

---

# THE CASCADED FORWARD ALGORITHM FOR NEURAL NETWORK TRAINING

---

Gongpei Zhao, Tao Wang, Yidong Li, Yi Jin, Congyan Lang

Beijing Jiaotong University

{csgpzha, twang, ydli, yjin, cylang}@bjtu.edu.cn

Haibin Ling

Stony Brook University

hling@cs.stonybrook.edu

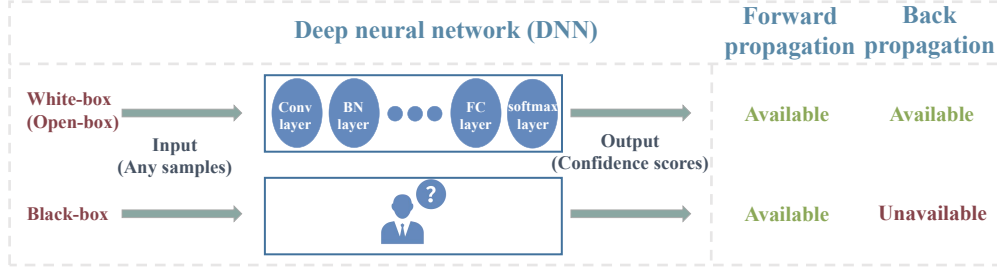
## ABSTRACT

Backpropagation algorithm has played a significant role in the development of deep learning. However, there exist some limitations associated with this algorithm, such as getting stuck in local minima and experiencing vanishing/exploding gradients, which have led to questions about its biological plausibility. To address these limitations, alternative algorithms to backpropagation have been preliminarily explored, with the Forward-Forward (FF) algorithm being one of the most well-known. In this paper we propose a new learning framework for neural networks, namely **Cascaded Forward (CaFo)** algorithm, which does not rely on BP optimization as that in FF. Unlike FF, CaFo directly outputs label distributions at each cascaded block and waives the requirement of generating additional negative samples. Consequently, CaFo leads to a more efficient process at both training and testing stages. Moreover, in our CaFo framework each block can be trained independently, allowing easy deployment to parallel acceleration systems. The proposed method is evaluated on four public image classification benchmarks, and the experimental results illustrate significant improvement in prediction accuracy in comparison with recently proposed baselines. Our code is available at: <https://github.com/Graph-ZKY/CaFo>.

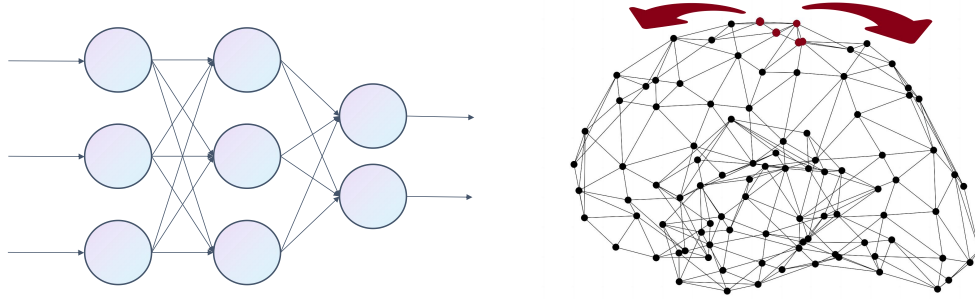
## 1 Introduction

Backpropagation (BP) algorithm [1] is a powerful technique that has proven to be effective for training deep neural networks on a wide range of tasks, including image classification [2], object detection [3], machine translation [4], text summarization [5] and graph learning [6]. The algorithm's ability of adjusting the weights based on the error between the prediction and the ground truth allows the network to learn and improve over time, making it a cornerstone of deep learning and artificial intelligence. Nevertheless, despite its effectiveness, BP suffers from several limitations in practical applications. These limitations include the problems of local minima [7], vanishing/exploding gradients [8], overfitting [9], slower convergence and non-convex optimization [10], which may negatively impact the training process. Additionally, BP relies on complete understanding of the computations performed during the forward pass in order to correctly calculate the derivatives. This characteristic makes it difficult to be generalized to the black-box systems where the internal working processes are not transparent, as shown in Figure 1 (a). For biological plausibility, it seems that backpropagation remains implausible as a model of how cortex learns, despite considerable effort to invent ways in which it could be implemented by real neurons. For example, as shown in Figure 1 (b), the connections between different cortical areas do not mirror the bottom-up connections of backpropagation-based deep learning models. Instead, they go in loops, in which neural signals traverse several cortical layers and then return to earlier areas. The bottom-up connections of backpropagation leads to the detachment of learning and inference, in which the training algorithm must stop inference to perform backpropagation in order to adjust the weights of a neural network. In contrast, the brain receives a constant stream of information, and the perceptual system needs to perform inference and learning in real time without stopping to perform backpropagation.

Due to the aforementioned limitations and the apparent difference in mechanisms between BP neural network and real cortical neurons as shown in Figure 1, some researchers [11–15] have raised concern about the biological plausibility of BP, questioning whether it has some other way of getting the gradients needed to adjust the weights on connections. This has prompted researchers to search for alternate algorithms to train deep neural networks. One of the most recent alternatives is the Forward-Forward (FF) algorithm [15], which replaces the forward and backward



(a)



(b)

Figure 1: Intuitive implausibilities of backpropagation algorithm: (a) BP requires access to the internal structure of the neural network, including the activation functions and the number of layer. This is not always possible in black box systems, where the internal workings of the system are not transparent. (b) The active cortex (indicated with red) processes information in multiple directions, while in neural networks, information typically flows in a single direction to optimize a specific objective.

passes of BP with two forward passes that work on positive and negative data with opposite optimization objectives. FF initiates preliminary and inspiring investigations in this field, demonstrating its potential as an alternative to BP. This opens up significant opportunities for further research and in-depth exploration.

Thus motivated, in this paper we present a flexible and effective learning framework for neural networks, namely **Cascaded Forward (CaFo)** algorithm, which offers a promising alternative to the traditional BP algorithm. Our CaFo framework consists of multiple cascaded neural blocks, with a layer-wise predictor attached to each block. Each neural block is composed of one or more convolutional layers, pooling layers, normalization layers, activation layers, etc., of which the aim is to project the input feature map into a new feature space. All the neural blocks are concatenated in sequence and pre-trained according to a non-backpropagation method. Each layer-wise predictor takes the feature map computed by the corresponding pre-trained neural block as input and outputs the prediction of the task, and it is trained independently without backpropagation according to the errors between the output and the ground truth. During the test period, the final prediction is a combination of the predictions from all predictors.

Roughly speaking, instead of performing backpropagation using the chain rule, the CaFo network performs only forward pass to directly estimate the prediction errors and updates the parameters *in situ*. Overall, our algorithm is a step forward on the basis of the FF algorithm and brings several advantages. Firstly, it eliminates the necessity for negative sampling, thereby reducing the impact of sampling quality on model performance and increasing stability. Secondly, our method directly outputs the prediction of multi-class distribution rather than a simple goodness estimation, and it is thus more suitable for multi-class prediction tasks. Finally, the pre-trained neural blocks enable independent training of the attached predictor without dependencies on the training outcomes of previous blocks. This design not only enhances portability and flexibility compared with FF but also facilitates straightforward deployment into parallel acceleration systems in certain scenarios.

For evaluation we test our algorithm on four public image classification benchmarks, in comparison with FF and other non-backpropagation algorithms. The experimental results show that our CaFo algorithm exhibits effectiveness and efficiency on image classification tasks, and outperforms all compared baselines by a remarkable margin.

In summary, with the proposed learning procedure for neural networks, we make contribution in three folds:

- We propose CaFo, a novel neural network training procedure that offers significant improvements in image classification task compared with FF and other state-of-the-art non-BP approaches.
- We develop distinct training strategies for the two core components of CaFo (i.e., neural blocks and predictors), respectively, to address the diverse requirements of both training efficiency and prediction effectiveness.
- Extensive experiments are conducted, and the experimental results demonstrate the superiority of our method against the state-of-the-art approaches.

This article is structured as follows: In Section II, we begin with a concise review of the well-known backpropagation algorithm and some representative non-backpropagation (non-BP) algorithms. Section III presents our proposed method, CaFo, encompassing notations, a pipeline overview, and the detailed introduction of neural block training, predictor training, and the inference process. In Section IV, we delve into discussions regarding the biological plausibility of CaFo and its relationship to some classic algorithms. Subsequently, in Section V, we present detailed experimental results and performance analyses obtained from real-world datasets. Finally, in Section VI, we present future work and provide concluding remarks.

## 2 Related Work

### 2.1 Backpropagation Algorithm

Backpropagation algorithm is a widely used training algorithm for deep neural networks, the goal of which is to optimize the parameters of the networks by minimizing the prediction error between the network’s output and the ground truth. BP uses the gradient descent algorithm for optimization. The gradient of the loss function is calculated using the chain rule [16] of differentiation, where the error is propagated backwards through the network and the parameters are updated according to the negative gradient.

BP is favored for its ease of implementation, computational efficiency, and effectiveness on a variety of problems. However, it also suffers from certain limitations that have been deeply studied [17–24]. For the overfitting problems in deep neural network, some studies propose regularization tricks such as L1, L2 regularization [25] and dropout [17] to enhance the generalizability of models. Other works propose to solve this problem by data augmentation [18] and training tricks (e.g. early stopping [19]). For the gradient problems during training, there are several techniques to address the vanishing and exploding gradients, including weight initialization [20], gradient clipping [21], non-linear activation functions [22], skip connections [23], and normalization techniques [24]. To address the problems of local minima and slow convergence, mini-batch training is encouraged, and some ingenious optimization algorithms for deep neural network is designed, such as Adam [26] and RMSProp [27].

Despite these efforts, problems with backpropagation still occur stochastically, highlighting the need for an in-depth understanding of deep learning and experience in model tuning. Due to the limitations mentioned above, some recent researches [15, 28] have even raised doubts about the biological plausibility of BP.

### 2.2 Non-backpropagation Algorithm

The biological implausibility of the BP mainly lies in the problems of weight transport [29], non-locality, freezing of neural activity, and update locking [30]. The feedback alignment algorithm (FA) [11] solves weight transport problem and demonstrates that using fixed random weights in the feedback pathway allows conveying useful error gradient information. Following this pipeline, the DFA [12], DRTP [13] and PEPITA [14] algorithms are proposed in sequence, which partially address the remaining three problems by introducing non-backpropagation updating algorithms for training. The most recent concern raised by Hinton [15] about the biological plausibility of the BP has spurred the introduction of the FF algorithm as an alternative to better align with neurophysiological findings. FF replaces the forward and backward passes of BP with two forward passes that work on positive and negative data with opposite optimization objectives. FF has been validated to be an effective alternative to BP in some tasks, due to its simplicity and flexibility for network architecture and its ability to handle non-differentiable components. Nevertheless, some potential limitations of FF make it hard to be integrated into existing learning systems, including the high computational demanding of two forward passes, the high dependence on positive and negative sampling strategies, and the rough approximation of optimization objective.

Our CaFo framework follows this pipeline and works without backpropagation. Moreover, it eliminates the need for negative sampling and updates the parameters directly according to the errors between the predicted output and the ground truth. It simplifies both training and testing processes, and gains significant improvement in prediction accuracy.

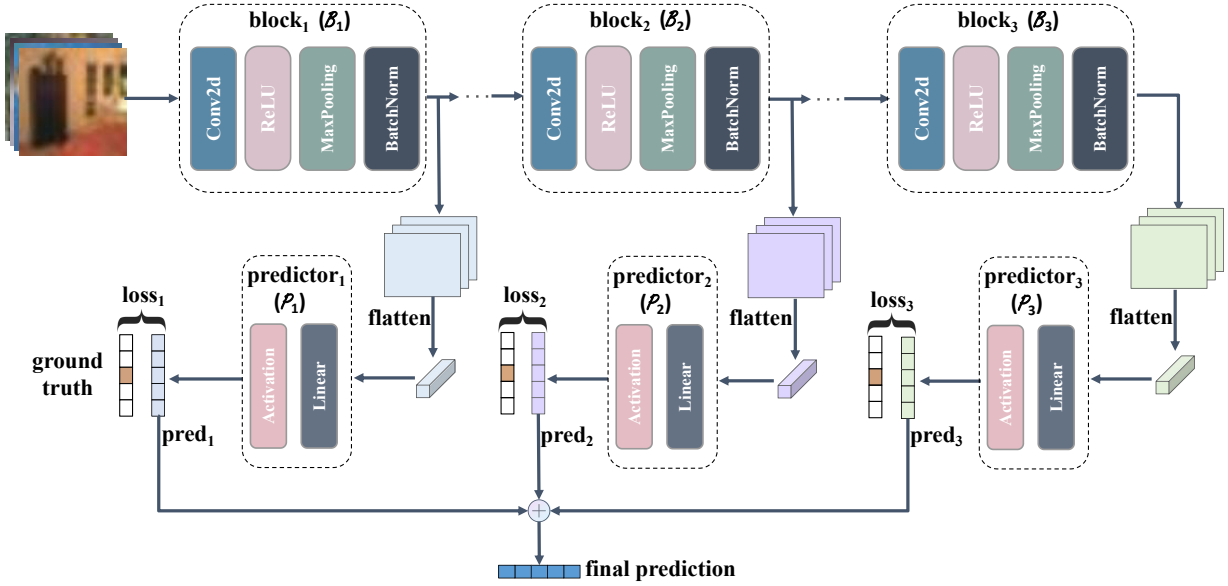


Figure 2: Overall architecture of the CaFo network. Two main components are cascaded in the network: (1) a neural block that extracts representation from source data (e.g., images) at a specific scale, and (2) a predictor that takes the intermediate representation extracted by the corresponding block as input and outputs the prediction results.

### 3 The Proposed Method

In this section, we describe the details of the proposed CaFo framework. Following the work of FF, we apply this method to the task of image classification, a well-established challenge in computer vision. To ensure a clear and consistent presentation throughout the paper, we firstly provide an overview of the notations and mathematical symbols used in our analysis. This is followed by an elaboration of the overall architecture of the CaFo framework. Finally we describe the training and inference processes of the proposed approach.

#### 3.1 Notations

To ensure clarity and consistency in our presentation, we provide a brief overview of the notations we used in this paper. The scalars used in our analysis are represented by italic lowercase letters (e.g.,  $m, n$ ), while the vectors used in our method are represented by bold lowercase letters (e.g.,  $\mathbf{x}, \mathbf{y}$ ). All vectors in this paper are assumed to be column vectors. Bold uppercase letters (e.g.,  $\mathbf{W}, \mathbf{Z}$ ) are used to represent matrices. Furthermore, the subscripts and superscripts in our notations represent the indices and exponents, respectively. In the rest of the paper, we will provide more detailed definitions and explanations of the notations as needed.

#### 3.2 Pipeline Overview

The pipeline of the proposed CaFo framework is depicted in Figure 2, which consists of multiple repeated stacks of two key components:

**Neural Block:** The neural block, abbreviated as *block* for convenience in the rest of the paper, is comprised of multiple element units, such as convolutional layer, pooling layer, normalization layer, activation layer, etc. Its architecture depends on the specific task. For image classification in our experiments, it consists of a convolutional layer followed by a ReLU activation function, a max-pooling layer and a batch normalization layer. Multiple blocks with the same or different architectures are cascaded to extract intermediate image representations at different scales. It is worth noting that the block can be either parametric or non-parametric, which provides flexibility in its construction.

**Predictor:** As illustrated in Figure 2, each block is equipped with a layer-wise predictor that consists of a fully connected layer followed by an activation function (e.g., softmax) if necessary. The predictor takes the intermediate

representation extracted by the block as input and outputs a prediction result for the task. Since there is only one fully connected layer, the predictor can be trained *in situ* without backpropagation.

As described below, each block is trained separately, independent of the predictors in our network.

### 3.3 Training of Neural Block

Inspired by DFA [12], we employ a non-BP approach to train the blocks. Firstly, we ignore the predictors and concatenate all the blocks in sequence. An additional output layer is introduced, performing a linear transformation followed by sigmoid activation after the last block. We firstly introduce the training strategy for the fully connected layer and subsequently extend it to cover the convolutional layer.

**Fully connected layer** Let  $(\mathbf{x}, \mathbf{y})$  be the input-output vectors of a sample we want the network to learn. For simplicity, we assume that the network has only two blocks and omit normalization layers. Denoting that  $\mathbf{V}_i$  and  $\mathbf{b}_i$  the weights and biases for the units in hidden layer  $i$ , the forward pass in the network is calculated as:

$$block_1 : \mathbf{a}_1 = (\mathbf{V}_1 \mathbf{x}) + \mathbf{b}_1, \quad \mathbf{h}_1 = f(\mathbf{a}_1), \quad (1)$$

$$block_2 : \mathbf{a}_2 = (\mathbf{V}_2 \mathbf{h}_1) + \mathbf{b}_2, \quad \mathbf{h}_2 = f(\mathbf{a}_2), \quad (2)$$

$$output\_layer : \mathbf{a}_y = (\mathbf{V}_3 \mathbf{h}_2) + \mathbf{b}_3, \quad \hat{\mathbf{y}} = \text{sigmoid}(\mathbf{a}_y), \quad (3)$$

where  $f(\cdot)$  the non-linearity used in hidden layers.

We update the parameters of blocks according to the classification error:  $\mathbf{e} = \hat{\mathbf{y}} - \mathbf{y}$ . For the linear output layer, the weight and bias updates are calculated as:

$$\mathbf{V}_3 \leftarrow \mathbf{V}_3 - \eta \Delta \mathbf{V}_3, \quad \Delta \mathbf{V}_3 = -\mathbf{e} \mathbf{h}_2^T, \quad (4)$$

$$\mathbf{b}_3 \leftarrow \mathbf{b}_3 - \eta \Delta \mathbf{b}_3, \quad \Delta \mathbf{b}_3 = -\mathbf{e}, \quad (5)$$

where  $\eta$  is the learning rate. For fully connected layer, in BP the changes of weights (i.e.,  $\Delta \mathbf{V}_i$ ) and biases (i.e.,  $\Delta \mathbf{b}_i$ ) depend on the value of  $\mathbf{V}_{i+1}$  and the error gradient  $\Delta \mathbf{h}_i$ . As for non-BP framework, we follow DFA to use a fixed random feedback matrix  $\mathbf{B}_i$  for each fully connected layer as follows:

$$\Delta \mathbf{h}_i = \mathbf{B}_i \mathbf{e}. \quad (6)$$

By employing this approach, the error gradient for each fully connected layer can be calculated in parallel, waiving the need for a backward pass of the error gradient.

**Covolutional layer** Let  $(\mathbf{X}, \mathbf{y})$  be the input-output tensors of a sample we want the network to learn. For  $i$ -th block we have:

$$block_i : \mathbf{a}_i = (\mathbf{V}_i * \mathbf{h}_{i-1}) + \mathbf{b}_i, \quad \mathbf{h}_i = f(\mathbf{a}_i), \quad (7)$$

where  $*$  denotes 2-D convolution operation, and  $f(\cdot)$  the non-linearity used in hidden layers. We initialize  $\mathbf{h}_0$  as  $\mathbf{X}$  for the first block, and for the sake of simplicity, we set the stride to one. The error gradient for each layer is calculated as

$$\Delta \mathbf{h}_{i-1} = \Delta \mathbf{a}_i * ROT180(\mathbf{V}_i), \quad (8)$$

where  $ROT180(\cdot)$  denotes a 180-degree rotation. Since  $\Delta \mathbf{a}_i$  is associated with the error term  $\mathbf{e} = \hat{\mathbf{y}} - \mathbf{y}$ , we represent this relationship as  $\Delta \mathbf{a}_i = g_i(\mathbf{e})$ . Following the approach of DFA, we employ a fixed random feedback matrix for each convolutional layer, which is outlined as follows:

$$\Delta \mathbf{h}_{i-1} = g_i(\mathbf{e}) * ROT180(\mathbf{B}_i). \quad (9)$$

For computational simplicity, we further employ Eq. 6 to approximate Eq. 9 during the training of convolutional layers.

Noting that, for simplicity the parameters of all blocks can be randomly parameterized without training. The experimental results in Tables 1 and 2 demonstrate that, despite simplicity, the blocks parameterized with Kaiming uniform [31] can produce discriminative representations for prediction, achieving better performance in comparison with baseline algorithms. Moreover, blocks with finely trained parameters excel in extracting intermediate features and enhance the model performance.

### 3.4 Training of Predictor

For each predictor, we optimize it by minimizing the errors between the prediction and the ground truth. Different loss functions can be adopted in our framework, which lead to different optimization strategies and algorithms. In experiments we compare three different loss functions: mean square error loss (MSE), cross-entropy loss (CE) and sparsemax loss (SL) [32]. Here we describe the corresponding formulations and optimization algorithms respectively.

**MSE loss** The optimization objective of the predictor can be expressed as follows:

$$\arg \min_{\mathbf{W}} \frac{1}{2m} \|\mathbf{HW} - \mathbf{Y}\|_F^2, \quad (10)$$

where  $m$  denotes the number of training samples,  $\mathbf{H} \in \mathbb{R}^{m \times d}$  the  $d$ -dimensional intermediate representation output by the block,  $\mathbf{Y} \in \mathbb{R}^{m \times c}$  the one-hot labels of training samples,  $\mathbf{W} \in \mathbb{R}^{d \times c}$  the parameters of predictor to be optimized,  $\|\cdot\|_F$  the Frobenius norm, and  $c$  indicates the number of classes. The close-form solution for Eq. 10 can be obtained by solving the following equation for  $\mathbf{W}$ :

$$\frac{\partial}{\partial \mathbf{W}} \left( \frac{1}{2m} \|\mathbf{HW} - \mathbf{Y}\|_F^2 \right) = \frac{1}{m} \mathbf{H}^T (\mathbf{HW} - \mathbf{Y}) = 0, \quad (11)$$

and the solution of Eq. 10 is:

$$\mathbf{W} = (\mathbf{H}^T \mathbf{H})^{-1} \mathbf{H}^T \mathbf{Y}. \quad (12)$$

**Cross-entropy loss** The optimization objective of the predictor is written as:

$$\arg \min_{\mathbf{W}} -\frac{1}{m} \sum_{i=1}^m \sum_{j=1}^c \mathbf{Y}_{i,j} \cdot \ln \mathbf{P}_{i,j}, \quad (13)$$

$$\mathbf{P}_{i,j} = \frac{\exp(\mathbf{H}_{i,:} \mathbf{W}_{:,j})}{\sum_{k=1}^c \exp(\mathbf{H}_{i,:} \mathbf{W}_{:,k})}, \quad (14)$$

where  $\mathbf{Y}_{i,j}$  is the value at  $i$ -th row and  $j$ -th column of  $\mathbf{Y}$ ,  $\mathbf{H}_{i,:}$  and  $\mathbf{W}_{:,j}$  respectively denote the  $i$ -th row of  $\mathbf{H}$  and the  $j$ -th row of  $\mathbf{W}$ , and  $\mathbf{P}_{i,j}$  the corresponding output of softmax.

As the close-form solution for Eq. 13 is unavailable, we use gradient descent algorithm to optimize it. To do this, we need to calculate the gradient for  $\mathbf{W}$  according to the Jacobian matrix, and set a fixed step size to optimize Eq. 13 iteratively. The Jacobian matrix is calculated by  $\mathbf{J} = \mathbf{P} - \mathbf{Y}$ , and the gradient for  $\mathbf{W}$  is computed as:

$$\frac{1}{m} \sum_{i=1}^m [\mathbf{H}_{i,:}]^T \otimes \mathbf{J}_{i,:}, \quad (15)$$

where  $\otimes$  denotes the Kronecker product.

**Sparsemax loss** Different than softmax, the sparsemax function [32] produces sparse output probabilities by enforcing a constraint that the output confidence vector has at most a certain number of non-zero elements. It encourages the model to only assign high probabilities to the most relevant classes, while setting all other probabilities to zero. The sparsemax function achieves this by projecting the input vector onto a simplex, which is a convex polytope whose vertices lie on the coordinate axes. The formulaic expression for sparsemax is as follows:

$$\text{sparsemax}(\mathbf{z}) := \arg \min_{\mathbf{p} \in \Delta^{c-1}} \|\mathbf{p} - \mathbf{z}\|^2, \quad (16)$$

where  $\Delta^{c-1} := \{\mathbf{p} \in \mathbb{R}^c | \mathbf{1}^T \mathbf{p} = 1, \mathbf{p} \geq \mathbf{0}\}$  is the  $(c-1)$ -dimensional simplex. The resulting projection is unique and can be computed efficiently using a sorting algorithm [32].

We also use the gradient descent algorithm to approximate the sparsemax loss optimization objective:

$$\arg \min_{\mathbf{W}} \frac{1}{m} \sum_{\mathbf{z} \in \mathbf{Z}} \sum_{k=1}^c L_{\text{sparse}}(\mathbf{z}; k), \quad (17)$$

$$L_{\text{sparse}}(\mathbf{z}; k) = -\mathbf{z}_k + \frac{1}{2} \sum_{j \in S(\mathbf{z})} (\mathbf{z}_j^2 - \tau^2(\mathbf{z})) + \frac{1}{2}, \quad (18)$$

where  $\mathbf{z} \in \mathbb{R}^c$  is the output logits,  $\mathbf{Z} = [\mathbf{z}_1; \mathbf{z}_2; \dots; \mathbf{z}_m]^T \in \mathbb{R}^{m \times c}$  the set of output logits of  $m$  training samples,  $S(\mathbf{z})$  the support of  $\text{sparsemax}(\mathbf{z})$ , and  $\tau(\mathbf{z}) = \frac{(\sum_{j \in S(\mathbf{z})} \mathbf{z}_j)^{-1}}{|S(\mathbf{z})|}$ . The Jacobian matrix is calculated by:

$$\mathbf{J} = \text{sparsemax}(\mathbf{Z}) - \mathbf{Y}, \quad (19)$$

and the gradient for  $\mathbf{W}$  is computed using the same formula as Eq. 15.

---

**Algorithm 1** The CaFo Algorithm

---

**Input:**

$\mathbf{X} = [\mathbf{x}_1; \mathbf{x}_2; \dots; \mathbf{x}_m]^T$ :  $m$  samples for training,  
 $\hat{\mathbf{X}} = [\hat{\mathbf{x}}_1; \hat{\mathbf{x}}_2; \dots; \hat{\mathbf{x}}_n]^T$ :  $n$  samples for inference,  
 $\mathbf{Y} = [\mathbf{y}_1; \mathbf{y}_2; \dots; \mathbf{y}_m]^T$ : one-hot labels of training samples,  
 $\mathcal{B}_i$ : the  $i$ -th block,  
 $\mathcal{P}_i$ : the predictor of  $i$ -th block.

**Manual factors:**

$r$ : the number of blocks of CaFo

**Output:**  $\hat{\mathbf{y}} = [\hat{y}_1, \hat{y}_2, \dots, \hat{y}_n]^T$

**Training:**

- 1:  $\mathbf{H}^0 = [\mathbf{h}_1^0; \mathbf{h}_2^0; \dots; \mathbf{h}_m^0]^T = \mathbf{X}$
- 2: **for**  $i = 1 \rightarrow r$  **do**
- 3: Get the intermediate representations of the  $i$ -th block:  
 $\mathbf{H}^i = [\mathbf{h}_1^i; \mathbf{h}_2^i; \dots; \mathbf{h}_m^i]^T = \mathcal{B}_i(\mathbf{H}^{i-1})$
- 4: Get the predictions of the  $i$ -th block:  
 $\tilde{\mathbf{Y}}^i = [\tilde{\mathbf{y}}_1^i; \tilde{\mathbf{y}}_2^i; \dots; \tilde{\mathbf{y}}_m^i]^T = \mathcal{P}_i(\mathbf{H}^i)$
- 5: Optimizing  $\mathcal{P}_i$  according to the errors between  $\tilde{\mathbf{Y}}^i$  and  $\mathbf{Y}$
- 6: **end for**

**Inference:**

- 1:  $\hat{\mathbf{H}}^0 = [\hat{\mathbf{h}}_1^0; \hat{\mathbf{h}}_2^0; \dots; \hat{\mathbf{h}}_n^0]^T = \hat{\mathbf{X}}$
  - 2: **for**  $i = 1 \rightarrow r$  **do**
  - 3: Get the intermediate representations of the  $i$ -th block:  
 $\hat{\mathbf{H}}^i = [\hat{\mathbf{h}}_1^i; \hat{\mathbf{h}}_2^i; \dots; \hat{\mathbf{h}}_n^i]^T = \mathcal{B}_i(\hat{\mathbf{H}}^{i-1})$
  - 4: Get the predictions of the  $i$ -th block:  
 $\hat{\mathbf{Y}}^i = [\hat{\mathbf{y}}_1^i; \hat{\mathbf{y}}_2^i; \dots; \hat{\mathbf{y}}_n^i]^T = \mathcal{P}_i(\hat{\mathbf{H}}^i)$
  - 5: **end for**
  - 6:  $\hat{\mathbf{Y}} = \sum_{i=1}^r \hat{\mathbf{Y}}^i$
  - 7:  $\hat{\mathbf{y}} = \arg \max \hat{\mathbf{Y}}$
  - 8: **return**  $\hat{\mathbf{y}}$
- 

### 3.5 Inference Process

In the inference phase each predictor outputs a prediction and the final prediction is determined by combining all the individual ones. In classification tasks this is typically done by summing the predictions of all the predictors and choosing the index of the maximum value as the predicted class.

The detailed optimization algorithm is shown in Algorithm 1. During the training stage, the model generates the prediction block by block at step 3 and step 4, where  $\mathbf{h}_j^i \in \mathbb{R}^d$  and  $\tilde{\mathbf{y}}_j^i \in \mathbb{R}^c$  are the intermediate representation and prediction of  $i$ -th block for  $j$ -th training sample. At step 5, the predictors are updated based on a specific loss function (e.g., Eqs. 10, 13, and 17), using the corresponding optimization strategy. The predictors consist of a fully connected layer for MSE, and an additional activation layer is followed for CE and SL. During the inference stage, each predictor outputs a prediction for the test set at step 3 and step 4, where  $\hat{\mathbf{y}}_j^i \in \mathbb{R}^c$  is the prediction of  $i$ -th block for the  $j$ -th test sample. Then the predictions of each block (i.e.,  $\hat{\mathbf{Y}}^i \in \mathbb{R}^{n \times c}$ ) are summed to derive the final prediction  $\hat{\mathbf{Y}} \in \mathbb{R}^{n \times c}$ . The predicted class for  $j$ -th test samples is then calculated at step 7 as:  $\hat{y}_j = \arg \max_k \hat{\mathbf{Y}}_{j,k}$ .

## 4 Discussion

### 4.1 Relationship to Previous Non-BP Algorithms

In CaFo, the training mechanism of neural blocks shares similarities with DFA [12]. Nevertheless, each predictor independently makes its greedy decision about the class to which the input belongs, according to its local loss. This approach differs from DFA, which relies solely on a global loss. The predictors are trained on intermediate features

Data	Model Parameters	Intermediate Representation	Loss
<b>X</b> : Input $X_{\text{pos}}$ : Positive input data $X_{\text{neg}}$ : Negative input data <b><math>\hat{Y}</math></b> : Output <b>Y</b> : Ground truth	$W_N$ : Forward weights $G_N$ : Weights of block $B_N$ : Weights of random feedback matrix	$h_N$ : Hidden states $h_N^{\text{pos}}$ : Hidden states of 2 <sup>nd</sup> forward pass on negative data $z_N$ : Intermediate outputs of cascaded blocks	<b>L</b> : Global loss $L_N$ : Local loss

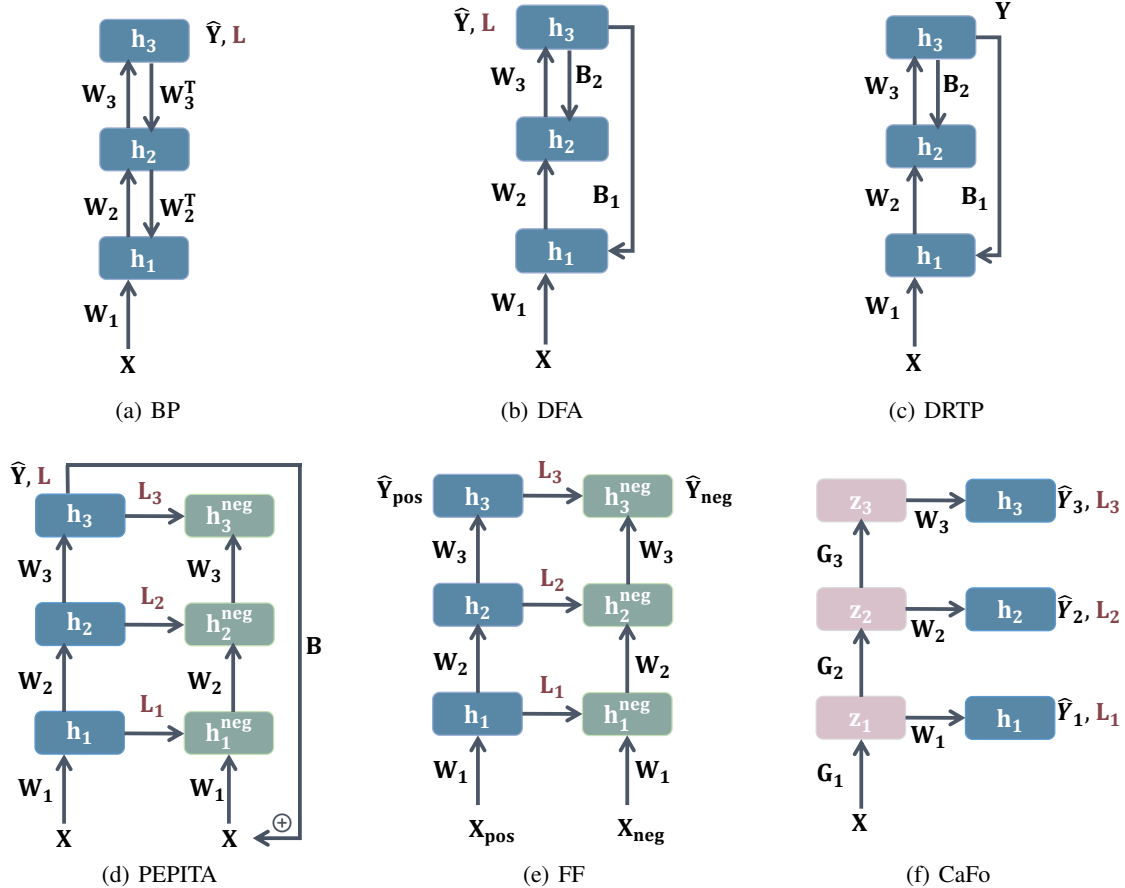


Figure 3: Comparison of Backpropagation (BP) algorithm, four non-BP algorithms (FF, DFA, DRTP, PEPITA) and the proposed Cascaded Forward (CaFo) algorithm.

extracted at various scales, enhancing the model performance by capturing more intricate patterns, as demonstrated in the experimental results depicted in Figure 5.

Similarly, DRTP [13] and PEPITA [14] update each layer in a forward direction. However, DRTP updates layers directly using ground truth without considering global or local loss calculations, potentially limiting the learning of class distributions. PEPITA updates layers based on the difference in intermediate representations between two forward passes, which is also not directly tied to classification loss.

Our CaFo algorithm builds upon the FF algorithm [15], and shares significant similarities with it on local optimization. They both have advantages over the traditional backpropagation methods, such as the ability to handle non-differentiable components and increased flexibility in network architecture. However, the differences between FF and CaFo are also significant. For FF, positive and negative samples are selected and samples are used to calculate a local loss for each layer. Once the local loss of the  $i$ -th layer is calculated (i.e.,  $L_i$ ), the weights of the  $i$ -th layer (i.e.,  $W_i$ ) will be updated according to  $L_i$ , and after that the  $W_i$  will be fixed during the next training stage for the  $(i + 1)$ -th layer. The loss for each layer is calculated in sequence, and the layers of FF are sequentially updated by greedy decision. In contrast, our CaFo has only one input stream, avoiding the negative sampling process. Although the sequential update of each block in our framework is similar to FF, the local loss in CaFo directly measures the error between the predicted



multi-class distribution and ground truth, rather than a simple goodness estimation. Moreover, since the pre-trained neural blocks enable independent training of the attached predictor without dependencies on the training outcomes of previous blocks, our method not only enhances portability and flexibility compared with FF but also facilitates straightforward deployment into parallel acceleration systems when the blocks are randomly parameterized without training.

The graphical depiction and comparison of BP, non-BP algorithms and proposed CaFo are shown in Figure 3.

## 4.2 Biological Plausibility

By introducing the aforementioned training strategies for neural blocks and predictors, our approach overcomes several of the biological unrealistic aspects associated with backpropagation. Firstly, CaFo avoids the weight transport problem since it does not utilize weight-specific feedback signals. Secondly, with regard to locality, the prescribed updates for each synapse are solely based on the activity of the connected nodes, avoiding any global errors that may perturb node activity. Thirdly, CaFo does not freeze the neural activity to propagate and apply the modulatory signal. Finally, in CaFo, updates are performed block by block in a sequential manner, partially mitigating the update locking issue. The weights of the first predictor can be updated right at the second forward pass, not requiring waiting for downstream predictor updates. This allows the forward computation for the next input sample at the first predictor to start in parallel with the update of the second predictor from the previous sample, suggesting that CaFo can be suitable for edge computing devices requiring fast processing of the input signals.

## 4.3 Relationship to Bagging

Bagging (Bootstrap Aggregating) [33] is an ensemble learning technique that combines multiple models to enhance the overall performance of a machine learning algorithm. The basic idea of bagging is to train some individual models on different subsets of the training data, and then combine their outputs by taking the average of majority vote. Bagging eliminates the variance of the individual models and prevents overfitting by reducing the impact of noisy data points or outliers. This is achieved by generating multiple random samples of the training data, each of which is adopted to train a different model. The results are then combined to make a final prediction.

Our method can be viewed as a special case of Bagging, where the predictors are trained on intermediate features extracted at various scales, rather than on the original images. In other words, our approach does not train individual models on different subsets of the training data, but assigns each individual model with a specific intermediate representations of the complete training data, which has the advantage of utilizing the information contained in intermediate features extracted at different scales, thus helping to capture more complex patterns to improve the performance.

## 4.4 Relationship to Co-training

Co-training [34] is a popular semi-supervised learning method that trains multiple classifiers on the same dataset with different data views to improve classification accuracy by leveraging unlabeled data. The process involves initializing several classifiers with different feature sets and training them on labeled data subsets. The classifiers then exchange predictions on unlabeled data and retrain on new labeled data subsets using these predictions. This iterative process continues until convergence, with classifiers refining their feature sets and predictions.

As a semi-supervised learning technique, co-training is usually applied in scenarios where the labeled data is scarce but there is a large amount of unlabeled data available. In contrast, our method is a fully-supervised approach that trains multiple predictors for each block in parallel, allowing for more efficient and stable optimization of the model in supervised learning scenarios. As our method and co-training share the same idea of training multiple modules to improve the overall performance, our approach can be seen as the first stage of a co-training process, and can be easily transformed into a co-training method if an augmentation strategy for training set is introduced.

# 5 Experiments

The aim of this paper is to introduce the CaFo algorithm and to show that it works in relatively small neural networks containing a few million connections. We evaluate the proposed algorithm on four benchmarks in comparison with four non-BP algorithms. Future research will focus on investigating the scalability of the proposed approach to handle large neural networks with a much greater number of connections.

Table 1: Error rate on MNIST and CIFAR-10. The best and the second best results on each dataset are indicated with red and blue colors, respectively.

Error rate (%)	MNIST		CIFAR-10	
	Training	Test	Training	Test
<b>DFA</b>	0.77	1.37	16.75	34.29
<b>DRTP</b>	0.62	1.41	32.55	39.50
<b>PEPITA</b>	1.81	1.84	12.14	36.55
<b>FF (Hinton)</b>	-	1.37	24.00	41.00
<b>FF (reprodu.)</b>	0.57	2.02	10.36	46.03
<b>CaFo+MSE (rand)</b>	0.71	1.55	12.56	34.79
<b>CaFo+CE (rand)</b>	0.04	1.30	9.22	32.57
<b>CaFo+SL (rand)</b>	0.06	1.20	9.58	33.74
<b>CaFo+MSE (dfa)</b>	0.48	1.24	10.99	33.46
<b>CaFo+CE (dfa)</b>	0.07	<b>1.18</b>	16.58	<b>30.52</b>
<b>CaFo+SL (dfa)</b>	0.17	<b>1.05</b>	11.42	<b>31.51</b>

Table 2: Error rate on CIFAR-100 and Mini-ImageNet

Error rate (%)	CIFAR-100		Mini-ImageNet	
	Training	Test	Training	Test
<b>DFA</b>	8.93	62.02	4.87	84.61
<b>DRTP</b>	12.07	63.15	80.51	87.91
<b>PEPITA</b>	1.32	61.26	0.16	90.87
<b>FF (reprodu.)</b>	29.32	78.67	27.84	91.71
<b>CaFo+MSE (rand)</b>	6.06	63.94	6.66	85.50
<b>CaFo+CE (rand)</b>	0.06	59.24	0.01	<b>78.20</b>
<b>CaFo+SL (rand)</b>	0.02	61.96	0.02	81.47
<b>CaFo+MSE (dfa)</b>	5.59	63.30	7.57	83.08
<b>CaFo+CE (dfa)</b>	2.72	<b>57.87</b>	0.01	<b>77.59</b>
<b>CaFo+SL (dfa)</b>	0.29	<b>58.97</b>	0.02	79.55

## 5.1 Datasets and Experimental Settings

The experiments are conducted on the MNIST [35], CIFAR-10 [36], CIFAR-100 [36] and Mini-ImageNet [37] datasets, all of which are widely-used in evaluating image classification algorithms. The standard splits of MNIST, CIFAR-10 and CIFAR-100 are directly used. In the case of Mini-ImageNet, which is often used for evaluating few-shot learning algorithms, all the classes are mixed, and for each class, 500 images are randomly sampled for training and 100 other images for testing.

We compare the proposed CaFo algorithm with four non-BP baseline algorithms: DFA [12], DRTP [13], PEPITA [14], and FF [15]. Specifically, for MNIST and CIFAR-10, we report the results of both the original version of FF [15] and our reproduced version. For CIFAR-100 and Mini-ImageNet, we only report the results of our reproduced version, as the source code of the original version is not publicly available.

For the proposed CaFo algorithm, we set the number of blocks  $r=3$ . Each block consists of a convolutional layer followed by a ReLU activation function, a max-pooling layer and a batch normalization layer. The convolutional layer of the three blocks have the same kernel size, padding and stride, which are set to be  $3 \times 3$ , 1 and 1, respectively; and the output channels are set to be 32, 128, and 512 respectively.

Table 3: Running time (seconds) for training and testing

	MNIST		CIFAR-10		CIFAR-100		Mini-ImageNet	
Time (seconds)	Training	Test	Training	Test	Training	Test	Training	Test
<b>FF (reprodu.)</b>	1655.22	0.02	2374.30	0.04	2376.77	0.42	3668.35	2.03
<b>DFA</b>	925.48	0.05	956.06	0.08	891.40	0.07	12942.45	0.31
<b>DRTP</b>	1054.29	0.06	955.01	0.08	892.88	0.08	12966.59	0.31
<b>PEPITA</b>	1551.94	0.10	1725.66	0.11	1717.39	0.11	41522.68	0.32
<b>CaFo+MSE (rand)</b>	0.62	0.06	1.21	0.07	1.40	0.08	22.55	0.32
<b>CaFo+CE (rand)</b>	612.22	0.05	873.56	0.07	1178.42	0.08	8101.42	0.52
<b>CaFo+SL (rand)</b>	837.56	0.06	1082.45	0.09	2192.56	0.10	14128.06	0.49
<b>CaFo+CE (dfa)</b>	1504.29	0.05	1770.82	0.07	2059.08	0.08	20842.13	0.52

As described in the previous section, different error measure functions can be adopted in our framework to guide the training. For analysis of the affects of different measure functions, in this section we report the results of three versions of our CaFo algorithms that adopt MSE (CaFo+MSE), cross-entropy (CaFo+CE) and sparsemax loss (CaFo+SL) as the measure function respectively. The predictor is a fully connected layer without bias, followed by a softmax layer for CE and SL. For CaFo+CE and CaFo+SL, each predictor is trained for 5000 epochs on MNIST and CIFAR-10, and for 1000 epochs on CIFAR-100 and Mini-ImageNet. We pre-initialize the parameters of each block with Kaiming uniform [31] (abbr. rand) and the pre-trained weights of DFA (abbr. dfa). All experiments are run on AMD EPYC 7542 32-Core Processor with an Nvidia GeForce RTX 3090 GPU.

## 5.2 Performance Comparison

The comparison of these algorithms on four datasets is summarized in Tables 1 and 2, in which the best and the second best results on each dataset are indicated with red and blue colors, respectively. It is obvious our method achieves overall improvements in test error rate compared with the baseline algorithms, obtaining the best or second best results on all datasets.

Even the simplest version (CaFo+MSE) demonstrates comparable test results with DFA, DRTP and PEPITA, while outperforming the reproduced FF on all datasets. For FF, we observe a significant discrepancy in its results between the two tables. FF has low training error when the dataset is easy to discriminate (e.g., CIFAR-10 and MNIST), but has high training error when tackling much more complex datasets (e.g., CIFAR-100 and Mini-ImageNet). The reason for the underfitting of FF probably lies in that the contrastive loss function based on positive and negative samples may not provide effective guidance. If the training samples for each class are scarce (e.g., CIFAR-10 vs. CIFAR-100) or the samples contain more irrelevant information (e.g., CIFAR-10 vs. Mini-ImageNet), the goodness estimation of positive and negative samples may not effectively help it to learn the potential distribution for each class because of the uncertain quality of negative samples and the lack of supervision information directly related to the labels. This may lead to the poor fitting for the training set. Similarly, the PEPITA algorithm that does not use classification loss as updating guidance also experiences consistently high training errors.

In contrast, our method exhibits a smaller discrepancy between these datasets and demonstrates better fitting ability in general. For test error rate, our CaFo outperforms all baseline algorithms by remarkable margins, especially on more complex datasets. For fair comparison, the results of CaFo are reported without deployment of any regularization trick, and thus the performance still has the potential to be improved by means of regularization strategies.

When comparing the three variants of our method that have randomly initialized block parameters, we observe that they obtain overall similar performance. CaFo+MSE has the highest test and train error on all the datasets, which suggests that MSE may not be good enough for error estimation in comparison with other two variants that adopt classification-oriented loss. In contrast, both CaFo+CE and CaFo+SL exhibit low training and test errors on all the four datasets, indicating the effectiveness of gradient descent optimization strategy and the superiority of CE and SL loss.

Given DFA’s effectiveness as a non-BP approach, we employ this method to train the parameters of each block and incorporate them into CaFo. The results show that CaFo (dfa) stands out as the better promising variant in comparison with CaFo (rand), achieving the best performance on all datasets. This intriguing observation suggests the feasibility

Table 4: Error rate with different training strategies for neural block

Error rate (%)	MNIST		CIFAR-10		CIFAR-100		Mini-ImageNet	
	Training	Test	Training	Test	Training	Test	Training	Test
<b>CaFo+MSE (rand)</b>	0.71	1.55	12.56	34.79	6.06	63.94	6.66	85.50
<b>CaFo+CE (rand)</b>	0.04	1.30	9.22	32.57	0.06	59.24	0.01	78.20
<b>CaFo+SL (rand)</b>	0.06	1.20	9.58	33.74	0.02	61.96	0.02	81.47
<b>CaFo+MSE (dfa)</b>	0.48	1.24	10.99	33.46	5.59	63.30	7.57	83.08
<b>CaFo+CE (dfa)</b>	0.07	1.18	16.58	30.52	2.72	57.87	0.01	77.59
<b>CaFo+SL (dfa)</b>	0.17	1.05	11.42	31.51	0.29	58.97	0.02	79.55
<b>CaFo+MSE (drtp)</b>	1.23	2.18	69.05	76.09	91.66	96.64	91.08	97.69
<b>CaFo+CE (drtp)</b>	1.33	1.64	30.40	37.58	23.82	64.50	0.03	81.92
<b>CaFo+SL (drtp)</b>	0.64	1.18	36.91	43.98	29.62	64.13	0.02	83.32
<b>CaFo+MSE (bp)</b>	0.51	1.74	1.98	29.52	2.41	62.51	6.41	82.18
<b>CaFo+CE (bp)</b>	0.17	1.06	1.74	26.30	0.03	56.52	0.01	74.95
<b>CaFo+SL (bp)</b>	0.21	0.96	0.33	27.17	0.02	58.60	0.01	76.58
<b>BP</b>	0.00	0.77	0.00	24.36	0.02	54.11	0.05	73.76

of integrating well-learned blocks trained by non-BP approaches into CaFo, potentially enhancing the representation ability of our model.

### 5.3 Time Comparison

Furthermore, we report the time required for training and testing our method as shown in Table 3. For each trial, we report the time required for the entire training process and the test time that includes the time for the forward pass of the test samples and the time for calculating the predicted category. The experimental settings are the same as those in Tables 1 and 2.

Comparing our three variants with baseline algorithms with whose blocks are randomly initialized, we observe that CaFo+MSE exhibits significant improvements in training efficiency, which has a much shorter training time than other reported approaches due to its direct calculation of the close-form solution without iterative training. Additionally, CaFo+CE and CaFo+SL also demonstrate comparable training efficiency when compared with baseline algorithms. Generally, our method shows significant improvements in training efficiency on MNIST, CIFAR-10, while it requires more time on CIFAR-100 and Mini-ImageNet. This is because we divide the training set of CIFAR-100 and Mini-ImageNet into 20 and 200 training batches, respectively, for the latter two variants due to memory limitation, which leads to the extra time consumption of data movement. However in the trials of MNIST and CIFAR-10 where the batch number equals one, we find the two CaFo variants are trained significantly faster than baseline algorithms. CaFo+CE (dfa) requires more time than other variants due to its time cost, which includes both the pre-initialization time of blocks and the training time of predictors. Despite this, the overall time cost is still acceptable when compared with baselines.

In terms of test time, our methods demonstrate overall comparable efficiency in comparison with DFA, DRTP and PEPITA due to their similar testing mechanism. Compared with FF, our method has comparable efficiency on MNIST and CIFAR-10, but significantly faster test time on CIFAR-100 and Mini-ImageNet. The reason for this lies in that our methods directly output the prediction of multi-class distribution, rather than compute the goodness estimation for each candidate label as that in FF. It tremendously improves the test efficiency when the number of categories is relatively large, such as in CIFAR-100 and Mini-ImageNet.

### 5.4 Investigation of different training strategies for neural block

In Table 4, we present the error rates of CaFo under various neural block training strategies: random initialization (abbr. rand), backpropagation (abbr. bp), direct feedback alignment (abbr. dfa), and direct random target projection (abbr. drtp). Additionally, we include results for an end-to-end BP method (abbr. BP) using the same architecture.

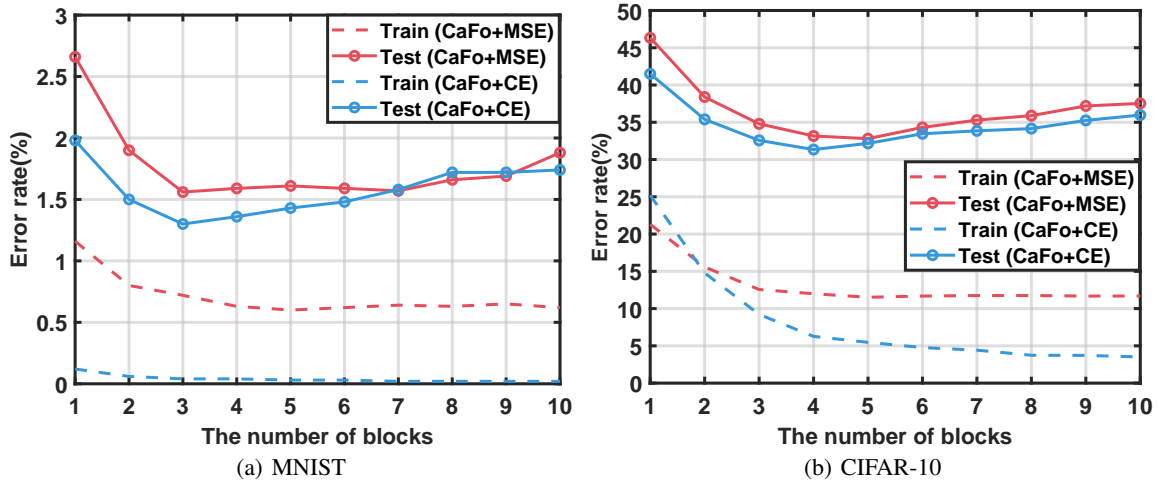


Figure 4: Error rates with different numbers of blocks.

CaFo trained with DFA displays an overall improvement compared with DRTP, highlighting DFA’s effectiveness as a non-BP method for training neural blocks. Surprisingly, even trials with random initialization show better performance than DRTP trials, underscoring the negative impact of an inappropriate training strategy on our method’s effectiveness.

In comparison with BP, although CaFo demonstrates superiority over competitive non-BP methods as evident in Tables 1 and 2, a notable performance gap still exists. When neural blocks are trained via backpropagation, CaFo, while not classified as a non-BP method, exhibits better sample fitting and a smaller error rate. Thus, achieving results comparable to backpropagation in non-BP environments remains a crucial research concern for CaFo, necessitating further investigation.

### 5.5 Investigation of the Number of Blocks

In Tables 1 and 2 we fix the number of blocks  $r=3$  for fair comparison with [15], in which the results of BP and FF with three layers are reported. However, the number of blocks may have diverse effects on different datasets. Here we investigate how the number of blocks influences the performance of our method.

In Figure 4 the error rates of CaFo+MSE and CaFo+CE with the number of blocks ranging from 1 to 10 are reported. As CIFAR-10 allows for a maximum of five  $2 \times 2$  max-pooling layers with  $stride=2$ , and four are available for MNIST, we remove the max-pooling layer from some blocks to keep the number of max-pooling layers within the upper bound if necessary. We observe that the training error rates (dashed lines) continuously decrease with the increasing number of blocks, indicating an enhancement of fitting ability as more learnable predictors are introduced. However, the test error curves (solid lines) show that the model is more or less affected by overfitting when the number of blocks is large.

### 5.6 Performance of each Predictor

To investigate the performance of each predictor, we report the test error rate of each predictor for CaFo+MSE and CaFo+CE, on four datasets. As shown in Figure 5, the test error of the final prediction is lower than those of previous individual predictors. Moreover, we find that a deep-layer predictor tends to obtain better performance than a shallow-layer one, demonstrating that the deep features play a more important role in enabling the predictor to make correct classifications. Overall, the results in Figure 5 well validate that the combination of all the predictors effectively contributes to the final prediction. Although both FF and our method illustrate that the straightforward sum of each block’s prediction (goodness) is sufficient for comparable performance, however, designing more effective fusion strategies for these non-backpropagation approaches remains to be studied.

## 6 Conclusion

In this paper, we propose a new learning framework, named CaFo, for deep neural networks that provides a viable alternative to the backpropagation algorithm. We show that CaFo has the feasibility of using only non-backpropagation

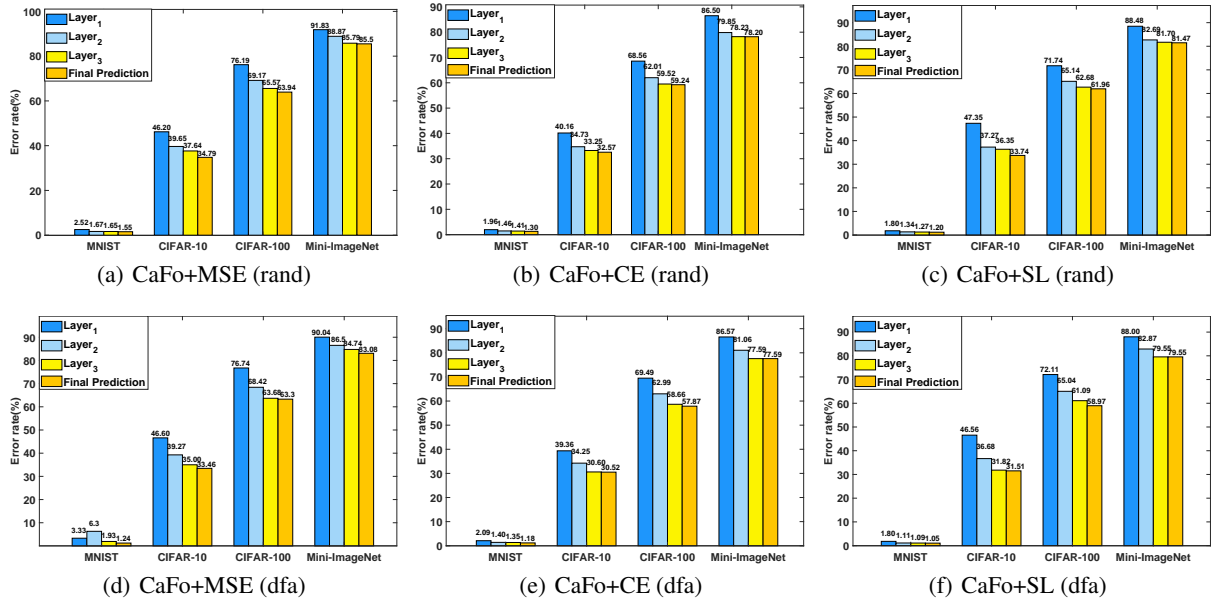


Figure 5: Test error rate of each predictor.

strategies without the need for backward pass in the training process of neural blocks and predictors, which offers significant improvements in prediction accuracy, training efficiency and simplicity compared with previous algorithms. Extensive experiments clearly validate the superior performance over the competitive methods on several public benchmarks.

While our proposed method has demonstrated remarkable performance, several challenges persist, particularly regarding scalability to large-scale datasets and more complex network architectures, which are common hurdles in neural network learning. In our future endeavors, we intend to investigate approaches to extend CaFo’s applicability to handle larger datasets and diverse architectures, such as graph neural networks (GNNs) and transformers. Furthermore, enhancing fusion strategies for CaFo is an area we plan to delve into, aiming to optimize the utilization of predictions from each predictor for improved overall performance.

## References

- [1] David E Rumelhart, Geoffrey E Hinton, and Ronald J Williams. Learning representations by back-propagating errors. *nature*, 323(6088):533–536, 1986.
- [2] Keke Tang, Yuexin Ma, Dingrui Miao, Peng Song, Zhaoquan Gu, Zhihong Tian, and Wenping Wang. Decision fusion networks for image classification. *IEEE Transactions on Neural Networks and Learning Systems*, 2022.
- [3] Zhong-Qiu Zhao, Peng Zheng, Shou-tao Xu, and Xindong Wu. Object detection with deep learning: A review. *IEEE transactions on neural networks and learning systems*, 30(11):3212–3232, 2019.
- [4] Biao Zhang, Deyi Xiong, Jun Xie, and Jinsong Su. Neural machine translation with gru-gated attention model. *IEEE transactions on neural networks and learning systems*, 31(11):4688–4698, 2020.
- [5] Min Yang, Chengming Li, Ying Shen, Qingyao Wu, Zhou Zhao, and Xiaojun Chen. Hierarchical human-like deep neural networks for abstractive text summarization. *IEEE Transactions on Neural Networks and Learning Systems*, 32(6):2744–2757, 2020.
- [6] Gongpei Zhao, Tao Wang, Yidong Li, Yi Jin, Congyan Lang, and Songhe Feng. Neighborhood pattern is crucial for graph convolutional networks performing node classification. *IEEE Transactions on Neural Networks and Learning Systems*, 2022.
- [7] Pierre Baldi and Kurt Hornik. Neural networks and principal component analysis: Learning from examples without local minima. *Neural networks*, 2(1):53–58, 1989.
- [8] Boris Hanin. Which neural net architectures give rise to exploding and vanishing gradients? *Advances in neural information processing systems*, 31, 2018.

- [9] Steve Lawrence and C Lee Giles. Overfitting and neural networks: conjugate gradient and backpropagation. In *Proceedings of the IEEE-INNS-ENNS International Joint Conference on Neural Networks. IJCNN 2000. Neural Computing: New Challenges and Perspectives for the New Millennium*, volume 1, pages 114–119. IEEE, 2000.
- [10] Yann N Dauphin, Razvan Pascanu, Caglar Gulcehre, Kyunghyun Cho, Surya Ganguli, and Yoshua Bengio. Identifying and attacking the saddle point problem in high-dimensional non-convex optimization. *Advances in neural information processing systems*, 27, 2014.
- [11] Timothy P Lillicrap, Daniel Cownden, Douglas B Tweed, and Colin J Akerman. Random synaptic feedback weights support error backpropagation for deep learning. *Nature communications*, 7(1):13276, 2016.
- [12] Arild Nøkland. Direct feedback alignment provides learning in deep neural networks. *Advances in neural information processing systems*, 29, 2016.
- [13] Charlotte Frenkel, Martin Lefebvre, and David Bol. Learning without feedback: Fixed random learning signals allow for feedforward training of deep neural networks. *Frontiers in neuroscience*, 15:629892, 2021.
- [14] Giorgia Dellaferera and Gabriel Kreiman. Error-driven input modulation: solving the credit assignment problem without a backward pass. In *International Conference on Machine Learning*, pages 4937–4955. PMLR, 2022.
- [15] Geoffrey Hinton. The forward-forward algorithm: Some preliminary investigations. *arXiv preprint arXiv:2212.13345*, 2022.
- [16] Julie M Clark, Francisco Cordero, Jim Cottrill, Bronislaw Czarnocha, David J DeVries, Denny St John, Georgia Tolias, and Draga Vidakovic. Constructing a schema: The case of the chain rule? *The Journal of Mathematical Behavior*, 16(4):345–364, 1997.
- [17] Nitish Srivastava, Geoffrey Hinton, Alex Krizhevsky, Ilya Sutskever, and Ruslan Salakhutdinov. Dropout: a simple way to prevent neural networks from overfitting. *The journal of machine learning research*, 15(1):1929–1958, 2014.
- [18] Zhun Zhong, Liang Zheng, Guoliang Kang, Shaozi Li, and Yi Yang. Random erasing data augmentation. In *Proceedings of the AAAI conference on artificial intelligence*, volume 34, pages 13001–13008, 2020.
- [19] Lutz Prechelt. Early stopping—but when? *Neural networks: tricks of the trade: second edition*, pages 53–67, 2012.
- [20] Xavier Glorot and Yoshua Bengio. Understanding the difficulty of training deep feedforward neural networks. In *Proceedings of the thirteenth international conference on artificial intelligence and statistics*, pages 249–256. JMLR Workshop and Conference Proceedings, 2010.
- [21] Jingzhao Zhang, Tianxing He, Suvrit Sra, and Ali Jadbabaie. Why gradient clipping accelerates training: A theoretical justification for adaptivity. *arXiv preprint arXiv:1905.11881*, 2019.
- [22] Jin Xu, Zishan Li, Bowen Du, Miaomiao Zhang, and Jing Liu. Reluplex made more practical: Leaky relu. In *2020 IEEE Symposium on Computers and communications (ISCC)*, pages 1–7. IEEE, 2020.
- [23] Kaiming He, Xiangyu Zhang, Shaoqing Ren, and Jian Sun. Deep residual learning for image recognition. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pages 770–778, 2016.
- [24] Sergey Ioffe and Christian Szegedy. Batch normalization: Accelerating deep network training by reducing internal covariate shift. In *International conference on machine learning*, pages 448–456. pmlr, 2015.
- [25] Sebahattin Bektaş and Yasemin Şişman. The comparison of l1 and l2-norm minimization methods. *International Journal of the Physical Sciences*, 5(11):1721–1727, 2010.
- [26] Diederik P Kingma and Jimmy Ba. Adam: A method for stochastic optimization. *arXiv preprint arXiv:1412.6980*, 2014.
- [27] Sebastian Ruder. An overview of gradient descent optimization algorithms. *arXiv preprint arXiv:1609.04747*, 2016.
- [28] Alexander Ororbias and Ankur Mali. The predictive forward-forward algorithm. *arXiv preprint arXiv:2301.01452*, 2023.
- [29] Stephen Grossberg. Competitive learning: From interactive activation to adaptive resonance. *Cognitive science*, 11(1):23–63, 1987.
- [30] Max Jaderberg, Wojciech Marian Czarnecki, Simon Osindero, Oriol Vinyals, Alex Graves, David Silver, and Koray Kavukcuoglu. Decoupled neural interfaces using synthetic gradients. In *International conference on machine learning*, pages 1627–1635. PMLR, 2017.

- [31] Kaiming He, Xiangyu Zhang, Shaoqing Ren, and Jian Sun. Delving deep into rectifiers: Surpassing human-level performance on imagenet classification. In *Proceedings of the IEEE international conference on computer vision*, pages 1026–1034, 2015.
- [32] Andre Martins and Ramon Astudillo. From softmax to sparsemax: A sparse model of attention and multi-label classification. In *International conference on machine learning*, pages 1614–1623. PMLR, 2016.
- [33] Leo Breiman. Bagging predictors. *Machine learning*, 24:123–140, 1996.
- [34] Avrim Blum and Tom Mitchell. Combining labeled and unlabeled data with co-training. In *Proceedings of the eleventh annual conference on Computational learning theory*, pages 92–100, 1998.
- [35] Yann LeCun, Léon Bottou, Yoshua Bengio, and Patrick Haffner. Gradient-based learning applied to document recognition. *Proceedings of the IEEE*, 86(11):2278–2324, 1998.
- [36] Alex Krizhevsky, Geoffrey Hinton, et al. Learning multiple layers of features from tiny images. 2009.
- [37] Oriol Vinyals, Charles Blundell, Timothy Lillicrap, Daan Wierstra, et al. Matching networks for one shot learning. *Advances in neural information processing systems*, 29, 2016.