sudeep Fadadu, Shreyash Pandey, Darshan Hegde, Yi Shi, Fang-Chieh Chou, Nemanja Djuric,  
315 and Carlos Vallespi-Gonzalez, “Multi-view fusion of sensor data for improved perception and  
316 prediction in autonomous driving,” in *Proceedings of the IEEE/CVF Winter Conference on*  
317 *Applications of Computer Vision*, 2022, pp. 2349–2357.
- 318 [14] Huu-Thiet Nguyen, Chien Chern Cheah, and Kar-Ann Toh, “An analytic layer-wise deep  
319 learning framework with applications to robotics,” *Automatica*, vol. 135, pp. 110007, 2022.
- 320 [15] Radouan Ait Mouha et al., “Deep learning for robotics,” *Journal of Data Analysis and*  
321 *Information Processing*, vol. 9, no. 02, pp. 63, 2021.
- 322 [16] Yinong Chen and Gennaro De Luca, “Technologies supporting artificial intelligence and  
323 robotics application development,” *Journal of Artificial Intelligence and Technology*, vol. 1, no.  
324 1, pp. 1–8, 2021.
- 325 [17] G. Menghani, “Efficient deep learning: A survey on making deep learning models smaller,  
326 faster, and better,” *arXiv preprint arXiv:2106.08962*, 2021.
- 327 [18] A. Krizhevsky, I. Sutskever, and G. E. Hinton, “ImageNet classification with deep convolutional  
328 neural networks,” *Advances in Neural Information Processing Systems*, vol. 25, pp. 1097–1105,  
329 2012.

- 330 [19] J. Deng, W. Dong, R. Socher, L.-J. Li, K. Li, and F.-F. Li, “ImageNet: A large-scale hierarchical  
331 image database,” in *IEEE CVPR*. Ieee, 2009, pp. 248–255.
- 332 [20] M. Denil, B. Shakibi, L. Dinh, M. Ranzato, and N. De Freitas, “Predicting parameters in deep  
333 learning,” in *Advances in Neural Information Processing Systems*, 2013, pp. 2148–2156.
- 334 [21] V. Sindhwani, T. Sainath, and S. Kumar, “Structured transforms for small-footprint deep  
335 learning,” in *Advances in Neural Information Processing Systems*, 2015, pp. 3088–3096.
- 336 [22] Y. Cheng, F. Yu, R. S. Feris, S. Kumar, A. Choudhary, and S.-F. Chang, “An exploration of  
337 parameter redundancy in deep networks with circulant projections,” in *IEEE International  
338 Conference on Computer Vision*, 2015, pp. 2857–2865.
- 339 [23] M. Moczulski, M. Denil, J. Appleyard, N. De Freitas, Z. Wang, M. Zoghi, F. Hutter, D. Matheson,  
340 and S. Reed, “ACDC: A structured efficient linear layer,” in *ICLR*, 2015, vol. 55.
- 341 [24] A. Prabhu, A. Farhadi, and M. Rastegari, “Butterfly transform: An efficient fft based neural  
342 architecture design,” in *IEEE/CVF Conference on Computer Vision and Pattern Recognition*,  
343 2020, pp. 12024–12033.
- 344 [25] C. Ding, S. Liao, Y. Wang, Z. Li, N. Liu, Y. Zhuo, C. Wang, X. Qian, Y. Bai, and G. Yuan,  
345 “Circnn: accelerating and compressing deep neural networks using block-circulant weight  
346 matrices,” in *Annual IEEE/ACM International Symposium on Microarchitecture*, 2017, pp.  
347 395–408.
- 348 [26] S. Liao and B. Yuan, “Circconv: A structured convolution with low complexity,” in *AAAI  
349 Conference on Artificial Intelligence*, 2019, vol. 33, pp. 4287–4294.
- 350 [27] B. Gong, B. Jou, F. Yu, and S.-F. Chang, “Tamp: A library for compact deep neural networks  
351 with structured matrices,” in *ACM International Conference on Multimedia*, 2016, pp. 1206–  
352 1209.
- 353 [28] T. N Sainath, B. Kingsbury, V. Sindhwani, E. Arisoy, and B. Ramabhadran, “Low-rank matrix  
354 factorization for deep neural network training with high-dimensional output targets,” in *IEEE  
355 ICASSP*, 2013, pp. 6655–6659.
- 356 [29] A. Novikov, D. Podoprikin, A. Osokin, and D.P. Vetrov, “Tensorizing neural networks,” in  
357 *Advances in Neural Information Processing Systems*, 2015, pp. 442–450.
- 358 [30] W. Wang, Y. Sun, B. Eriksson, W. Wang, and V. Aggarwal, “Wide compression: Tensor ring  
359 nets,” in *IEEE CVPR*, 2018, pp. 9329–9338.
- 360 [31] V. Lebedev, Y. Ganin, M. Rakhuba, I. Oseledets, and V. Lempitsky, “Speeding-up convolutional  
361 neural networks using fine-tuned cp-decomposition,” *ICLR*, 2015.
- 362 [32] M. Yin, S. Liao, X.-Y. Liu, X. Wang, and B. Yuan, “Towards extremely compact rnns for video  
363 recognition with fully decomposed hierarchical tucker structure,” in *IEEE CVPR*, 2021, pp.  
364 12085–12094.
- 365 [33] J. Verbraeken, M. Wolting, J. Katzy, J. Kloppenburg, T. Verbelen, and J. S. Rellermeier, “A  
366 survey on distributed machine learning,” *ACM Computing Surveys (CSUR)*, vol. 53, no. 2, pp.  
367 1–33, 2020.
- 368 [34] Vipul Gupta, Dhruv Choudhary, Peter Tang, Xiaohan Wei, Xing Wang, Yuzhen Huang, Arun  
369 Kejariwal, Kannan Ramchandran, and Michael W Mahoney, “Training recommender systems  
370 at scale: Communication-efficient model and data parallelism,” in *Proceedings of the 27th ACM  
371 SIGKDD Conference on Knowledge Discovery & Data Mining*, 2021, pp. 2928–2936.
- 372 [35] Xiangyu Ye, Zhiquan Lai, Shengwei Li, Lei Cai, Ding Sun, Linbo Qiao, and Dongsheng Li,  
373 “Hippie: A data-paralleled pipeline approach to improve memory-efficiency and scalability for  
374 large dnn training,” in *50th International Conference on Parallel Processing*, 2021, pp. 1–10.
- 375 [36] A. L. Gaunt, M. A. Johnson, M. Riechert, D. Tarlow, R. Tomioka, D. Vytiniotis, and S. Webster,  
376 “AMPNet: Asynchronous model-parallel training for dynamic neural networks,” *arXiv preprint  
377 arXiv:1705.09786*, 2017.

- 378 [37] An Xu, Zhouyuan Huo, and Heng Huang, “On the acceleration of deep learning model  
379 parallelism with staleness,” in *Proceedings of the IEEE/CVF Conference on Computer Vision  
380 and Pattern Recognition*, 2020, pp. 2088–2097.
- 381 [38] Kabir Nagrecha, “Model-parallel model selection for deep learning systems,” in *Proceedings of  
382 the 2021 International Conference on Management of Data*, 2021, pp. 2929–2931.
- 383 [39] E. P. Xing, Q. Ho, P. Xie, and D. Wei, “Strategies and principles of distributed machine learning  
384 on big data,” *Engineering*, vol. 2, no. 2, pp. 179–195, 2016.
- 385 [40] M. Abadi, P. Barham, J. Chen, Z. Chen, A. Davis, J. Dean, M. Devin, S. Ghemawat, G. Irving,  
386 M. Isard, et al., “TensorFlow: A system for large-scale machine learning,” in *USENIX  
387 Symposium on Operating Systems Design and Implementation*, 2016, pp. 265–283.
- 388 [41] E. Kernfeld, M. Kilmer, and S. Aeron, “Tensor–tensor products with invertible linear transforms,”  
389 *Linear Algebra and its Applications*, vol. 485, pp. 545–570, 2015.
- 390 [42] M. E. Kilmer, L. Horesh, H. Avron, and E. Newman, “Tensor-tensor algebra for optimal  
391 representation and compression of multiway data,” *Proceedings of the National Academy of  
392 Sciences*, vol. 118, no. 28, 2021.
- 393 [43] X.-Y. Liu and X. Wang, “Fourth-order tensors with multidimensional discrete transforms,”  
394 *arXiv preprint arXiv:1705.01576*, pp. 1–37, 2017.
- 395 [44] E. Newman, L. Horesh, H. Avron, and M. Kilmer, “Stable tensor neural networks for rapid  
396 deep learning,” *arXiv preprint arXiv:1811.06569*, 2018.
- 397 [45] M.E. Kilmer and C.D Martin, “Factorization strategies for third-order tensors,” *Linear Algebra  
398 and its Applications*, vol. 435, no. 3, pp. 641–658, 2011.
- 399 [46] M.E. Kilmer, K. Braman, N. Hao, and R.C. Hoover, “Third-order tensors as operators on  
400 matrices: A theoretical and computational framework with applications in imaging,” *SIAM  
401 Journal on Matrix Analysis and Applications*, vol. 34, no. 1, pp. 148–172, 2013.
- 402 [47] Z. Li, F. Liu, W. Yang, S. Peng, and J. Zhou, “A survey of convolutional neural networks:  
403 analysis, applications, and prospects,” *IEEE Transactions on Neural Networks and Learning  
404 Systems*, 2021.
- 405 [48] L. Deng, “The MNIST database of handwritten digit images for machine learning research,”  
406 *IEEE Signal Processing Magazine*, vol. 29, no. 6, pp. 141–142, 2012.
- 407 [49] A. Krizhevsky and G. Hinton, “Learning multiple layers of features from tiny images,” *Master’s  
408 thesis, University of Tront*, 2009.
- 409 [50] J. Choquette et al., “NVIDIA A100 tensor core GPU: Performance and innovation,” *IEEE  
410 Micro*, vol. 41, no. 2, pp. 29–35, 2021.
- 411 [51] A. Paszke, S. Gross, F. Massa, A. Lerer, J. Bradbury, and et al, “PyTorch: An imperative  
412 style, high-performance deep learning library,” in *Advances in Neural Information Processing  
413 Systems*, 2019, pp. 8026–8037.
- 414 [52] D.P Kingma and J. Ba, “Adam: A method for stochastic optimization,” *ICLR*, 2015.
- 415 [53] S. Chen, E. Xie, C. Ge, D. Liang, and P. Luo, “CycleMLP: A MLP-like architecture for dense  
416 prediction,” *ICLR*, 2022.
- 417 [54] Mitchell Wortsman, Gabriel Ilharco, and et al, “Model soups: averaging weights of multiple fine-  
418 tuned models improves accuracy without increasing inference time,” *preprint arXiv:2203.05482*,  
419 2022.

420 **Checklist**

- 421 1. For all authors...
- 422 (a) Do the main claims made in the abstract and introduction accurately reflect the paper’s  
423 contributions and scope? [Yes]
- 424 (b) Did you describe the limitations of your work? [No]
- 425 (c) Did you discuss any potential negative societal impacts of your work? [No]
- 426 (d) Have you read the ethics review guidelines and ensured that your paper conforms to  
427 them? [Yes]
- 428 2. If you are including theoretical results...
- 429 (a) Did you state the full set of assumptions of all theoretical results? [N/A]
- 430 (b) Did you include complete proofs of all theoretical results? [N/A]
- 431 3. If you ran experiments...
- 432 (a) Did you include the code, data, and instructions needed to reproduce the main experi-  
433 mental results (either in the supplemental material or as a URL)? [Yes] **See Section**  
434 **4.1.**
- 435 (b) Did you specify all the training details (e.g., data splits, hyperparameters, how they  
436 were chosen)? [Yes] **See Section 4.2 and 4.3.**
- 437 (c) Did you report error bars (e.g., with respect to the random seed after running experi-  
438 ments multiple times)? [Yes] We do not supply experiments of recurrent spectral tensor  
439 layer.
- 440 (d) Did you include the total amount of compute and the type of resources used (e.g., type  
441 of GPUs, internal cluster, or cloud provider)? [Yes] **See Section 4.1.**
- 442 4. If you are using existing assets (e.g., code, data, models) or curating/releasing new assets...
- 443 (a) If your work uses existing assets, did you cite the creators? [Yes]
- 444 (b) Did you mention the license of the assets? [Yes]
- 445 (c) Did you include any new assets either in the supplemental material or as a URL? [Yes]  
446 **See Section 4.3.**
- 447 (d) Did you discuss whether and how consent was obtained from people whose data you’re  
448 using/curating? [Yes] **All the consent is open-source.**
- 449 (e) Did you discuss whether the data you are using/curating contains personally identifiable  
450 information or offensive content? [No] **We do not use any personally identifiable**  
451 **information or offensive content.**
- 452 5. If you used crowdsourcing or conducted research with human subjects...
- 453 (a) Did you include the full text of instructions given to participants and screenshots, if  
454 applicable? [N/A]
- 455 (b) Did you describe any potential participant risks, with links to Institutional Review  
456 Board (IRB) approvals, if applicable? [N/A]
- 457 (c) Did you include the estimated hourly wage paid to participants and the total amount  
458 spent on participant compensation? [N/A]