

## Appendix

### A The Quantity Comparison Between Activations and Weights

As shown in Table 4,  $D_l/C_l$  is the size/channel of output of layer  $l$  and  $D_{l-1}/C_{l-1}$  is the size/channel of input to layer  $l$ ;  $k$  is the kernel size. Normally,  $D_l \approx D_{l-1}$  and  $D_l \gg k$ . Therefore, weights occupy greatly more memory overheads than the corresponding activations for FC; for LC, the weights also occupy more memory overheads than the corresponding activations because  $(C_l \times k \times k) > (C_{l-1} + C_l)$ . On the contrary, for CNN, the activations occupy greatly more memory overheads than the corresponding weights.

	Activations	Weights
FC	$D_l + D_{l-1}$	$D_l \times D_{l-1}$
LC	$C_l \times D_l + C_{l-1} \times D_{l-1}$	$C_l \times k \times k \times D_l$
CNN	$C_l \times D_l + C_{l-1} \times D_{l-1}$	$C_l \times k \times k$

Table 4: The quantity comparison between activations and weights for TinyFoA in various DNN architectures.

## B Equations for TinyFoA and Other Algorithms

Steps	Forward Pass	Gradient Calculation
BP (Rumelhart, Hinton, and Williams 1986)	$\mathbf{h}_l = \sigma_l(\mathbf{W}_l \mathbf{h}_{l-1} + \mathbf{b}_l)$	$\delta \mathbf{W}_l = \left[ (\mathbf{W}_{l+1}^T (\mathbf{W}_{l+2}^T \dots (\mathbf{W}_L^T \delta \mathbf{h}_L \odot \sigma'_L) \dots \odot \sigma'_{l+2}) \odot \sigma'_{l+1}) \right] \odot \sigma'_l \otimes \mathbf{h}_{l-1}^T$
DRTP (Frenkel, Lefebvre, and Bol 2021)	$\mathbf{h}_l = \sigma_l(\mathbf{W}_l \mathbf{h}_{l-1} + \mathbf{b}_l)$ $\mathbf{y}_{one-hot}$ is the one-hot-encoded label	$\delta \mathbf{W}_l = (\mathbf{B}_l^T \mathbf{y}_{one-hot}) \odot \sigma'_l \otimes \mathbf{h}_{l-1}^T$ where $\mathbf{B}_l$ is a fixed random connectivity matrix
PEPITA (Dellaferrera et al. 2022)	$\mathbf{h}_l = \sigma_l(\mathbf{W}_l \mathbf{h}_{l-1} + \mathbf{b}_l)$ $\mathbf{h}_l^{err} = \sigma_l(\mathbf{W}_l \mathbf{h}_{l-1}^{err} + \mathbf{b}_l)$	$\delta \mathbf{W}_l = (\mathbf{h}_l - \mathbf{h}_l^{err}) \otimes (\mathbf{h}_{l-1}^{err})^T$ , where $\mathbf{h}_0^{err} = x + Fe$ $e$ denotes the error and $F$ denotes the fixed random matrix
FF (Hinton 2022)	$\mathbf{h}_l^{pos} = \sigma_l(\mathbf{W}_l \mathbf{h}_{l-1}^{pos} + \mathbf{b}_l)$ $\mathbf{h}_l^{neg} = \sigma_l(\mathbf{W}_l \mathbf{h}_{l-1}^{neg} + \mathbf{b}_l)$ $g_l^{(pos)} = (\mathbf{h}_l^{(pos)})^T \mathbf{h}_l^{(pos)}$ , $g_l^{(neg)} = (\mathbf{h}_l^{(neg)})^T \mathbf{h}_l^{(neg)}$	$\delta \mathbf{W}_l = \mathbf{h}_l^{(pos)} \cdot \delta g_l^{(pos)} \odot \sigma'_l \otimes (\mathbf{h}_{l-1}^{(pos)})^T + \mathbf{h}_l^{(neg)} \cdot \delta g_l^{(neg)} \odot \sigma'_l \otimes (\mathbf{h}_{l-1}^{(neg)})^T$ $\delta g_l^{(pos)} = \text{Sigmoid}(g_l^{(pos)} - \theta) - 1$ , $\delta g_l^{(neg)} = \text{Sigmoid}(g_l^{(neg)} - \theta) - 0$ where $\theta$ is the threshold
TinyFoA-V-BA-BW	$\mathbf{h}_l = \sigma_l(\mathbf{W}_l \mathbf{h}_{l-1} + \mathbf{b}_l)$ $\mathbf{a}_l = \mathbf{B}_l \mathbf{h}_l$	$\delta \mathbf{W}_l = (\mathbf{B}_l^T \delta \mathbf{a}_l) \odot \sigma'_l \otimes \mathbf{h}_{l-1}^T$ where $\mathbf{B}_l$ is a fixed random matrix

Table 5: The formulation of forward pass and gradient calculation for various algorithms, where  $\sigma_l$  is the ReLU activation function,  $\mathbf{h}_0$  is the input  $\mathbf{x}$ .  $\odot$  is the element-wise product and  $\otimes$  is the Kronecker Product. For parameters update,  $\mathbf{W}_l = \mathbf{W}_l - \eta \delta \mathbf{W}_l$  and  $\mathbf{b}_l = \mathbf{b}_l - \eta \delta \mathbf{b}_l$ , where  $\eta$  is the learning rate.

## C The Derivation and Proof for Gradient and Parameters Update of TinyFoA

Let us consider a DNN with  $L$  layers. The input  $\mathbf{x} \in \mathbb{R}^{D_x \times 1}$  and the target  $\mathbf{y} \in \mathbb{R}^{N_c \times 1}$  are considered for training the DNN. The activations of the hidden layer  $l$  of the DNN are denoted as  $\mathbf{h}_l$ , where  $\mathbf{h}_0 = \mathbf{x}$ . For the hidden layer  $l$ , the activations  $\mathbf{h}_l \in \mathbb{R}^{D_l \times 1}$  based on binary activations  $\mathbf{h}_{l-1}^b \in \mathbb{R}^{D_{l-1} \times 1}$ , binary weights  $\mathbf{W}_l^b \in \mathbb{R}^{D_l \times D_{l-1}}$  and biases  $\mathbf{b}_l \in \mathbb{R}^{D_l \times 1}$ , are presented as follows:

$$\mathbf{h}_l = \sigma_l(\mathbf{W}_l^b \mathbf{h}_{l-1}^b + \mathbf{b}_l) = \sigma_l(\mathbf{z}_l). \quad (6)$$

For each  $\mathbf{h}_l^b$ , a fixed random matrix  $\mathbf{B}_l \in \mathbb{R}^{N_c \times D_l}$  is employed to project the  $\mathbf{h}_l^b \in \mathbb{R}^{D_l \times 1}$  to the output vector  $\mathbf{a}_l \in \mathbb{R}^{N_c \times 1}$ . Then, the projected vector  $\mathbf{a}_l$  is used to calculate the error and the gradient, as follows:

$$\begin{aligned} \mathbf{a}_l &= \mathbf{B}_l \mathbf{h}_l^b, \\ \mathbf{p}_l &= \sigma_{\text{output}}(\mathbf{a}_l), \\ \delta \mathbf{h}_l &= \frac{\partial \text{loss}_l(\mathbf{p}_l, \mathbf{y})}{\partial \mathbf{a}_l} \frac{\partial \mathbf{a}_l}{\partial \mathbf{h}_l^b} \frac{\partial \mathbf{h}_l^b}{\partial \mathbf{h}_l}, \end{aligned}$$

where  $\mathbf{p}_l \in \mathbb{R}^{N_c \times 1}$  is the probability distribution vector of  $\mathbf{a}_l$  and  $\sigma_{\text{output}}$  is the Softmax activation function.  $\text{loss}_l$  is the cross-entropy loss for the hidden layer  $l$ , and  $\delta \mathbf{a}_l = \frac{\partial \text{loss}_l(\mathbf{p}_l, \mathbf{y})}{\partial \mathbf{a}_l}$  is the gradient of  $\text{loss}_l$  with respect to  $\mathbf{a}_l$ .  $\delta \mathbf{h}_l \in \mathbb{R}^{D_l \times 1}$  is the gradient for updating weights  $\mathbf{W}_l^b$  and biases  $\mathbf{b}_l$ .

For  $\frac{\partial \text{loss}_l(\mathbf{p}_l, \mathbf{y})}{\partial \mathbf{a}_l}$ ,  $\delta \mathbf{a}_l = \mathbf{p}_l - \mathbf{y}$  due to the cross-entropy loss. Additionally,  $\frac{\partial \mathbf{a}_l}{\partial \mathbf{h}_l^b}$  is  $\mathbf{B}_l^T$  and  $\frac{\partial \mathbf{h}_l^b}{\partial \mathbf{h}_l}$  is  $\mathbf{1}_{|\mathbf{h}_l| \leq 1}$  (means output 1 if  $|\mathbf{h}_l| \leq 1$ ; otherwise output 0). Therefore, the gradient of  $\text{loss}_l(\mathbf{p}_l, \mathbf{y})$  with respect to  $\mathbf{h}_l$  is calculated as follows:

$$\begin{aligned} \delta \mathbf{h}_l &= \frac{\partial \text{loss}_l(\mathbf{p}_l, \mathbf{y})}{\partial \mathbf{a}_l} \frac{\partial \mathbf{a}_l}{\partial \mathbf{h}_l^b} \frac{\partial \mathbf{h}_l^b}{\partial \mathbf{h}_l} \\ &= \mathbf{B}_l^T (\mathbf{p}_l - \mathbf{y}) \odot \mathbf{1}_{|\mathbf{h}_l| \leq 1}. \end{aligned}$$

Moreover, considering the vertical layer-wise training, we only need to calculate the partial gradients  $\delta \mathbf{W}_l^{(i,j)}$  to reduce the training memory overheads. We divide the input activation  $\mathbf{h}_{l-1} \in \mathbb{R}^{D_{l-1} \times 1}$  and the output activation  $\mathbf{h}_l \in \mathbb{R}^{D_l \times 1}$  into  $M$  slices, respectively. The sliced input activations are denoted as  $\mathbf{h}_{l-1}^{(i)} \in \mathbb{R}^{\frac{D_{l-1}}{M} \times 1}$  and the sliced output activations are denoted as  $\mathbf{h}_l^{(j)} \in \mathbb{R}^{\frac{D_l}{M} \times 1}$ , where  $i = 1, 2, \dots, M$  and  $j = 1, 2, \dots, M$ . Therefore, the gradient of  $\text{loss}_l(\mathbf{p}_l, \mathbf{y})$  with respect to  $\mathbf{h}_l^{(j)}$  is calculated as follows:

$$\begin{aligned} \delta \mathbf{h}_l^{(j)} &= \frac{\partial \text{loss}_l(\mathbf{p}_l, \mathbf{y})}{\partial \mathbf{a}_l} \frac{\partial \mathbf{a}_l}{\partial \mathbf{h}_l^{b,(j)}} \frac{\partial \mathbf{h}_l^{b,(j)}}{\partial \mathbf{h}_l^{(j)}} \\ &= (\mathbf{B}_l^{(j)})^T (\mathbf{p}_l - \mathbf{y}) \odot \mathbf{1}_{|\mathbf{h}_l^{(j)}| \leq 1}, \end{aligned}$$

where  $\mathbf{B}_l^{(j)} \in \mathbb{R}^{N_c \times \frac{D_l}{M}}$ .  $\mathbf{1}_{|\mathbf{h}_l^{(j)}| \leq 1} = 1$  if  $|\mathbf{h}_l^{(j)}| \leq 1$  and  $\mathbf{1}_{|\mathbf{h}_l^{(j)}| \leq 1} = 0$  otherwise.

Furthermore, we introduce the derivation for the sliced parameter (weights and biases) updates. The gradient of  $\text{loss}_l(\mathbf{p}_l, \mathbf{y})$  with respect to  $\mathbf{W}_l^{b,(i,j)}$  and  $\mathbf{b}_l^{(j)}$  are calculated as follows:

$$\begin{aligned} \delta \mathbf{W}_l^{b,(i,j)} &= \delta \mathbf{h}_l^{(j)} \frac{\partial \mathbf{h}_l^{(j)}}{\partial \mathbf{z}_l^{(j)}} \frac{\partial \mathbf{z}_l^{(j)}}{\partial \mathbf{W}_l^{b,(i,j)}} \\ &= (\mathbf{B}_l^{(j)})^T (\mathbf{p}_l - \mathbf{y}) \odot \mathbf{1}_{|\mathbf{h}_l^{(j)}| \leq 1} \odot \sigma_l'(\mathbf{z}_l^{(j)}) \otimes (\mathbf{h}_{l-1}^{b,(i)})^T \\ \delta \mathbf{b}_l^{(j)} &= \delta \mathbf{h}_l^{(j)} \frac{\partial \mathbf{h}_l^{(j)}}{\partial \mathbf{z}_l^{(j)}} \frac{\partial \mathbf{z}_l^{(j)}}{\partial \mathbf{b}_l^{(j)}} \\ &= (\mathbf{B}_l^{(j)})^T (\mathbf{p}_l - \mathbf{y}) \odot \mathbf{1}_{|\mathbf{h}_l^{(j)}| \leq 1} \odot \sigma_l'(\mathbf{z}_l^{(j)}) \cdot 1, \end{aligned}$$

where  $\odot$  is the Hadamard Product, and  $\otimes$  is the Kronecker Product.  $\sigma_l'(\mathbf{z}_l^{(j)}) = \mathbf{1}_{\mathbf{z}_l^{(j)} > 0}$  because  $\sigma_l$  is ReLU activation function here.

Finally, the gradient of  $\text{loss}_l(\mathbf{p}_l, \mathbf{y})$  with respect to  $\mathbf{W}_l^{(i,j)}$  is denoted as follow:

$$\begin{aligned} \delta \mathbf{W}_l^{(i,j)} &= \delta \mathbf{W}_l^{b,(i,j)} \odot \mathbf{1}_{|\mathbf{W}_l^{(i,j)}| \leq 1} \\ &= \delta \mathbf{h}_l^{(j)} \odot \mathbf{1}_{\mathbf{z}_l^{(j)} > 0} \otimes (\mathbf{h}_{l-1}^{b,(i)})^T \odot \mathbf{1}_{|\mathbf{W}_l^{(i,j)}| \leq 1} \end{aligned}$$

## D Hyperparameter Details

In this section, we present the details of hyperparameters, including input size, category, batch size, epochs, learning rate, optimizer, the length of hidden layers, the activation function of hidden layers, and loss function, for all the algorithms and various datasets evaluated in this paper.

Table 6 demonstrates the hyperparameter details for all the algorithms, which means these hyperparameters apply to all four datasets in our experiments. Moreover, Table 7 demonstrates the hyperparameters details for these algorithms with respect to various datasets. Table 8 demonstrates the hyperparameters details for TinyFoA-LC.

Algorithms	Optimizer	Length/number of Hidden Layers	Activation Function	Loss	$\eta$
BP (Rumelhart, Hinton, and Williams 1986)	Adam	2000/4	ReLU	CCE	1e-4
DRTP (Frenkel, Lefebvre, and Bol 2021)	NAG	2000/4	Tanh	BCE	1e-4
PEPITA (Dellaferrera et al. 2022)	Mom	2000/3	ReLU	CCE	1e-3/1e-4
FF (Hinton 2022)	SGD	2000/4	ReLU	BCE	1e-3
TinyFoA-FC	Adam	2000/4	ReLU	CCE	1e-4

Table 6: The hyperparameter details for all the algorithms.

Hyperparameters	MNIST (LeCun 1998)	CIFAR-10 (Krizhevsky 2009)	CIFAR-100 (Krizhevsky 2009)	MIT-BIH (Mark et al. 1982)
Input Size	28x28x1	32x32x3	32x32x3	13x13x1
Category	10	10	100	5
Batchsize	100	100	100	100
Epochs	100	100	100	100
$\eta$ of PEPITA	1e-3	1e-3	1e-3	1e-4

Table 7: The hyperparameter details for all the algorithms with respect to various datasets.

Hyperparameters	Channel	kernel	Stride	Epochs	Optimizer	Activation Function	Loss	$\eta$	Batchsize
TinyFoA-LC	[16,32,64,64]	3	1	200	Adam	ReLU	CCE	1e-4	100

Table 8: The hyperparameter details for TinyFoA-LC.