# Lecture 11:
# Hardware Accelerator for DNN Training

# Notes

- Project meeting, please sign up!
- Final project presentation
  - I will send out a form for you to select the date and time options.
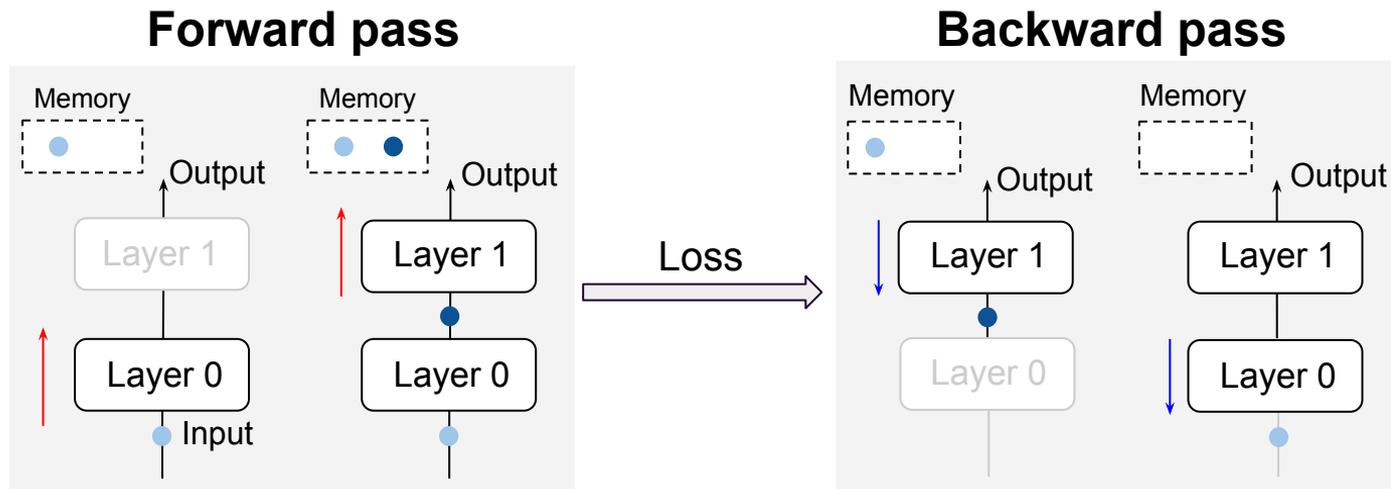
NYU SAI LAB

# Recap

- Matrix Multiplication with Transposition
- Hardware design for Nonlinear Blocks
- System optimization of LLMs
- Popular transformer accelerator design
  - SpAtten
  - EdgeBert
  - Olive

NYU SAI LAB
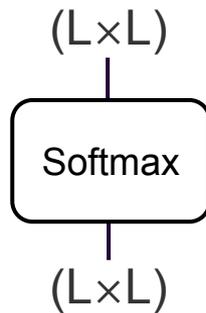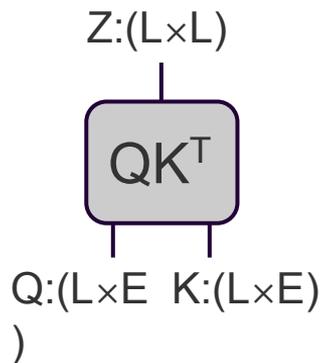
# Topics
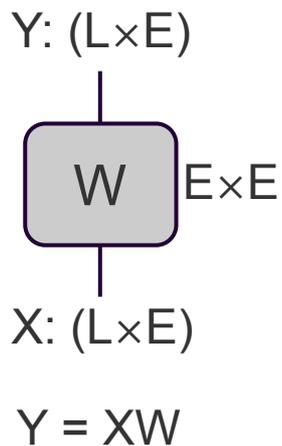
- <span style="color:red">Computation during backward propagation</span>
- Hardware architecture for backward propagation design
- Popular DNN training accelerator design
  - FAST
  - CAMEL
  - ADA-GP
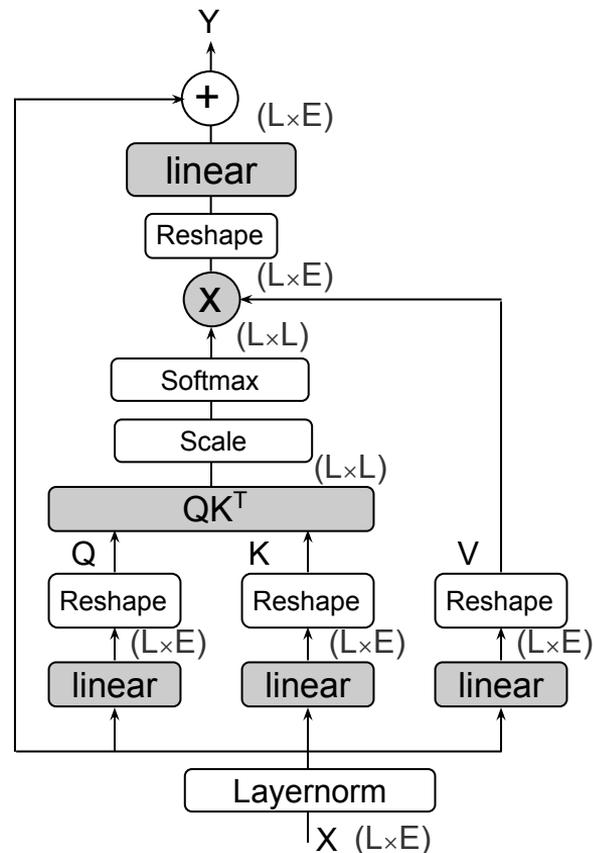  - Procrustes

# Neural Network Training



- The peak memory grow linearly as the layer depth increases.
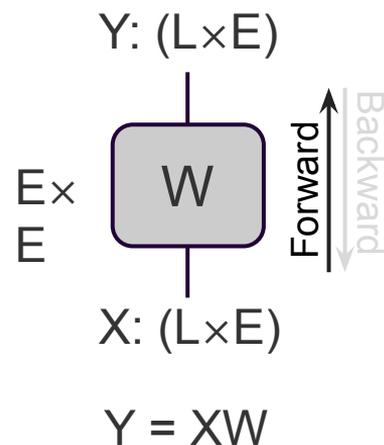- The backward propagation involves more computations

# Training Flow

Y: (L×E)

W $E×E$

X: (L×E)

Y = XW

Z:(L×L)

$QK^T$

Q:(L×E)  K:(L×E)
)

(L×L)

Softmax

(L×L)
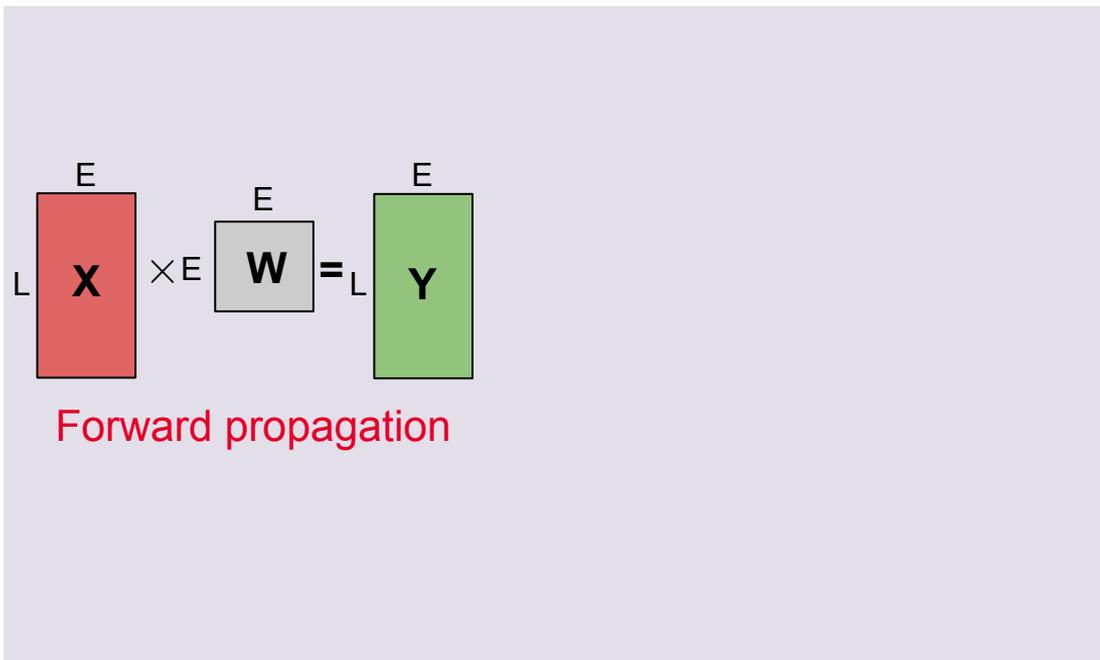
L: number of tokens
E: embedding dimension



NYU SAI LAB

6

# Training Flow: Linear Layer

$$E$$
$$E$$
$$E$$

$L$   **X**   $\times E$   **W**   =   $L$   **Y**

Forward propagation

Y: (L×E)

E×
E    **W**    Forward / Backward

X: (L×E)

Y = XW

# Training Flow: Linear Layer



Forward propagation

Backward propagation: weight gradient computation

Backward propagation: input gradient computation

Y: (L×E)

X: (L×E)

Y = XW

# Training Flow: Self-Attention



Forward propagation

Backward propagation

# Training Flow: Softmax

$$s_i = \frac{e^{z_i}}{\sum_{j=0}^{N-1} e^{z_j}} \; For \; i = 1, 2, \cdots, N$$

Forward propagation

$$\frac{ds}{dz} = diag(s) - ss^T$$

$$\frac{ds}{dz} = \begin{bmatrix} s_1 - s_1^2 & -s_1 \cdot s_2 & -s_1 \cdot s_3 \\ -s_2 \cdot s_1 & s_2 - s_2^2 & -s_2 \cdot s_3 \\ -s_3 \cdot s_1 & -s_3 \cdot s_2 & s_3 - s_3^2 \end{bmatrix}$$

Backward propagation

S: L×L

Softmax

Z: L×L

Forward

Backward

NYU SAI LAB

# Training Flow: Normalization

**Forward propagation:**

$$s = \alpha \frac{z - \mu_z}{\sigma_z} + \beta$$
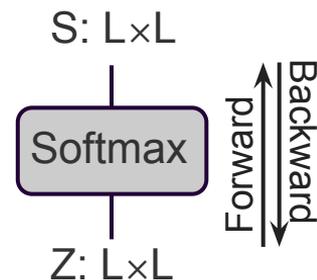
$$s_i = \alpha \bar{z}_i + \beta$$

$$\mu_z = \frac{\sum_i z_i}{N}$$

$$\sigma_z = \sqrt{\frac{\sum_i (z_i - \mu_z)^2}{N}}$$

Forward propagation

**Backward propagation:**

$$\frac{dL}{d\beta} = \sum_i \frac{dL}{ds_i}$$

$$\frac{dL}{d\alpha} = \sum_i \frac{dL}{ds_i} \bar{z}_i$$

Backward propagation

LayerNorm — Forward / Backward

# Topics

- Computation during backward propagation
- <span style="color:red">Hardware architecture for backward propagation design</span>
- Popular DNN training accelerator design
  - FAST
  - CAMEL
  - ADA-GP
  - Procrustes

NYU SAI LAB

# In-place Transposed Matrix Multiplication



Forward propagation

Backward propagation: weight gradient computation

Backward propagation: input gradient computation

# Forward Pass for Convolutional Layer

**Convolution View**

**Matrix View**

Forward Pass
Compute output **Y**



- Assume a weight kernel size of 1✶1.

# Backward Pass for Convolutional Layer

Backward Pass
Compute Activation gradients ∇**X**

# Backward Pass for Convolutional Layer

Backward Pass
Compute weight gradients $\nabla \mathbf{W}$

# In-place Transposed Matrix Multiplication

- In the training of neural networks, we need to perform transposed matrix multiplication
- Instead of using a separate hardware for matrix transposition, transposed matrix multiplication can be performed using a systolic array.

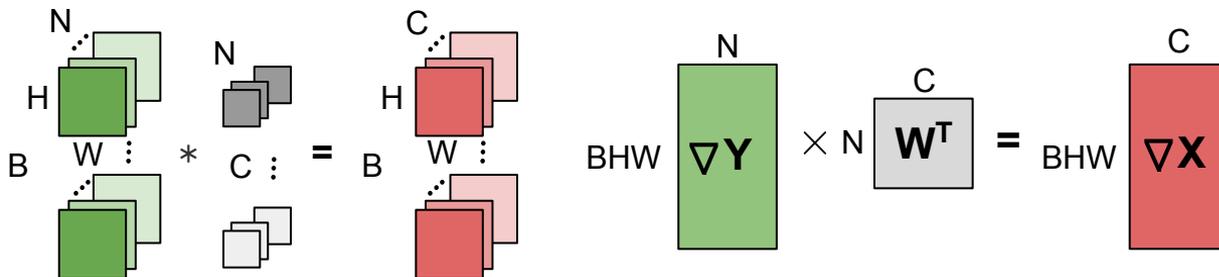$$\text{BHW} \begin{bmatrix} X \end{bmatrix}^C \times {}^C\begin{bmatrix} W \end{bmatrix}^N = \text{BHW} \begin{bmatrix} Y \end{bmatrix}^N$$

$$\text{BHW} \begin{bmatrix} \nabla Y \end{bmatrix}^N \times {}^N\begin{bmatrix} W^T \end{bmatrix}^C = \text{BHW} \begin{bmatrix} \nabla X \end{bmatrix}^C$$

$$^C\begin{bmatrix} X^T \end{bmatrix}^{BHW} \times \text{BHW}\begin{bmatrix} \nabla Y \end{bmatrix}^N = {}^C\begin{bmatrix} \nabla W \end{bmatrix}^N$$

NYU SAI LAB

# In-place Transposed Matrix Multiplication

$$\begin{bmatrix} 1 & 4 \\ 5 & 2 \end{bmatrix} \begin{bmatrix} 2 & 3 \\ 0 & 1 \end{bmatrix} = \begin{bmatrix} 2 & 7 \\ 10 & 17 \end{bmatrix}$$

**X**    **W**    **Y**

- Weight stationary, input from bottom, accumulation from left to right



- The weights are preloaded into the systolic array, while the input matrix is streamed into the array from bottom to top.
- The output is produced at the right.

18

# In-place Transposed Matrix Multiplication

$$\begin{bmatrix} 3 & 4 \\ 1 & 2 \end{bmatrix} \begin{bmatrix} 2 & 0 \\ 3 & 1 \end{bmatrix} = \begin{bmatrix} 18 & 4 \\ 8 & 2 \end{bmatrix}$$

$\nabla Y \qquad W^T \qquad \nabla X$

- Weight stationary, input from left, accumulation upwards

- To compute the input gradient, the data gradient is fed into the systolic array from the left, and the output is produced at the top.

# Systolic Array: Weight-Stationary Version

$$z = w \cdot x + y$$
$$v = x$$

(diagram: cell with inputs y (left), x (bottom), outputs z (right), v (top))

- Takes data (x and y) as input
- w stays in the systolic cell
- Performs a multiply-accumulate operation

# Systolic Array: Accumulation-Stationary Version



z = y
v = x
q = x·y + q

- Takes data (x and y) as input
- Accumulated result q stays in the systolic cell
- Performs a multiply-accumulate operation

# Systolic Cell



v

y=2 →

z = y
v = x
q = x·y + q

→ z

↑ x=4

q is stored in the MAC

NYU SAI LAB

# Systolic Cell



v

y=2 →

z = y
v = x
q = x·y + 0

→ z

↑ x=4

q is stored in the MAC

# Systolic Cell



q is stored in the MAC

# Systolic Cell

# In-place Transposed Matrix Multiplication

$$\begin{bmatrix} 1 & 5 \\ 4 & 2 \end{bmatrix} \begin{bmatrix} 3 & 4 \\ 1 & 2 \end{bmatrix} = \begin{bmatrix} 8 & 14 \\ 14 & 20 \end{bmatrix}$$
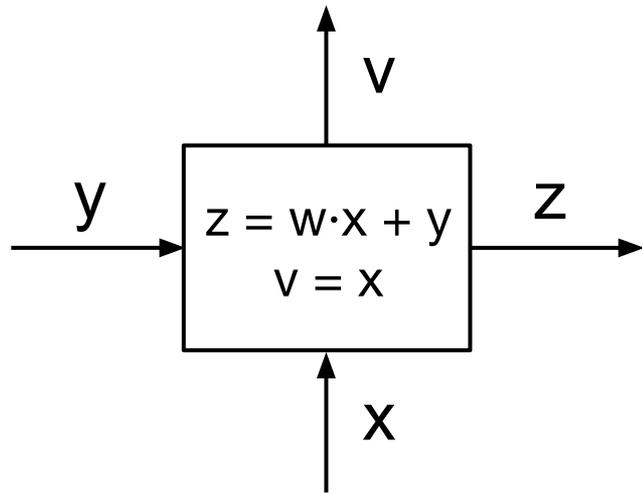
$\mathbf{X^T}$    $\nabla\,\mathbf{Y}$    $\nabla\mathbf{W}$

- Input from left and bottom, accumulation stationary.



- To compute the weight gradient, the data gradient is input from the left side of the systolic array, while the input activations are fed from the bottom. The resulting weight gradients are accumulated and stored within the systolic cells.

26

# Topics

- Computation during backward propagation
- Hardware architecture for backward propagation design
- Popular dnn training accelerator design
  - FAST
  - CAMEL
  - Tensordash
  - Procrustes

# FAST



- The Dot product is performed using BFP format.
- Compare with floating point (FP) format, block floating point (BFP) perform exponent additions and mantissa alignments only at the group level, rather than at the individual elements level.

# Fast First, Accurate Second Training

- Previous literature has demonstrated that adding zero-mean Gaussian noise to the weight gradient $\nabla W$ can reduce overfitting and improve the convergence of training.

- Decreasing the variance of the noise over iterations achieves better performance than using fixed Gaussian noise throughout training.

- We hypothesize that a similar effect can be achieved by adjusting the BFP precision of weights, activations, and gradients from low to high precisions over training.

Zhang, Sai Qian, Bradley McDanel, and H. T. Kung. "Fast: Dnn training under variable precision block floating point with stochastic rounding." *2022 IEEE International Symposium on High-Performance Computer Architecture (HPCA)*. IEEE, 2022.

NYU SAI LAB

# Variable Precision Training

Temporal precision change



Layerwise precision change



- Under the Temporal High-to-Low scheme, we use FP32 for weights, activations, and gradients for the first part of training, and lower-precision BFP for the second part of training

- Under the Layerwise High-to-Low scheme, we use FP32 precision for the first ten layers, and lower-precision BFP for later layers

# Variable Precision Training

**BFP precision increases across layer depth and training iterations**



- We progressively increase the BFP precision of weights, activations, and gradients along both layer depth and training iterations

- We name this approach **FAST** (Fast First, Accurate Second Training)

# FAST System Design



- Major components of FAST system:
    - Systolic array with FAST multiplier and accumulator (fMAC)
    - BFP converter
    - Accumulator and systolic array input generator
    - Memory subsystem

# FAST System Design



- fMAC operates on chunks of BFP mantissas (e.g., 2-bit chunks) to support variable-width mantissas in 2-bit increments

# Evaluation



- FAST progressively increases the BFP precision across both layer depth and iterations during the training process

# Evaluation



ResNet-18 (ImageNet) TTA of 68%

- We use Time-to-Accuracy (TTA) as the evaluation metric to compare different approaches
- Our FAST approach achieves the lowest TTA across all the numeric formats

# Topics

- Computation during backward propagation
- Hardware architecture for backward propagation design
- <span style="color:red">Popular dnn training accelerator design</span>
  - FAST
  - <span style="color:red">CAMEL</span>
  - Tensordash
  - Procrustes

# Memory Efficient Neural Network Training



**Forward pass**

**Backward pass**

Normalized Number of Parameters
Storage during DNN Training

- The memory footprint grows proportional with the layer depth.
- On top of this, small edge devices typically have limited on-chip storage, leading to frequent and costly accesses to off-chip memories.

Zhang, Sai Qian, et al. "CAMEL: Co-Designing AI Models and eDRAMs for Efficient On-Device Learning." *2024 IEEE International Symposium on High-Performance Computer Architecture (HPCA)*. IEEE, 2024.

# Memory Efficient Neural Network Training

**Residual Architecture**

**Reversible Architecture**



**Forward pass:**
$y_2 = F_1(x_1) + x_2$
$y_1 = F_2(y_2) + x_1$

**Backward Pass:**
$x_1 = y_1 - F_2(y_2)$
$x_2 = y_2 - F_1(x_1)$

- A reversible residual network (RevNet) is a variant of the canonical residual neural network (ResNet).

Zhang, Sai Qian, et al. "CAMEL: Co-Designing AI Models and eDRAMs for Efficient On-Device Learning." *2024 IEEE International Symposium on High-Performance Computer Architecture (HPCA)*. IEEE, 2024.

# Memory Efficient Neural Network Training

### Architecture



**Operations**

**Forward pass:**
$y_2 = F_1(x_1) + x_2$
$y_1 = F_2(y_2) + x_1$

**Backward Pass:**
$x_1 = y_1 - F_2(y_2)$
$x_2 = y_2 - F_1(x_1)$

**1. Recompute the input**

**2. Compute input gradient**

**3. Compute weight gradient and update**

- The reversible architecture enables the backward pass computations to be performed without the need to store the input activations.
- Given the output y, the input activations are first recomputed. Afterwards, the input and weight gradients are computed with standard backward pass operations.

Zhang, Sai Qian, et al. "CAMEL: Co-Designing AI Models and eDRAMs for Efficient On-Device Learning." *2024 IEEE International Symposium on High-Performance Computer Architecture (HPCA)*. IEEE, 2024.

# Memory Efficient Neural Network Training



Duplex DNN

- This approach in turn imposes higher compute demands.

- We propose to judiciously train a subset of the model parameters to minimize training.

- The backbone DNN is frozen during the backward pass of the DNN.

- The normalization layers are removed from the branch DNN to facilitate the training process.

Zhang, Sai Qian, et al. "CAMEL: Co-Designing AI Models and eDRAMs for Efficient On-Device Learning." *2024 IEEE International Symposium on High-Performance Computer Architecture (HPCA)*. IEEE, 2024.

# Procrustes

**init:** $W^{(0)}$ with $W^{(0)} \sim N(0, \sigma)$

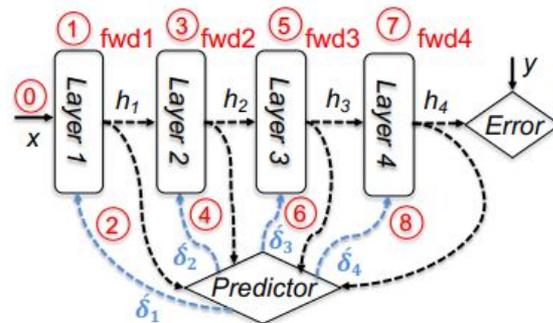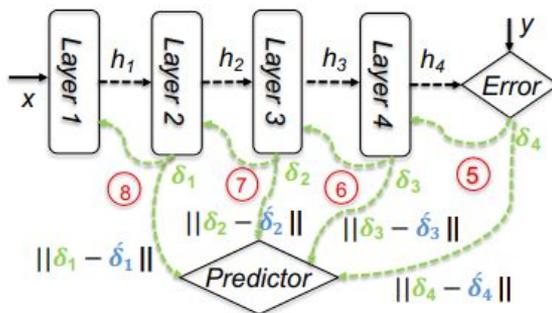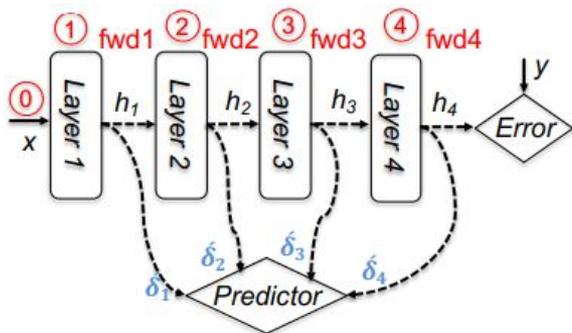**output:** $W^{(t)}$

**while** *not converged* **do**

$$T = \left\{ \left| \sum_{i=0}^{t-1} \frac{\eta \partial f(W^{(i-1)}; x^{(i-1)})}{\partial w} \right| \text{ s.t. } w \in W_{trk} \right\}$$

$$P = \left\{ \left| \frac{\eta \partial f(W^{(i-1)}; x^{(i-1)})}{\partial w} \right| \text{ s.t. } w \in W_{prn} \right\}$$

$$S = \text{sort}(T \cup P)$$

$$mask = \mathbb{1}(S > S[k])$$

$$W^{(t)} =$$
$$mask \cdot \left( W^{(t-1)} - \eta \nabla f \left( W^{(t-1)}; x^{(t-1)} \right) \right) + \overline{mask} \cdot W^{(0)}$$

$$t = t + 1$$

- We adapt a sparse training algorithm to be amenable to hardware acceleration; we then develop dataflow, data layout, and load balancing techniques to accelerate it.
- Only a fixed percentage of the parameters (e.g., 10%) are ever allowed to change
- During the each training iteration, the weights with the highest accumulated gradient survive

Yang, Dingqing, et al. "Procrustes: a dataflow and accelerator for sparse deep neural network training." *2020 53rd Annual IEEE/ACM International Symposium on Microarchitecture (MICRO)*. IEEE, 2020.

NYU SAI LAB

# ADA-GP



- We propose ADA-GP, which uses gradient prediction adaptively to speed up DNN training while maintaining accuracy.
- ADA-GP works by incorporating a small neural network to predict gradients for different layers of a DNN model.

Janfaza, Vahid, et al. "ADA-GP: Accelerating DNN Training By Adaptive Gradient Prediction." *Proceedings of the 56th Annual IEEE/ACM International Symposium on Microarchitecture*. 2023.

# Presentations

- On-Device Training Under 256KB Memory (Ankit)
- FlashDecoding++: Faster Large Language Model Inference on GPUs (Akshay, Su)
- CAMEL: Co-Designing AI Models and eDRAMs for Efficient On-Device Learning (Jahnavi, Isha)
- Procrustes: a Dataflow and Accelerator for Sparse Deep Neural Network Training (Mohnish, Dae Sung)