



Lecture 05: DNN Quantization

Recap

- Why pruning?
 - Running cost of CNNs and Transformers
- Sparse matrix encoding
- General pruning techniques
- Transformer pruning
- Large model pruning

Topics

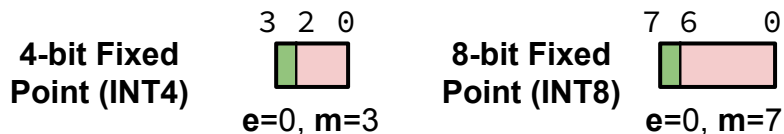
- Basic Data Formats
 - Fixed point (INT)
 - Floating point (FP)
 - Block floating point (BFP)
- Quantization methods
 - Taxonomy of Quantization
 - Learnable adaptive quantization scheme
 - Quantization for LLM

Topics

- **Basic Data Formats**
 - Fixed point (INT)
 - Floating point (FP)
 - Block floating point (BFP)
- **Quantization methods**
 - Taxonomy of Quantization
 - Learnable adaptive quantization scheme
 - Quantization for LLM

Fixed-Point Arithmetic (INT)

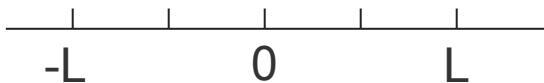
Fixed Point Formats



- Hyperparameter associated with the fixed-point format:
 - Clipping range $(-L, L)$: usually symmetrical around 0
 - Bitwidth (b)
- Quantization with Fixed-point format is called **Fixed point quantization** or **INT quantization**.

Fixed-Point Format (Symmetrical)

- How to convert a number x to INT representation?
 - Set the clipping range: $(-L, L)$, bitwidth: b
 - Compute the scale: $s = 2L / (2^b - 2)$
 - Clip the input x : $x_c = \text{Clip}(x, L, -L)$
 - Calculate the INT representation: $x_{int} = \text{round}(x_c / s)$
 - Rescale: $x_q = Sx_{int}$
- Have a uniform representation power within the clipping range.
- All the computations can be performed using x_{int}



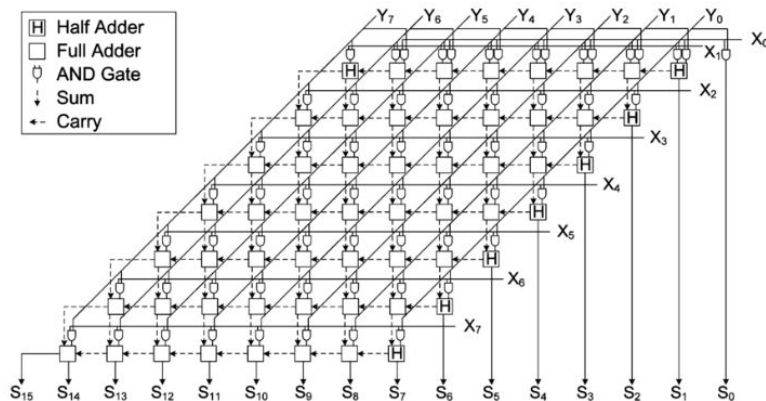
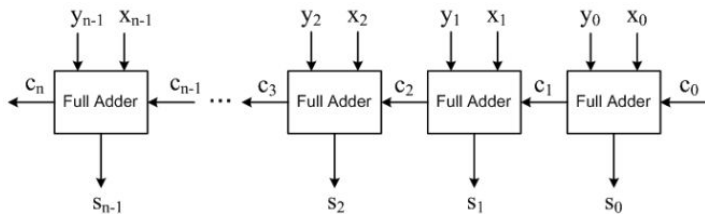
Example

- $X = [1.1, 2.4, -0.3, 0.8]$, bitwidth = 3, $L = 2$
- How to convert a number x to INT representation?
 - Set the clipping range: $(-L, L)$, bitwidth: b $b=3, L=2$
 - Compute the scale: $s = 2L / (2^b - 2)$ $s = 4/6 = 2/3$
 - Clip the input x : $x_c = \text{Clip}(x, L, -L)$ $X_c = [1.1, 2, -0.3, 0.8]$
 - Calculate the INT representation: $x_{int} = \text{round}(x_c / s)$ $X_{int} = [2, 3, 0, 1]$
 - Rescale: $x_q = s x_{int}$ $X_q = [1.33, 2.0, 0.0, 0.67]$

Computation with Fixed-Point Format

- Addition/Subtraction: $x_q \pm y_q = s(x_{int} \pm y_{int})$
- Multiplication: $x_q \times y_q = s^2(x_{int} \times y_{int})$
- Division: $x_q/y_q = x_{int}/y_{int}$

If the scales are the same



Computation with Fixed-Point Format

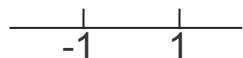
- If we try to compute the matrix multiplication between X and Y:

$$\begin{array}{|c|c|} \hline x_{q,1} & x_{q,2} \\ \hline \end{array} \times \begin{array}{|c|} \hline y_{q,1} \\ \hline y_{q,2} \\ \hline \end{array}$$

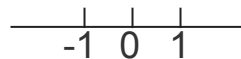
All elements within the tensors are quantized using the same scale

$$x_{q,1} \times y_{q,1} + x_{q,2} \times y_{q,2} = s_x s_y (x_{int,1} \times y_{int,1} + x_{int,2} \times y_{int,2})$$

INT Quantization with Low Precision



Binary neural network



Ternary neural network

- Binary and Ternary neural networks are both multiplication-free DNN.

Fixed Point Format (Unsymmetrical)

- How to convert a number to INT8 representation?
 - Set the clipping range: $(-L, L)$, bitwidth: b
 - Compute the scale: $s = (L_{max} - L_{min}) / (2^b - 1)$
 - Clip the input x : $x_c = \text{Clip}(x, L_{min}, L_{max})$
 - Calculate the fixed-point representation:
$$x_{int} = \text{round}((x_c - L_{min}) / s)$$
 - Rescale: $x_q = s x_{int} + L_{min}$

Example

- $X = [1.1, 2.4, -0.3, 0.8]$, bitwidth = 3, $L = 2$
- How to convert a number to INT8 representation?
 - Set the clipping range: $(-L, L)$, bitwidth: b $b=3, L_{max}=2, L_{min}=-0.5$
 - Compute the scale: $s = (L_{max} - L_{min}) / (2^b - 1)$ $s = 0.357$
 - Clip the input x : $x_c = Clip(x, L_{min}, L_{max})$ $X_c = [1.1, 2, -0.3, 0.8]$
 - Calculate the fixed-point representation:
 $x_{int} = round((x_c - L_{min}) / s)$ $X_{int} = [4, 7, 1, 4]$
 - Rescale: $x_q = sx_{int} + L_{min}$ $X_q = [0.93, 2.0, -0.14, 0.93]$

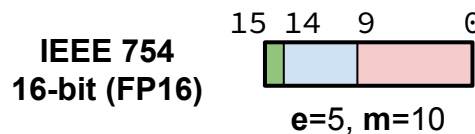
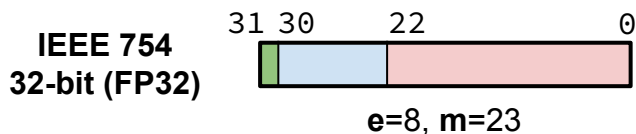
Computation with Fixed-Point Format

- Addition/Subtraction: hard to implement
- Multiplication (needs additional computation):

$$\mathbf{x}_q \times \mathbf{y}_q = \mathbf{s}_x \mathbf{s}_y (\mathbf{x}_{int} \times \mathbf{y}_{int}) + L_{min,x} \mathbf{y}_q \mathbf{s}_y + L_{min,y} \mathbf{x}_q \mathbf{s}_x + L_{min,x} L_{min,y}$$

- Division: hard to implement

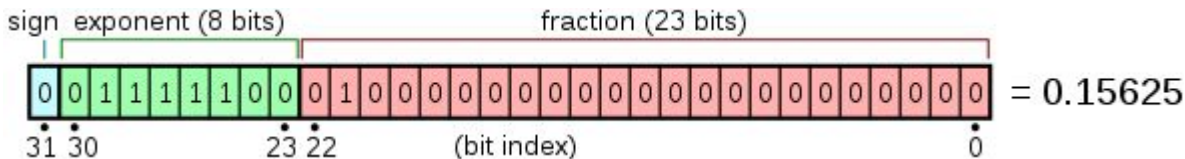
Floating-Point Arithmetic



■ Sign field ■ Exponent (e) ■ Mantissa (m)

- The floating-point number has three fields:
 - Sign (s)
 - Exponent (e)
 - Mantissa (m)

Floating-Point Arithmetic



- Every real number can be converted in the following format:

$$x = (-1)^s \times 2^{e-bias} \times m \quad \text{where } 1 \leq m < 2$$

$$m = (1.b_0b_1b_2\dots b_{22})_2$$

There typically exists a predefined bias: bias = 127 for IEEE 754 FP32.

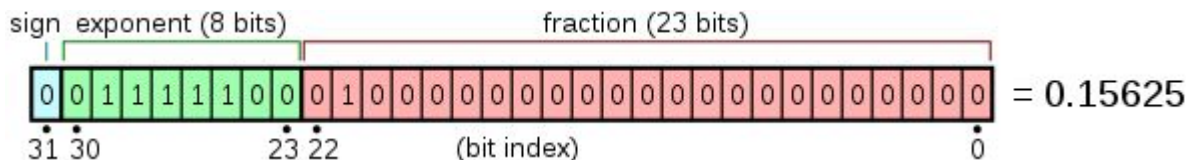
- For example:

- $5.5 = (-1)^0 \times 2^{129-127} \times (1.011)_2$ $s = 0, e = 129, m = 011$

- $-71 = (-1)^1 \times 2^{133-127} \times (1.000111)_2$

- $0.34375 = (-1)^1 \times 2^{125-127} \times (1.011)_2$

Floating-Point Arithmetic



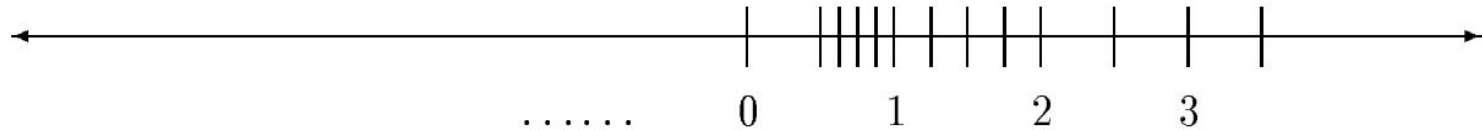
- IEEE-754 standard:

$$x = (-1)^s \times 2^{e-bias} \times m \quad \text{where } 1 \leq m < 2$$

$$m = (1.b_0b_1b_2\dots b_{22})_2$$

- The exponent field is unsigned.
- We need some special representation:
 - A bit stream of all zeros represents 0

Floating Point Arithmetic



- Have better representation power for values with small magnitudes.
- How to convert a real number x to FP representation?

$$x = |x| \quad s = \text{sign}(x)$$

$$a = \lfloor \log_2 x \rfloor \quad e = a + \text{bias} \quad m = \frac{x}{2^a} - 1$$

Example

$x = -13.24$, bias=127

$x = |x|$ $s = \text{sign}(x)$

$$a = \lfloor \log_2 x \rfloor \quad e = a + \text{bias} \quad m = \frac{x}{2^a} - 1$$

$a = 3$, $e = 130$, $m = 0.655$

$s = (0)_2$, $e = (1000010)_2$, $m = (0.101001111010111000010000)_2$

Computation with FP Representation

- Addition/Subtraction:

- Need to align the exponent

$$\begin{array}{ccccccc} 011010 & + & 001111 & = & 011010 & + & 011001 & = & 011011 \\ \underbrace{\quad\quad\quad} & & \underbrace{\quad\quad\quad} & & & & & & \\ s_1 e_1 m_1 & & s_2 e_2 m_2 & & & & & & \end{array} \quad \text{Alignment}$$

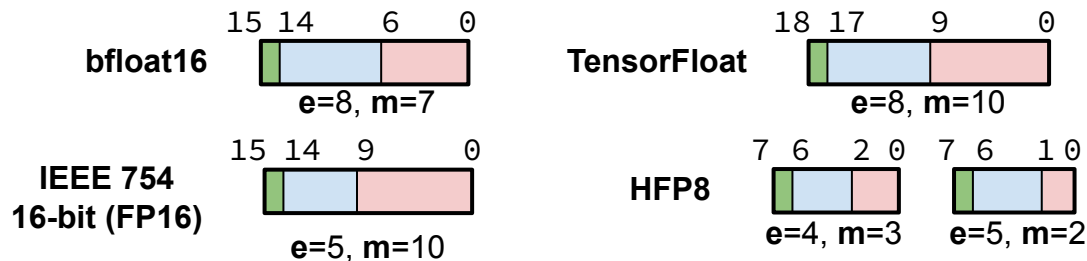
- Multiplication/Subtraction:

- Sum the exponent, multiply the mantissa

$$\begin{array}{ccccccc} 011010 & \times & 001111 & & e = e_1 + e_2 \\ \underbrace{\quad\quad\quad} & & \underbrace{\quad\quad\quad} & & \\ s_1 e_1 m_1 & & s_2 e_2 m_2 & & m = 1.m_1 \times 1.m_2 \end{array}$$

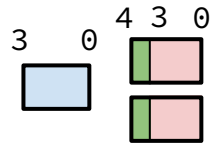
- Addition and subtraction is expensive for FP.

Customized FP Representation

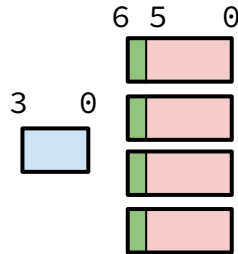


- Numerous customized FP representations have been developed to facilitate DNN execution.

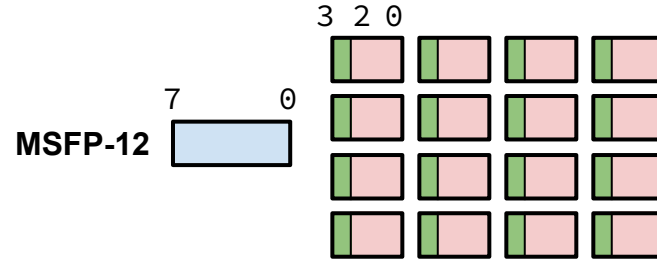
Block Floating Point (BFP)



$g=2, e=4, m=4$



$g=4, e=4, m=6$



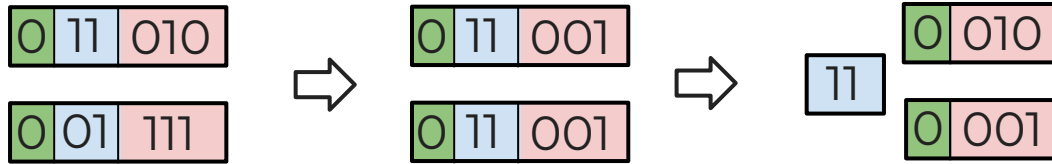
MSFP-12

$g=16, e=8, m=3$

■ Sign field ■ Exponent (e) ■ Mantissa (m)

- BFP formats offer a middle ground between FP and INT formats, by enforcing that a group of values share a common exponent while maintaining individual mantissas.

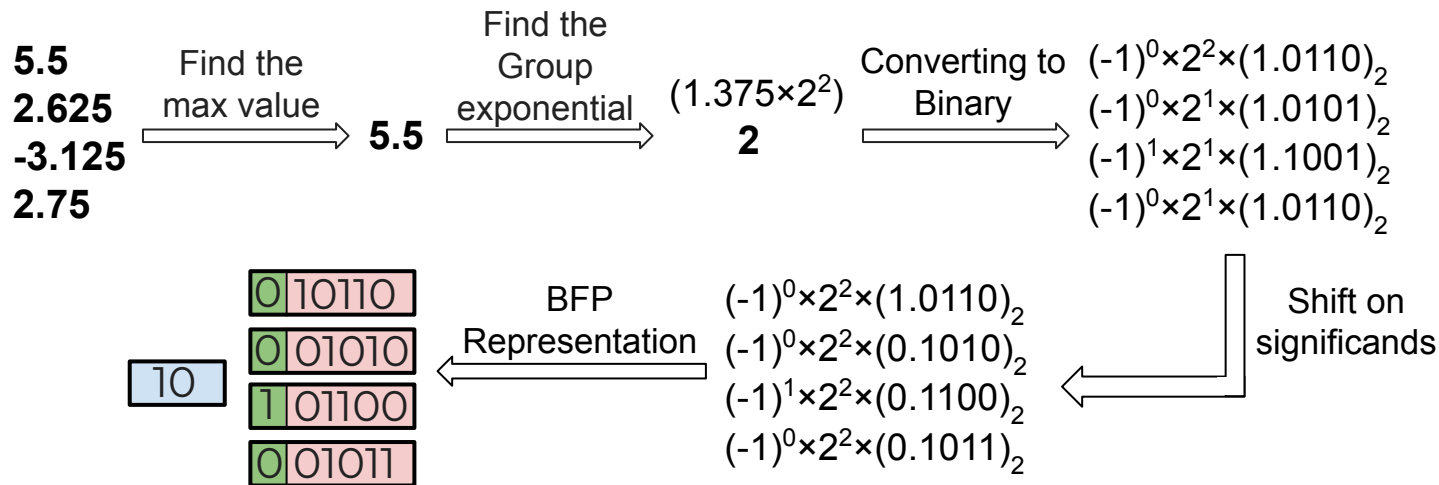
Block-Floating Arithmetics (BFP)



- Inner-group operations are performed using fixed-point arithmetic.
- Cross-group operations are performed using floating-point arithmetic.
- Each group exponent also includes a bias, which is shared across all the groups.

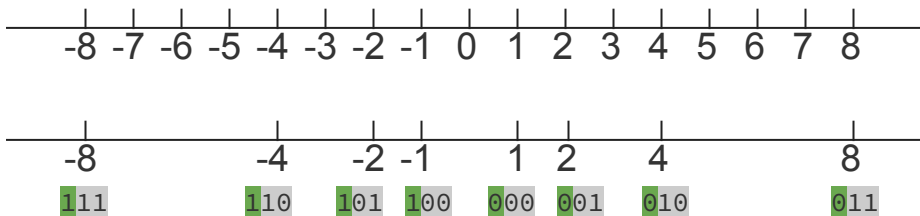
$$x = (-1)^s \times 2^{e-bias} \times m \quad \text{where } 1 \leq m < 2$$
$$m = (b_0.b_1b_2b_3\dots b_{22})_2$$

Example



Logarithm Arithmetics

- Only quantize the floating-point number to the nearest power-of-two values.
- Hardware multiplication is cheaper for power-of-two values.



- A total of 8 numbers, 3 bits are needed to encode the bits.

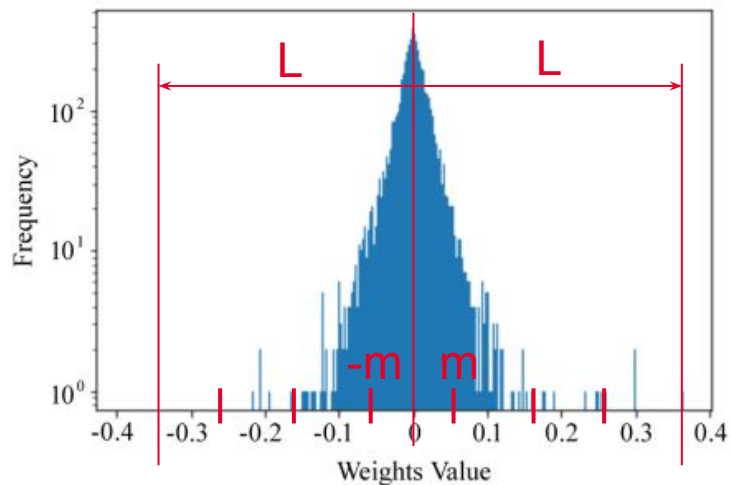
Topics

- Basic Data Formats
 - Fixed point (INT)
 - Floating point (FP)
 - Block floating point (BFP)
- Quantization methods
 - Taxonomy of Quantization
 - Learnable adaptive quantization scheme
 - Quantization for LLM

Taxonomy of Quantization

- Quantization techniques can be classified from different perspectives:
 - **Weight quantization, activation quantization**
 - Quantization aware training, post training quantization
 - Tensor-based quantization, vector-based quantization, group-based quantization
 - Quantization for inference/training
 - Deterministic quantization, stochastic quantization

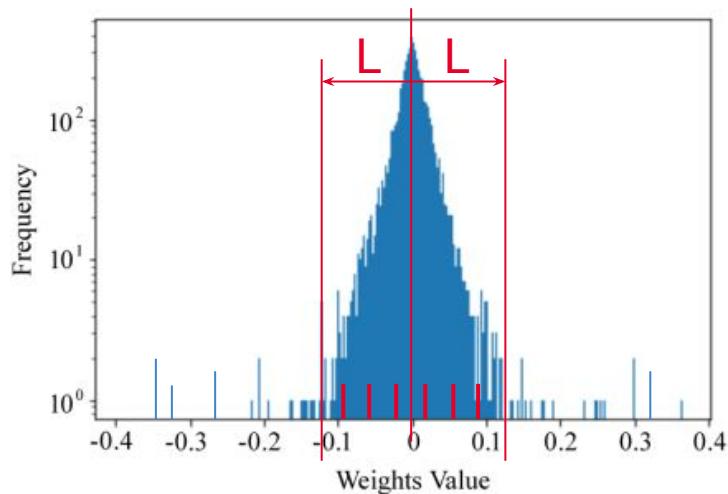
Weight Quantization



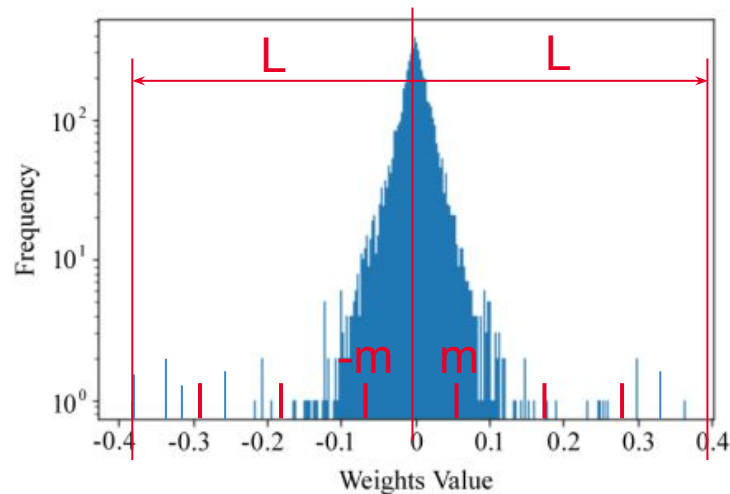
Weight distribution in ResNet

- The weight distribution follows a gaussian-like distribution.
- The outlier will lead to large quantization error.
- A good selection on the clip range L is critical for accuracy performance.

Weight Quantization



- Large truncation error
- Low quantization error for small values

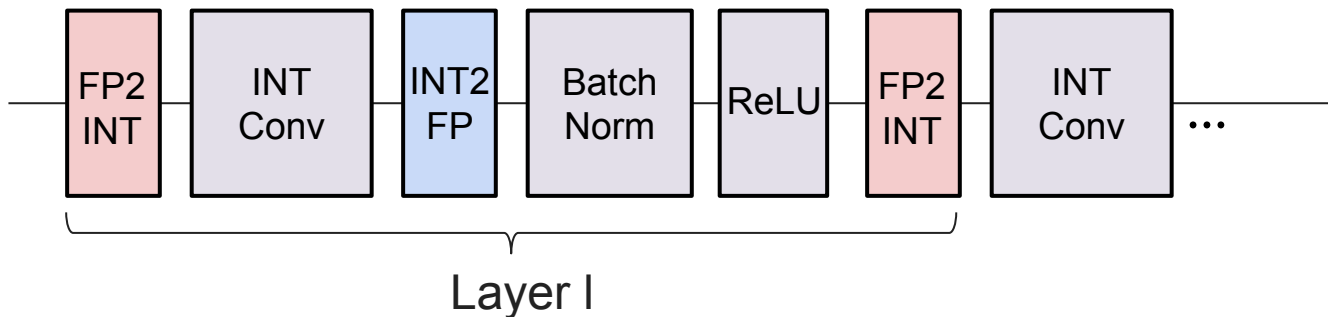


- Small truncation error
- Large quantization error for small values

- $L = 0.9 \times \max(|W|)$, $L = 0.95 \times \max(|W|)$

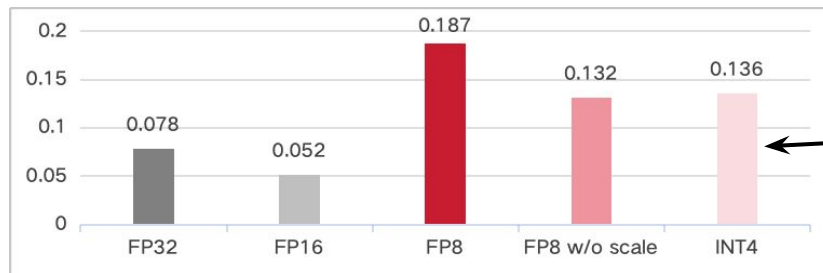
Activation Quantization

- Quantization on activation needs to be performed dynamically. This will introduce additional compute overhead.
- Also the activation will pass the nonlinear functions, requantization is required to convert the quantized number.

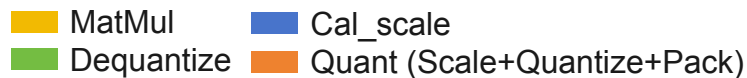


Activation Quantization

$(577 \times 1024) \times$
 (1024×1024)
Projection Layer:
Input: 577×1024
Weight: 4096×1024



On 4090 GPU



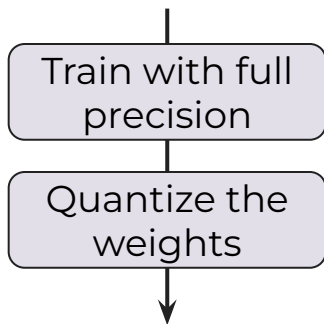
- For low-precision quantization, the quantization process may cause more computation than the computational savings achieved by using low-precision quantization.

Taxonomy of Quantization

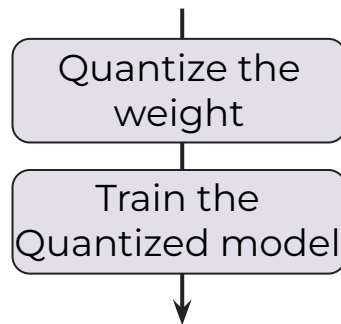
- Quantization techniques can be classified from different perspectives
 - Weight quantization, activation quantization
 - **Post training quantization, quantization aware training**
 - Tensor-based quantization, vector-based quantization, group-based quantization
 - Quantization for inference/training
 - Deterministic quantization, stochastic quantization

When to Quantize?

Post-training quantization (PTQ)

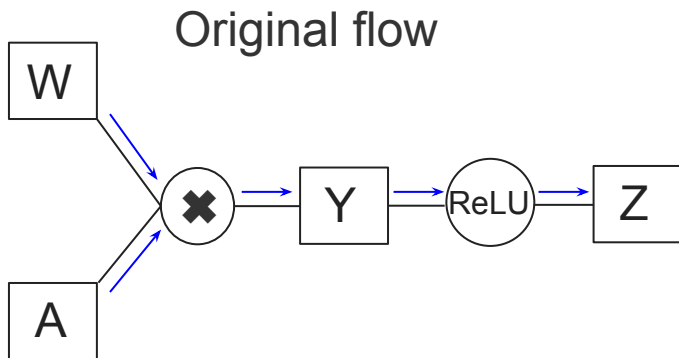


Quantization-aware Training (QAT)



- PTQ has lower computational cost, but accuracy is also lower.
- For the model which is expensive to train (LLM), PTQ is applied to facilitate their implementations.

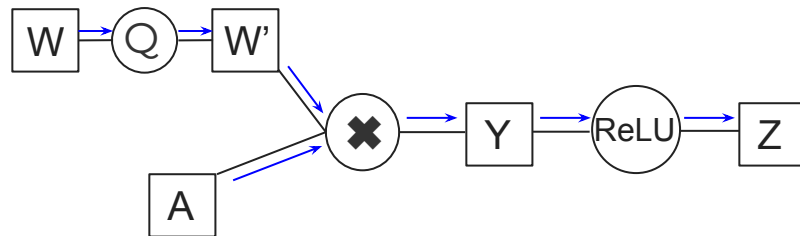
Another Way to Look at Quantization



$$Y = WA, Z = \text{ReLU}(Y)$$

$$\frac{\partial L}{\partial W} = \frac{\partial L}{\partial Z} \frac{\partial Z}{\partial Y} \frac{\partial Y}{\partial W}$$

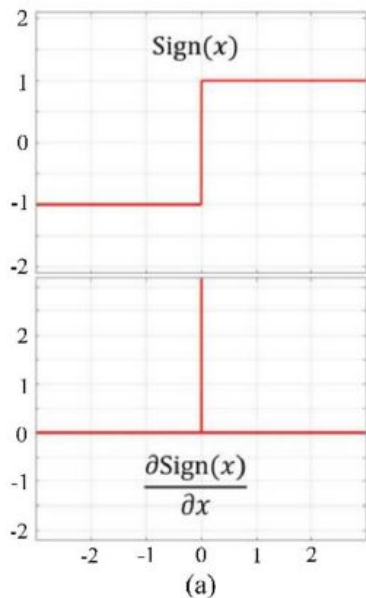
Flow with quantization



$$\frac{\partial L}{\partial W} = \frac{\partial L}{\partial Z} \frac{\partial Z}{\partial Y} \frac{\partial Y}{\partial W'} \frac{\partial W'}{\partial W}$$

How to compute $\frac{\partial W'}{\partial W}$?

Straight Through Estimator (STE)



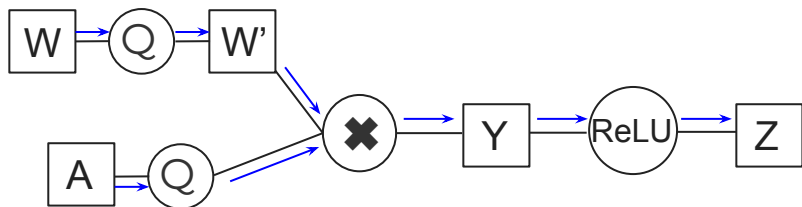
- Staircase function has a derivative of 0 at most of the values. This will make the DNN not trainable.
- We instead use STE to estimate the gradient of a non-differentiable quantized function in the backward pass.

$$\frac{\partial W'}{\partial W} = 1$$

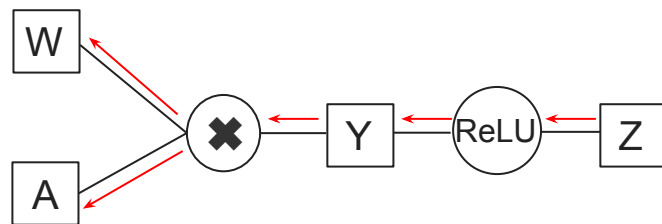
- During the forward pass, apply quantization, for backprop, ignore it.

Another Way to Look at Quantization

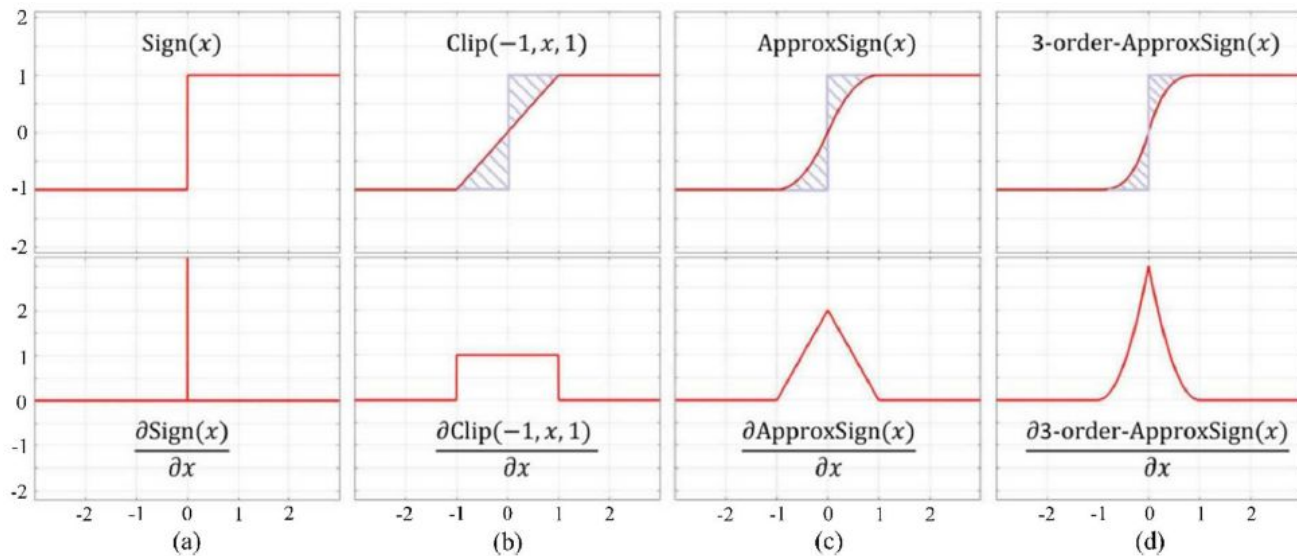
Forward pass



Backward pass



Other Ways to Approximate Quantization



Pytorch Implementation of Quantization

```
def forward(self, x):  
    y = F.conv2d(self.w, x)  
    return y
```

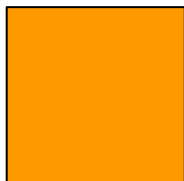
```
def forward(self, x, b, L):  
    self.quantized_w = Q(self.w, b, L)  
    y = F.conv2d(self.quantized_w, x)  
    return y  
def Q(w, b, L):  
    L = 0.9 * w.abs().max()  
    w = torch.clip(w, min=-L, max=L)  
    scale = 2L / (2**b - 2)  
    wq = (w / scale).round() * scale  
    return wq
```

Taxonomy of Quantization

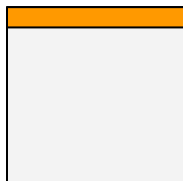
- Quantization techniques can be classified from different perspectives
 - Weight quantization, activation quantization
 - Post training quantization, quantization aware training
 - **Tensor-based quantization, vector-based quantization, group-based quantization**
 - Quantization for inference/training
 - Deterministic quantization, stochastic quantization

Granularity of Quantization

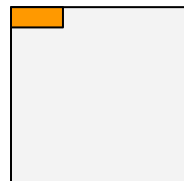
- The weight can be quantized with different granularity:
 - Tensor-based quantization
 - Vector-based quantization
 - Group-based quantization
- A higher quantization granularity will lead to a lower quantization error and a higher hardware implementation cost.



Tensor-based
quantization



Vector-based
quantization

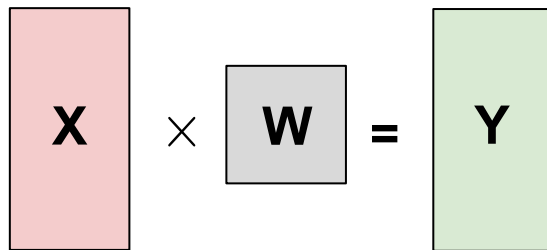


Group-based
quantization

Taxonomy of Quantization

- Quantization techniques can be classified from different perspectives
 - Weight quantization, activation quantization
 - Post training quantization, quantization aware training
 - Tensor-based quantization, vector-based quantization, group-based quantization
 - **Quantization for inference/training**
 - Deterministic quantization, stochastic quantization

Quantization During Training



X: input

W: weight filters

Y: output

- The forward propagation is very similar to the inference operation, where the input **X** is multiplied by weight **W**, generating the output **Y**.

Quantization During Training

Data gradient
Computation

$$\nabla Y \times W^T = \nabla X$$

X: input
 ∇X : input gradient

Weight gradient
Computation

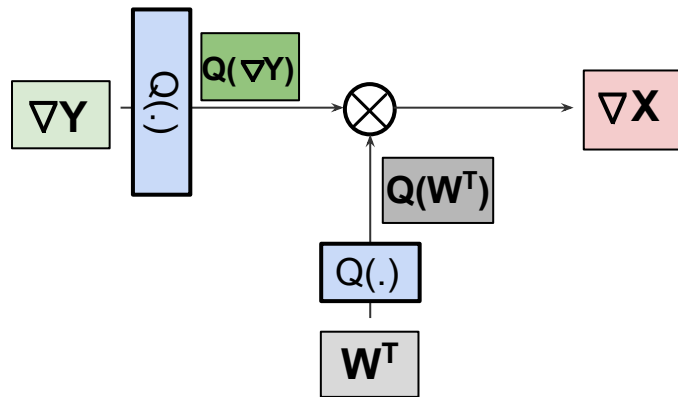
$$X^T \times \nabla Y = \nabla W$$

W: weight filters
 ∇W : weight gradient

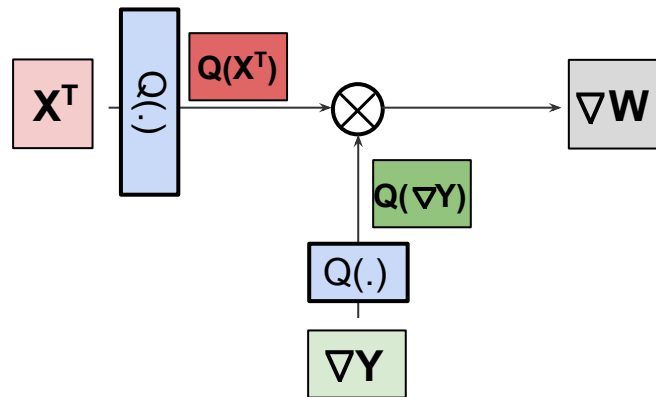
Y: output
 ∇Y : output gradient

Quantization During Training

Quantized Data Gradient Computation



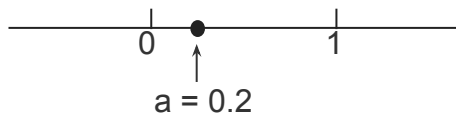
Quantized Weight Gradient Computation



Taxonomy of Quantization

- Quantization techniques can be classified from different perspectives
 - Weight quantization, activation quantization
 - Post training quantization, quantization aware training
 - Tensor-based quantization, vector-based quantization, group-based quantization
 - Quantization for inference/training
 - **Deterministic quantization, stochastic quantization**

Deterministic and Stochastic Quantization

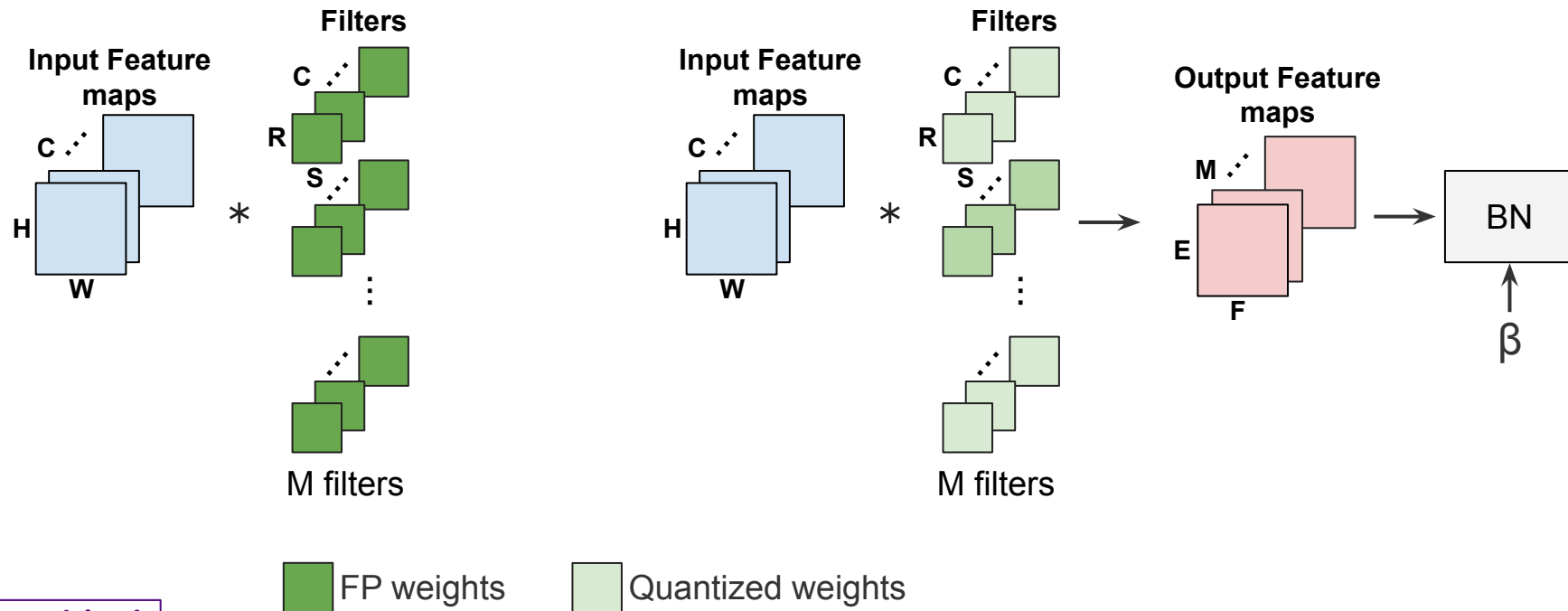


- To quantize a , conventional linear quantization will make $q(a) = 0$. However, this will cause a bias.
- With stochastic quantization:

$$q(a) = \begin{cases} 1 & \text{for } p = 0.2 \\ 0 & \text{for } p = 0.8 \end{cases}$$

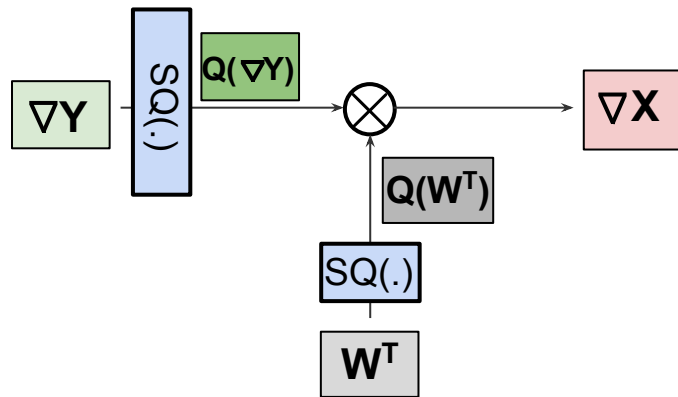
- For quantization during the forward pass of DNN training, the bias will not cause any problem, due to the existence of bias in BN.
- Stochastic quantization is extremely useful when applying quantization to accelerate DNN training.

Deterministic and Stochastic Quantization

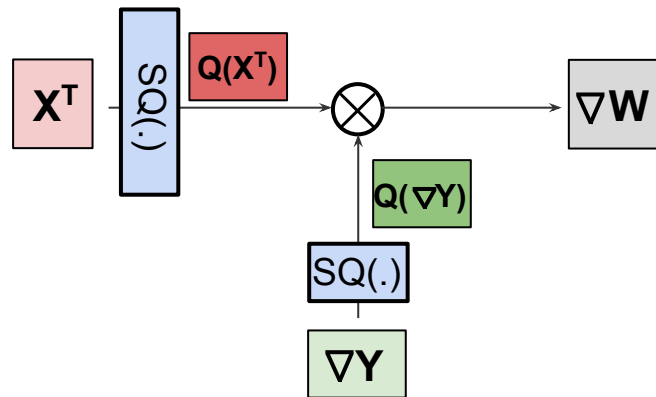


Quantization During Training

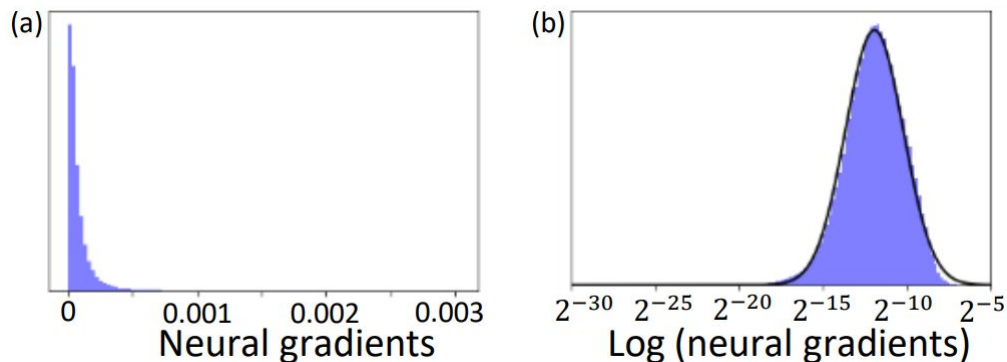
Quantized Data Gradient Computation



Quantized Weight Gradient Computation



DNN Gradient Distribution



- DNN gradient is much hard to quantize and very sensitive to quantization error.

Topics

- Basic Data Formats
 - Fixed point (INT)
 - Floating point (FP)
 - Block floating point (BFP)
- Quantization methods
 - Taxonomy of Quantization
 - Learnable adaptive quantization scheme
 - Quantization for LLM

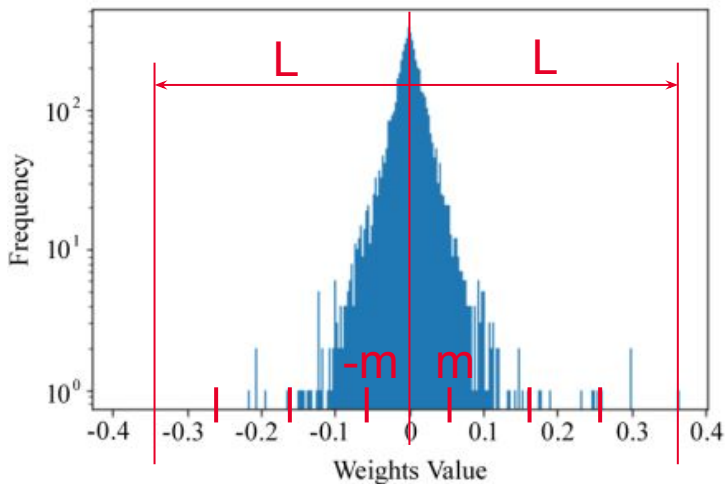
Learnable Quantization

- Multiple methods have been proposed to learn the quantization hyperparameters:
 - PACT
 - QIL
 - Quantization network
 - LQ-Net

Learnable Quantization

- How to convert a number to INT8 representation?
 - Set the clipping range: $(-L, L)$, bitwidth: b
 - Compute the scale: $s = (L_{max} - L_{min}) / (2^b - 1)$
 - Clip the input x : $x_c = \text{Clip}(x, L_{min}, L_{max})$
 - Calculate the fixed-point representation:
$$x_{int} = \text{round}((x_c - L_{min}) / s)$$
 - Rescale: $x_q = s x_{int} + L_{min}$

Learnable Quantization



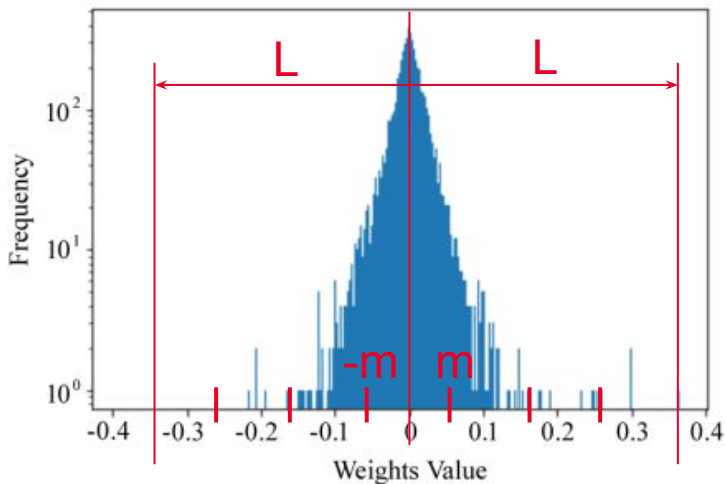
Weight distribution in ResNet

- How to convert a number to INT8 representation?
 - Set the clipping range: $(-l, l)$, bitwidth: b
 - Compute the scale: $s = (l_{\max} - l_{\min}) / (2^b - 1)$
 - Clip the input x : $x_c = \text{Clip}(x, l_{\max}, l_{\min})$
 - Calculate the fixed-point representation:
$$x_{\text{int}} = \text{round}((x_c - l_{\min}) / s)$$
 - Rescale: $x_q = s x_{\text{int}} + l_{\min}$

$$l = 0.9 \times \max(|W|), l = 0.95 \times \max(|W|)$$

Can learn by learnt during training?

Learnable Quantization



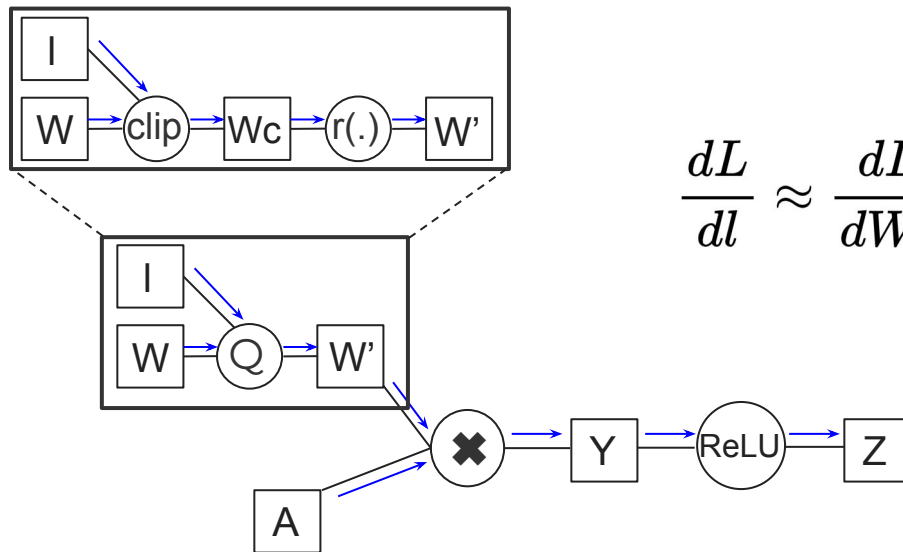
- First we need to apply CLIP function to the input x , where the clip function has a range of $(-l, l)$.

- $$x_c = \text{Clip}(x, l) = \begin{cases} l, & \text{if } x \geq l \\ x, & -l \leq x \leq l \\ -l, & x \leq -l \end{cases}$$

$$x_q = \text{round}\left(\frac{x_c}{s}\right) \times s$$

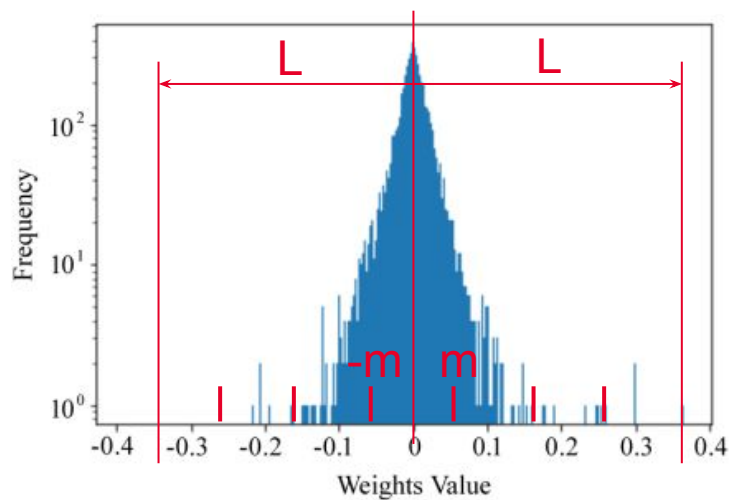
- Can we learn l ?
$$\frac{dL}{dl} = \frac{dL}{dx_q} \frac{dx_q}{dx_c} \frac{dx_c}{dl} \approx \frac{dL}{dx_q} \frac{dx_c}{dl}$$

Learnable Quantization



$$\frac{dL}{dl} \approx \frac{dL}{dW'} \frac{dW'}{dW_c} \frac{dW_c}{dl}$$

Learnable Quantization



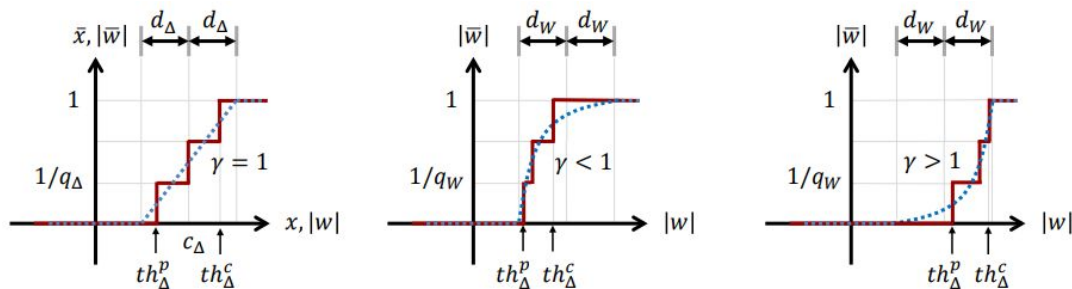
$$\text{Clip}(x, l) = \begin{cases} l, & \text{if } x \geq l \\ x, & -l \leq x \leq l \\ -l, & x \leq -l \end{cases}$$

$$\frac{d\text{Clip}(x, l)}{dx} = \begin{cases} 0, & \text{if } x \geq l \\ 1, & -l \leq x \leq l \\ 0, & x \leq -l \end{cases}$$

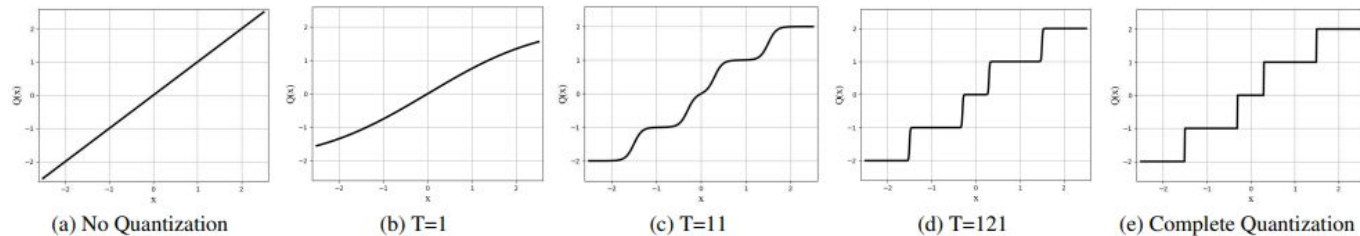
$$\frac{d\text{Clip}(x, l)}{dl} = \begin{cases} 1, & \text{if } x \geq l \\ 0, & -l \leq x \leq l \\ -1, & x \leq -l \end{cases}$$

L can be learnable

Learnable Quantization



QIL

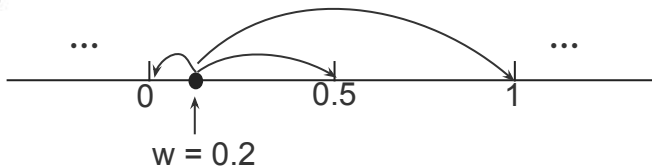
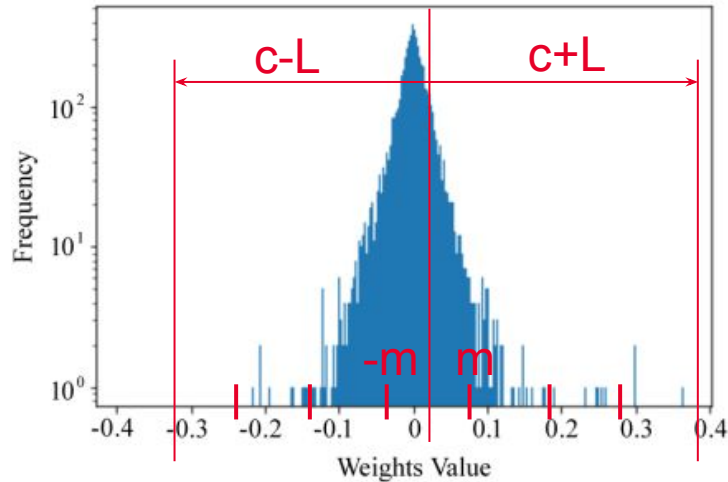
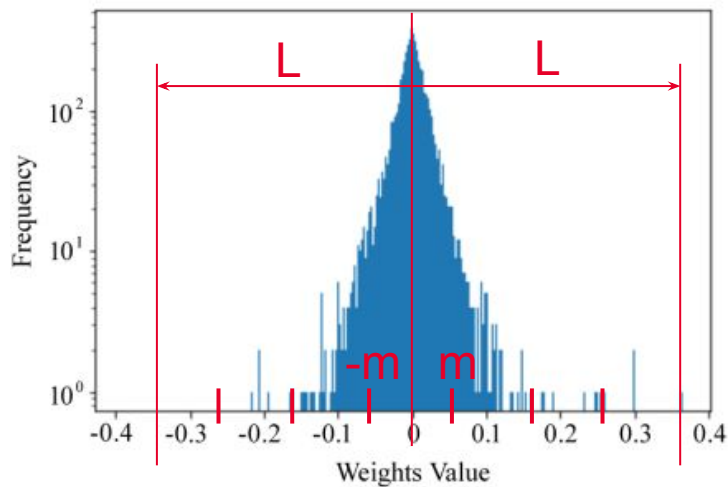


QN

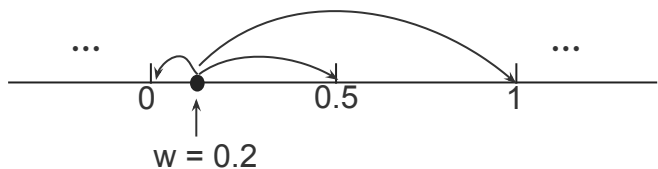
Jung, Sangil, et al. "Learning to quantize deep networks by optimizing quantization intervals with task loss." *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*. 2019.

Yang, Jiwei, et al. "Quantization networks." *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*. 2019.

Quantization Interval Learning (QIL)



Quantization Interval Learning (QIL)



- To achieve this rounding flexibility, we combine a learnable function with quantization.

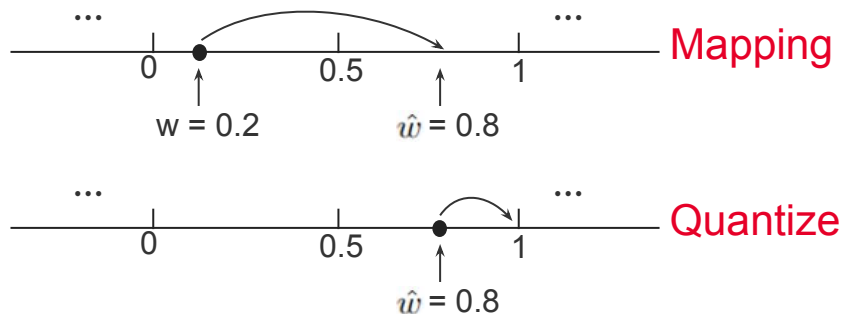
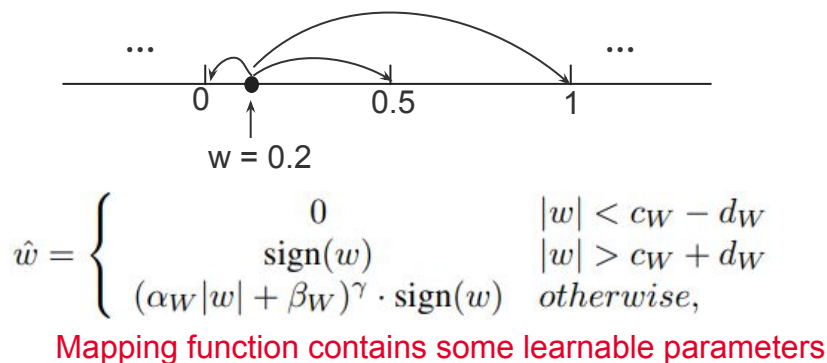
$$w_q = Q(w) \longrightarrow w_q = Q(F(w))$$

- $F(\cdot)$ is a function which contains learnable hyperparameters.

$$\hat{w} = \begin{cases} 0 & |w| < c_W - d_W \\ \text{sign}(w) & |w| > c_W + d_W \\ (\alpha_W |w| + \beta_W)^\gamma \cdot \text{sign}(w) & \text{otherwise,} \end{cases}$$

Quantization Interval Learning (QIL)

- QIL offers flexibility to round the FP weights.



- $w_q = Q(F(w))$ are stored for inference after the training process finished.
- We can not apply this techniques over the activation, due to its large computational overhead.

LQ-Nets

$$q(x) = \left\langle \begin{bmatrix} \mathbf{v} \\ 1 \\ 2 \\ \dots \\ 2^{K-1} \end{bmatrix}, \begin{bmatrix} \mathbf{e}_x \\ b_1 \\ b_2 \\ \dots \\ b_K \end{bmatrix} \right\rangle$$

$Q(x) = \mathbf{v}^T \mathbf{e}_x$, \mathbf{e}_x is a binary vector

- \mathbf{V} can be learnable.
- The resultant quantization can still facilitate MAC computation.

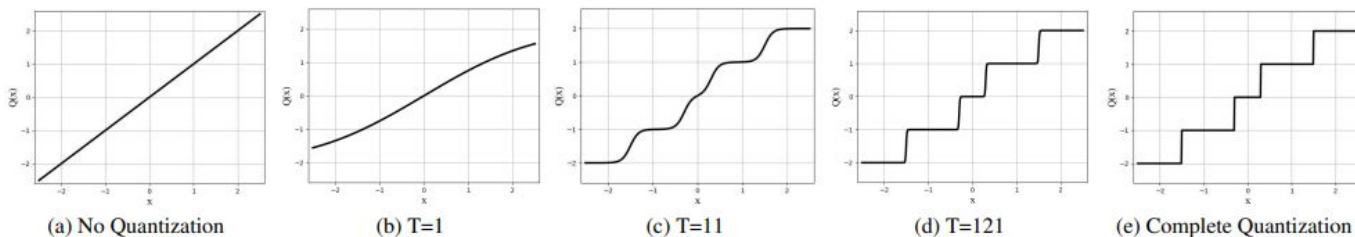
$$Q_{\text{ours}}(\mathbf{w}, \mathbf{v}^w)^T Q_{\text{ours}}(\mathbf{a}, \mathbf{v}^a) = \sum_{i=1}^{K_w} \sum_{j=1}^{K_a} v_i^w v_j^a (\mathbf{b}_i^w \odot \mathbf{b}_j^a)$$

- $v_i^w v_j^a$ can be computed at low cost.
- $\mathbf{b}_i^w \odot \mathbf{b}_j^a$ can be pre-computed.

$$\mathbf{v}^*, B^* = \arg \min_{\mathbf{v}, B} \|\mathbf{B}^T \mathbf{v} - \mathbf{x}\|_2^2, \quad s.t. B \in \{-1, 1\}^{K \times N}$$

Quantization Networks

- We propose a novel perspective of interpreting and implementing neural network quantization by formulating low-bit quantization as a differentiable non-linear function.



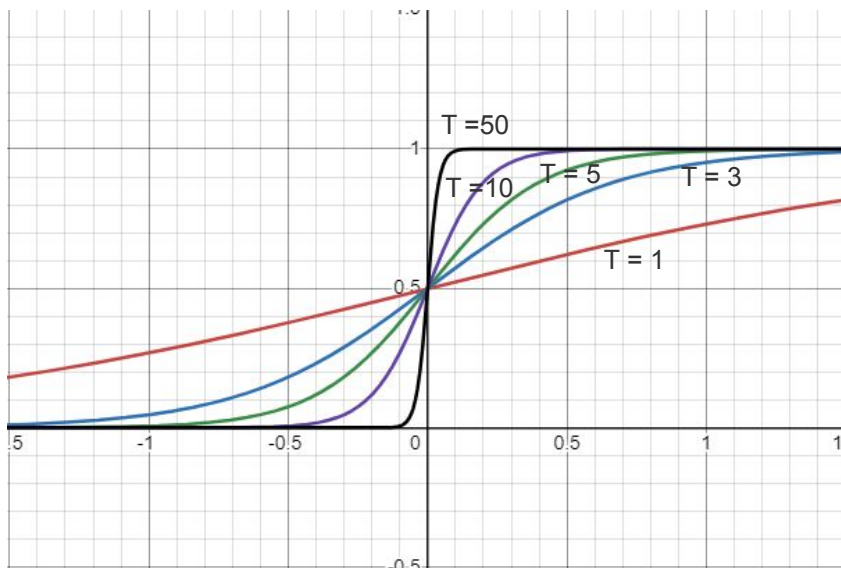
$$y = \alpha \left(\sum_{i=1}^n s_i \mathcal{A}(\beta x - b_i) - o \right)$$

$$\mathcal{A}(x) = \begin{cases} 1 & x \geq 0, \\ 0 & x < 0. \end{cases}$$

- $n + 1$ is the number of quantization intervals
- β is the scale factor of inputs
- s_i and b_i are the scales and biases for the unit step functions

Quantization Networks

$$\mathcal{A}(x) = \begin{cases} 1 & x \geq 0, \\ 0 & x < 0. \end{cases} \quad \sigma(Tx) = \frac{1}{1 + \exp(-Tx)}$$



- We can replace the staircase function with a sigmoid function.
- We can progressively increase T during the training process.

Presentation

- Trained ternary quantization
- Incremental network quantization: Towards lossless cnns with low-precision weights
- Quantization and training of neural networks for efficient integer-arithmetic-only inference
- Smoothquant: Accurate and efficient post-training quantization for large language models

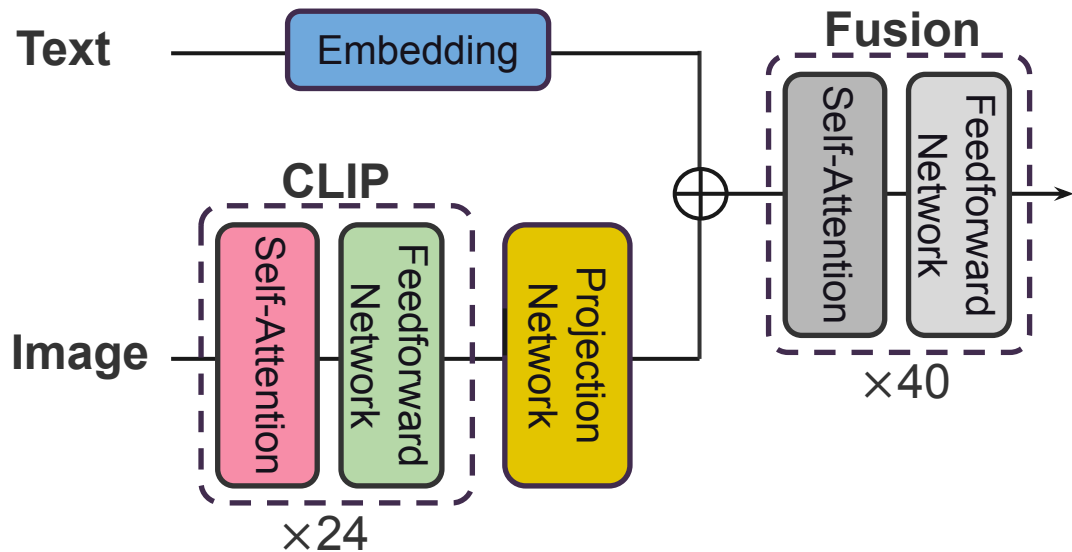
Topics

- Basic Data Formats
 - Fixed point (INT)
 - Floating point (FP)
 - Block floating point (BFP)
- Quantization methods
 - Taxonomy of Quantization
 - Learnable adaptive quantization scheme
 - **Quantization for LLM**

Post Training Quantization

- Several Methods have been proposed to efficient post-training quantization.
- Given the large size of the modern LLM, it is beneficial to applied the quantization on the model directly without the need of finetuning.

Model Architecture: Llava



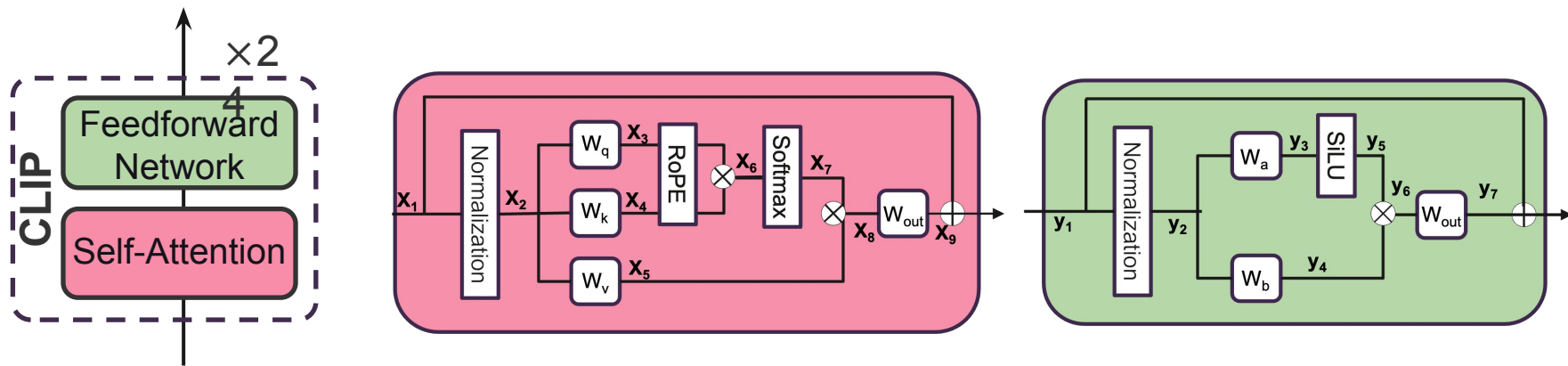
\oplus Concat in token dim

Visual input example, Tricky Question:



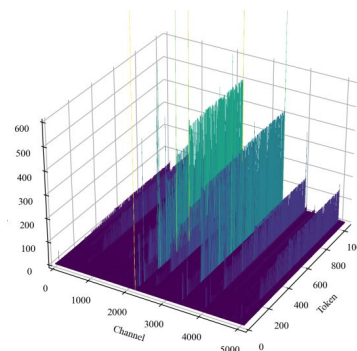
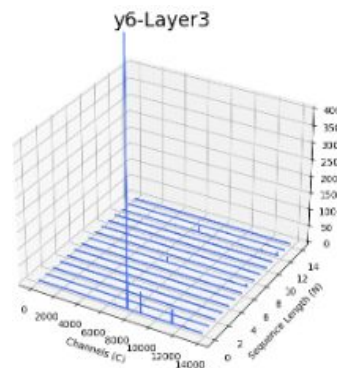
User	If there are factual errors in the questions, point it out; if not, proceed answering the question. What's happening in the desert?
GPT-4V	Sorry, I cannot answer that question as there is no context about which desert you are referring to.
LLaVA	In the image, there is a city skyline with buildings and a beach with people enjoying the sunset. The scene is set in a desert, which is unusual for a beach setting.
LLaVA-1.5	There are no deserts in the image. The image features a beach with palm trees, a city skyline, and a large body of water.

CLIP Architecture



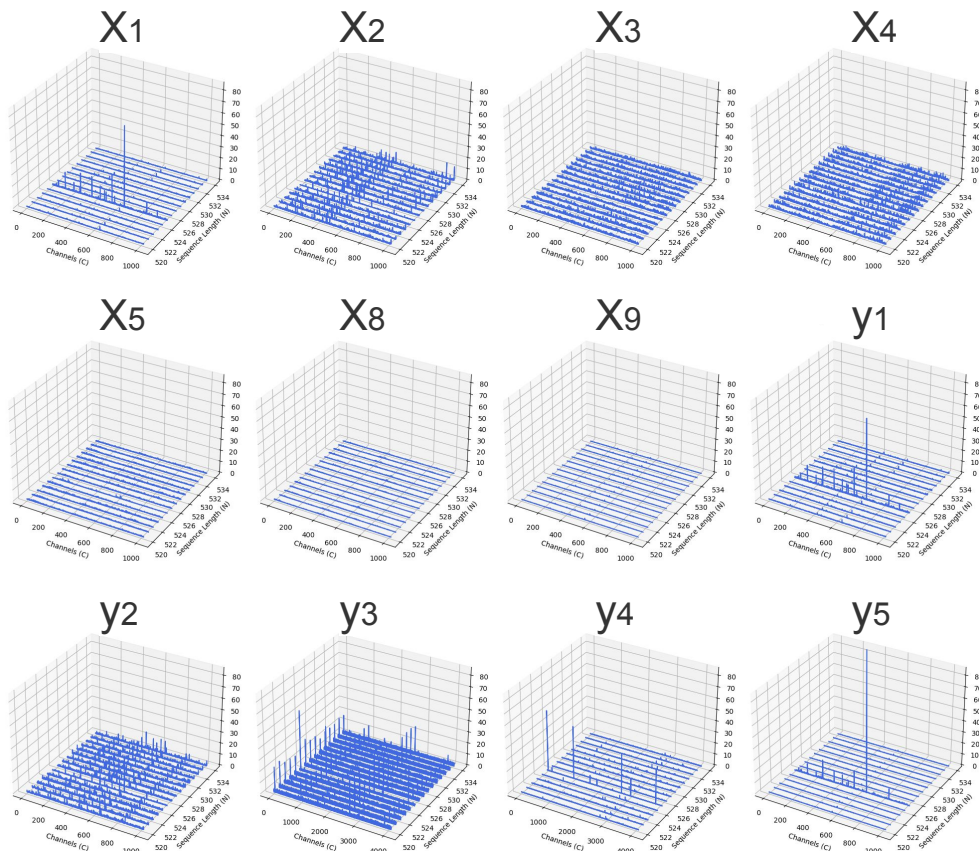
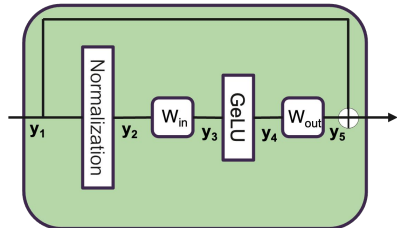
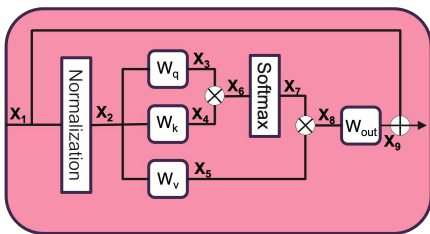
Types of Outlier

- Massive Activation:
 - For an activation matrix A , an massive activation is an element A_{ij} within it that satisfies:
 - $A_{ij} > \eta \times \text{mean}(|A|)$
 - $A_{ij} > \gamma$
 - $\eta=300, \gamma=50$
- Channelwise Outlier:
 - $\text{mean}(A_i) > \eta \times \text{std}(A) + \text{mean}(|A|)$
 - $\text{std}(A_i) < \beta$
 - $\eta=3, \beta=0.6$



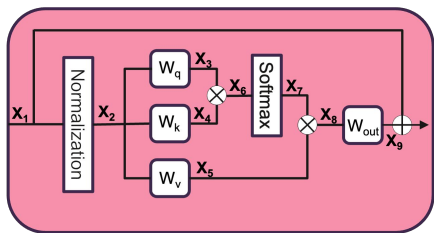
Outlier Study: CLIP

- 3D activation within **layer 12**
MA is produced on y_5
MA is propagated on x_1, y_1 from layer11



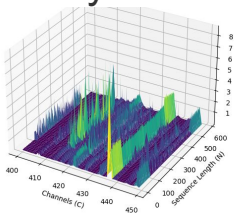
Outlier Study: CLIP

- 3D plots of X2 across layers.

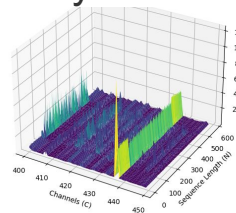


- x2 exhibits channel wise outlier

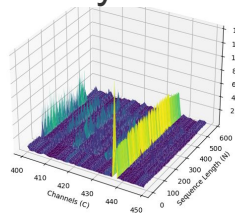
Layer 1



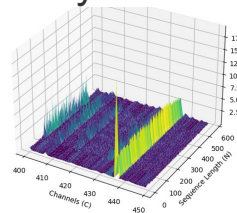
Layer 2



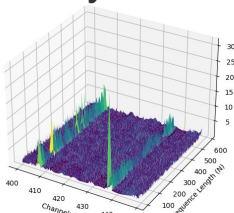
Layer 3



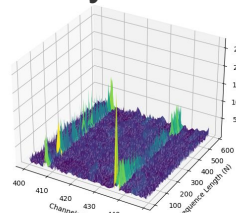
Layer 4



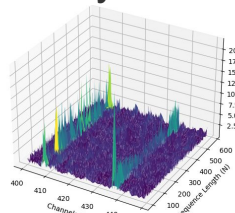
Layer 11



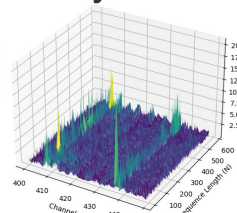
Layer 12



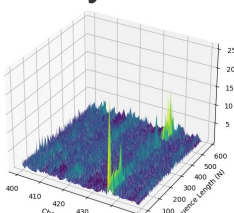
Layer 13



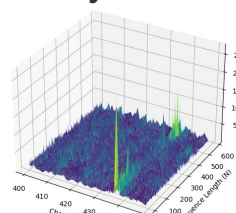
Layer 14



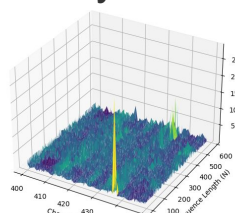
Layer 19



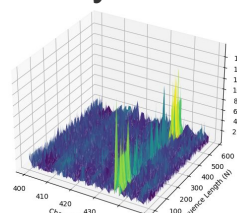
Layer 20



Layer 21

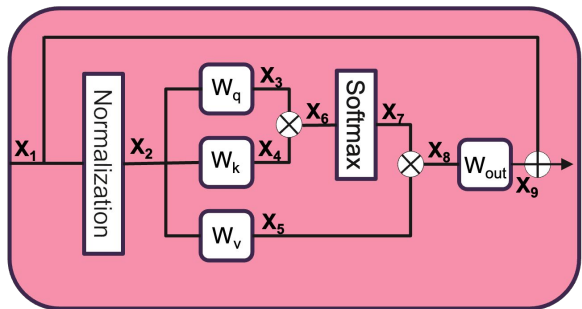


Layer 23

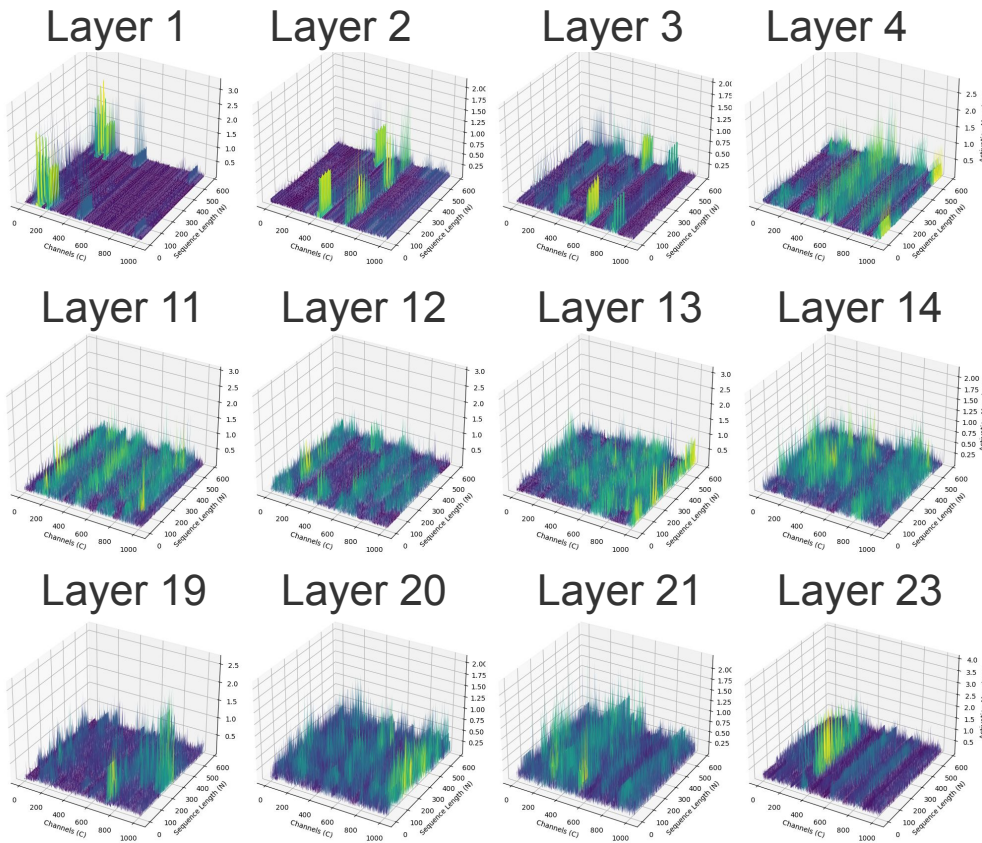


Outlier Study: CLIP

- 3D plots of x_8 across layers.

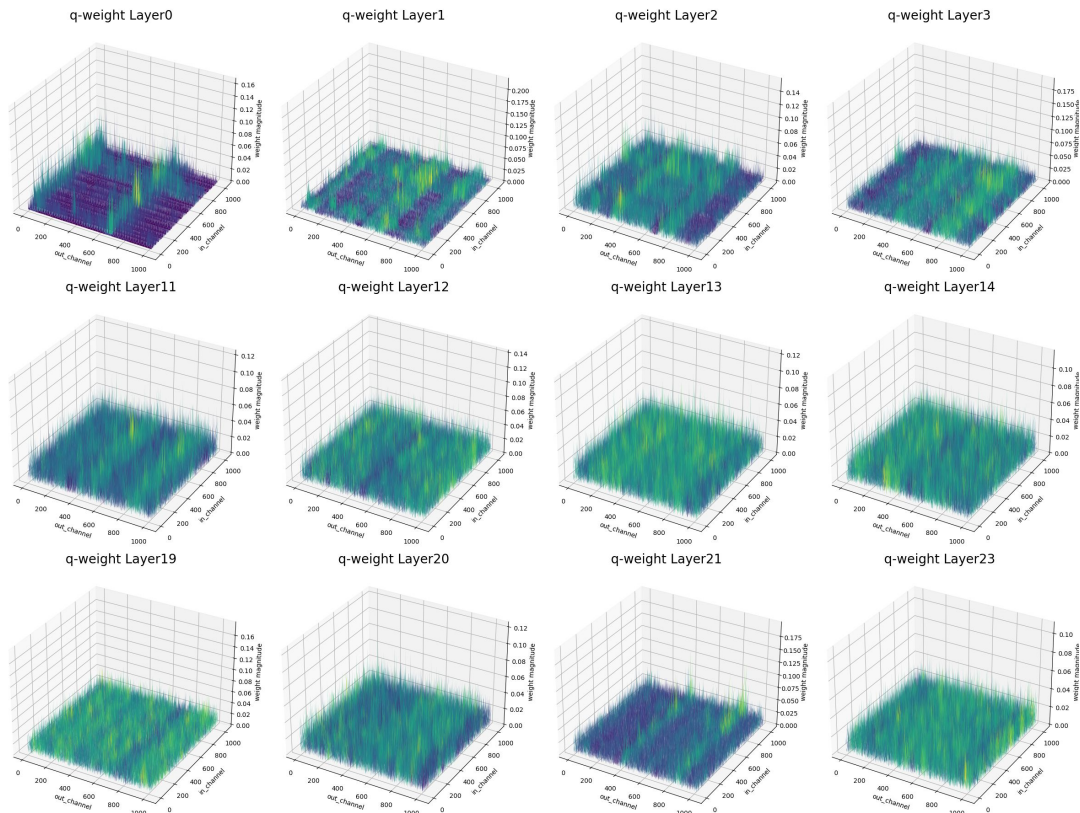
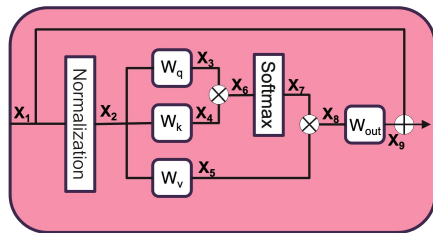


- x_8 exhibits channel wise outlier



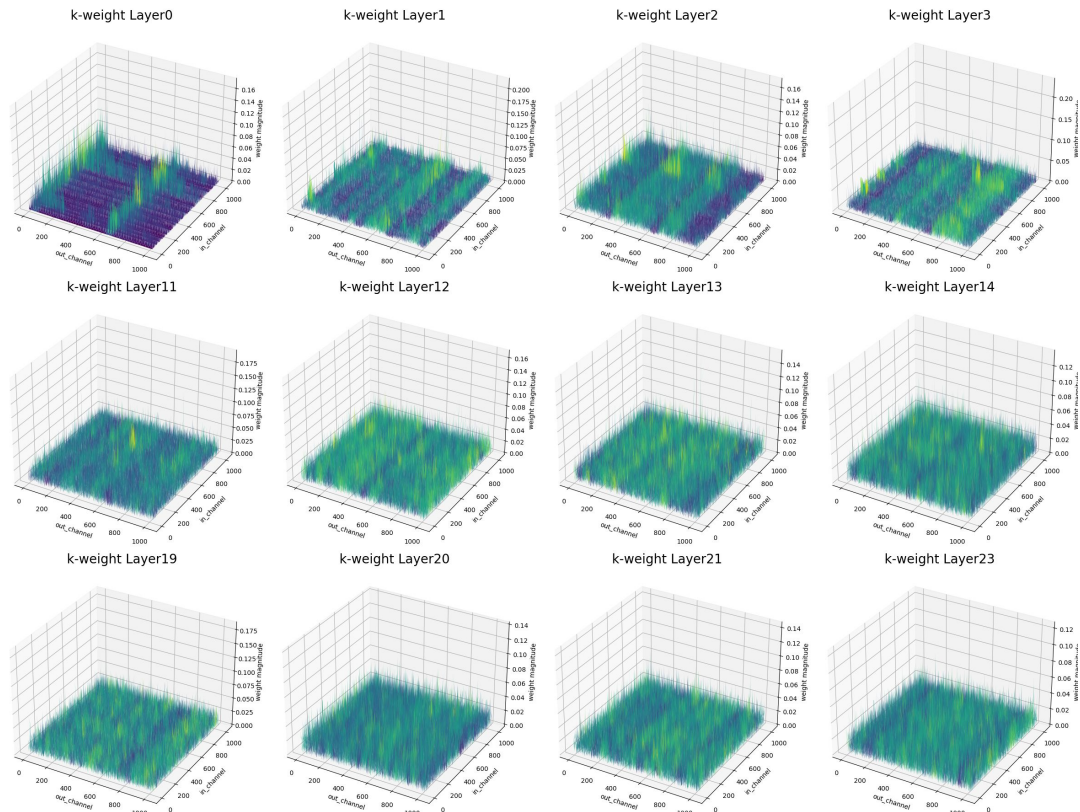
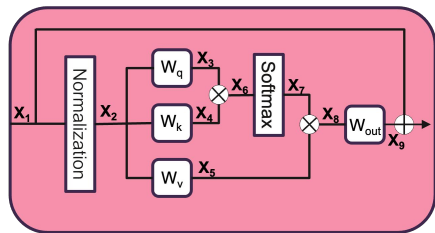
Outlier Study: CLIP Weights

- W_q across CLIP layers.



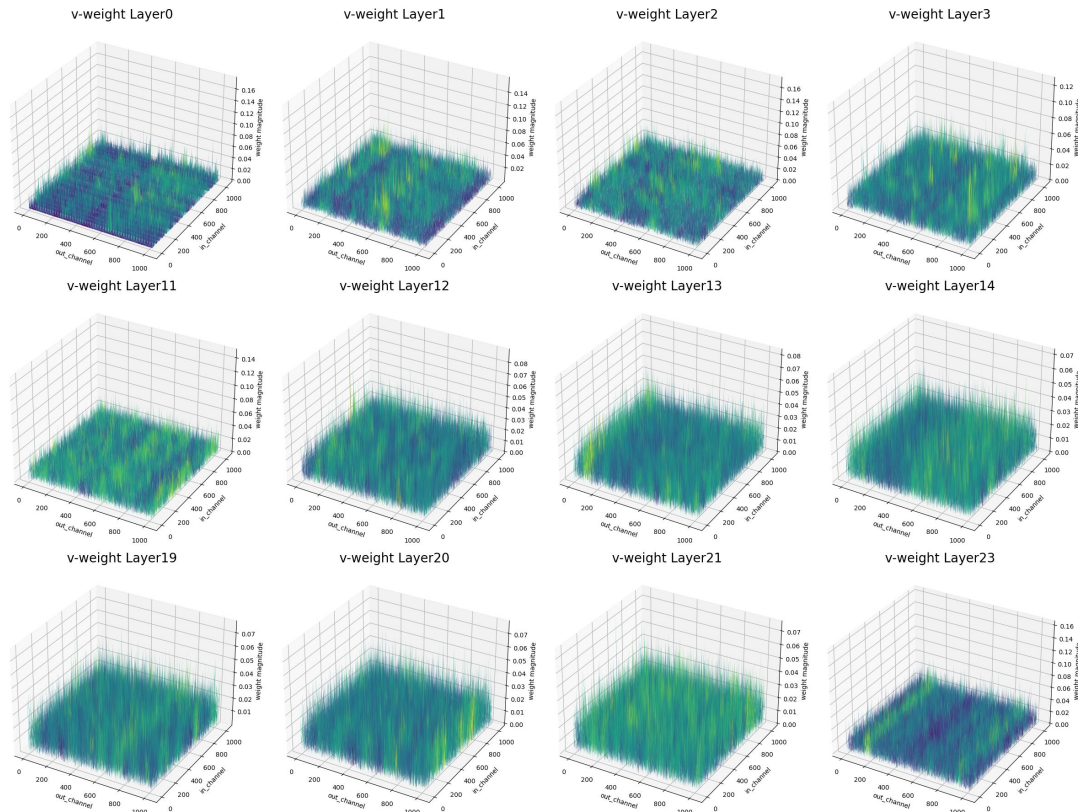
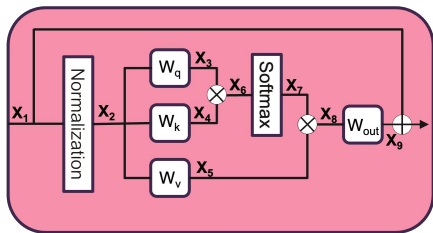
Outlier Study: CLIP Weights

- W_k across CLIP layers.



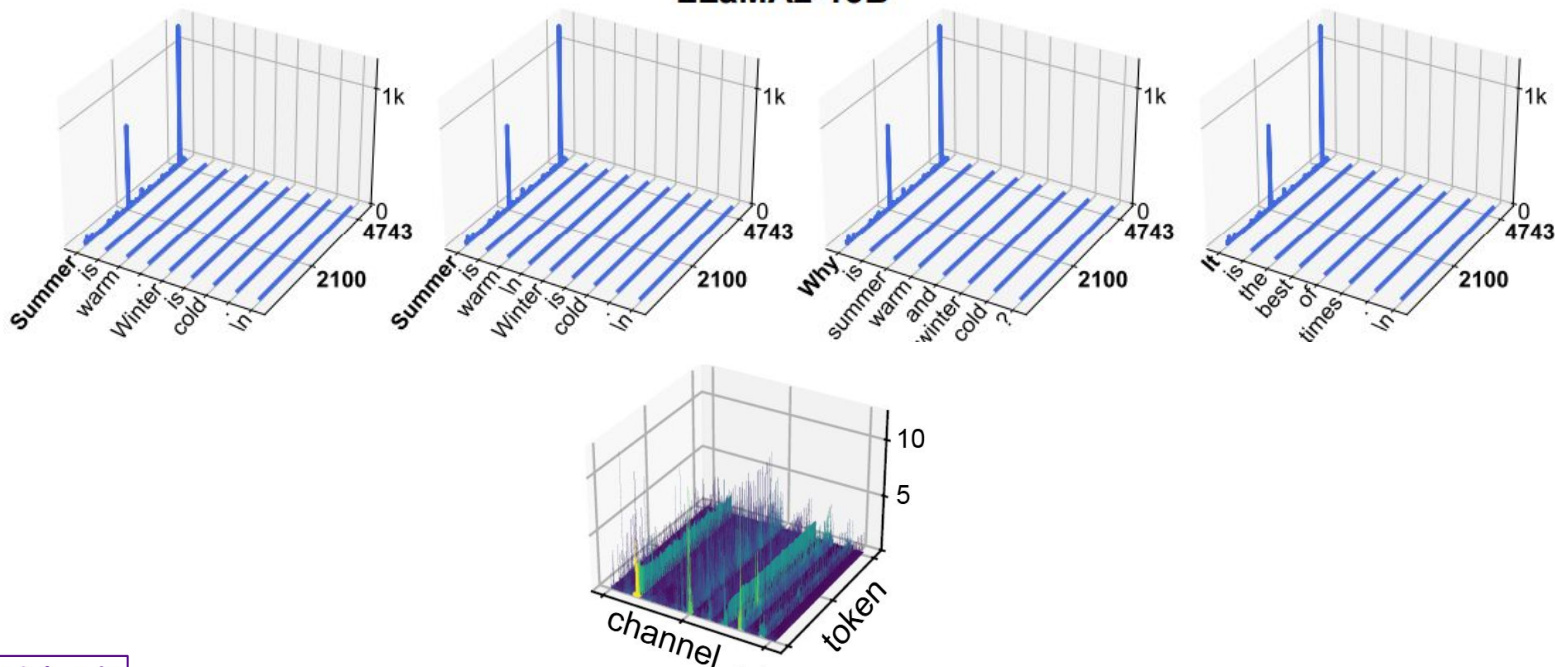
Outlier Study: CLIP Weights

- W_v across CLIP layers.



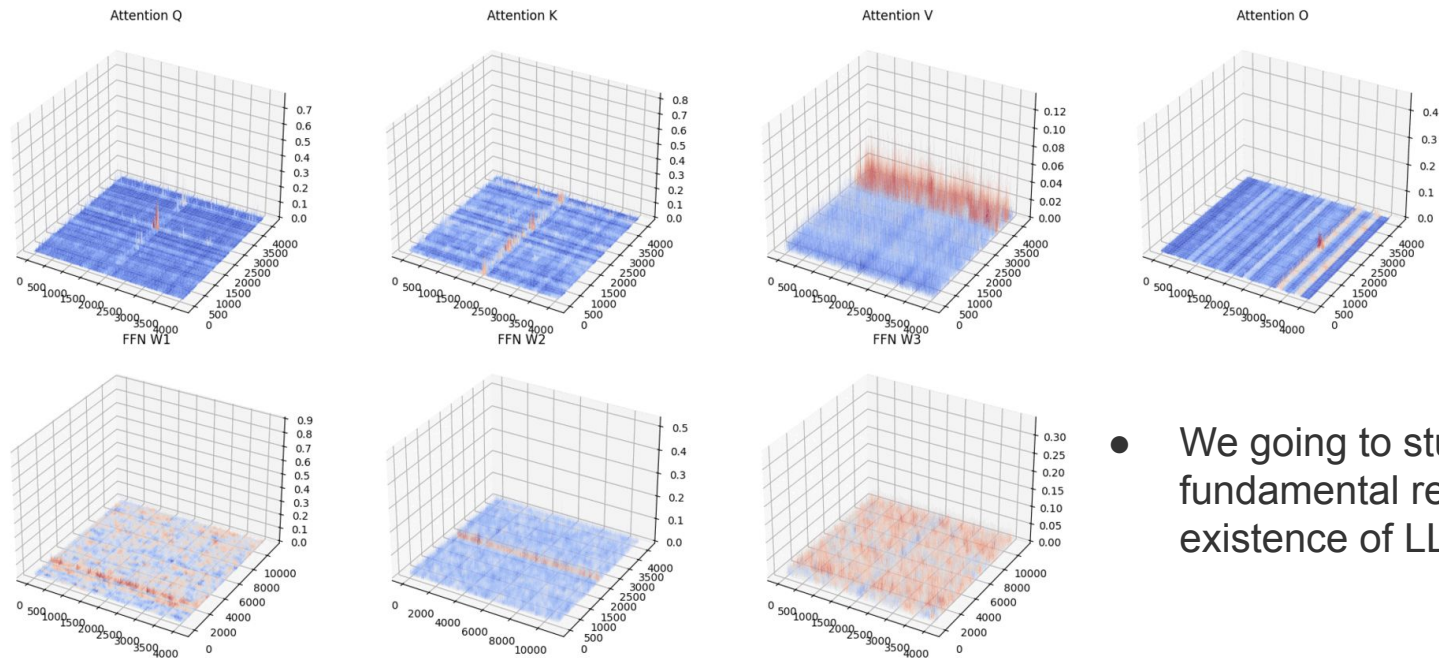
Outlier Study: LLaMA Activations

LLaMA2-13B



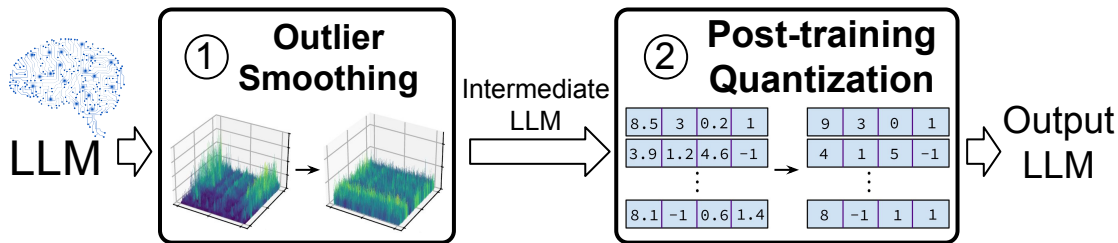
Study the Reason of LLM Outliers

Layer 0 Weights



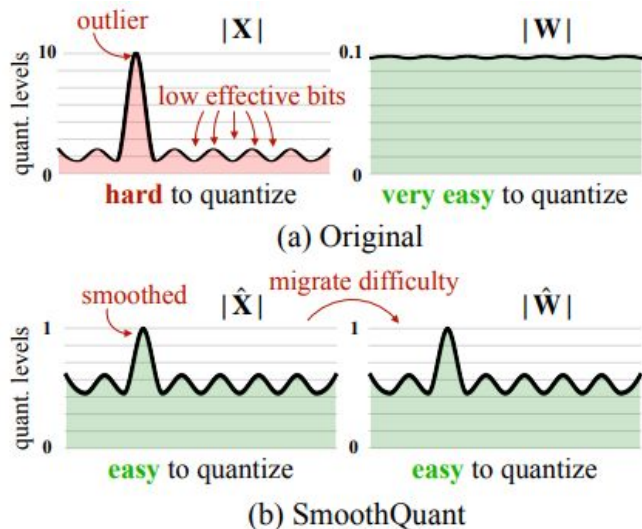
- We going to study the fundamental reason of the existence of LLM outliers.

Outlier Smoothing



- When performing post-training quantization on a LLM, it's common to include a step of outlier smoothing prior to the quantization process.

SmoothQuant

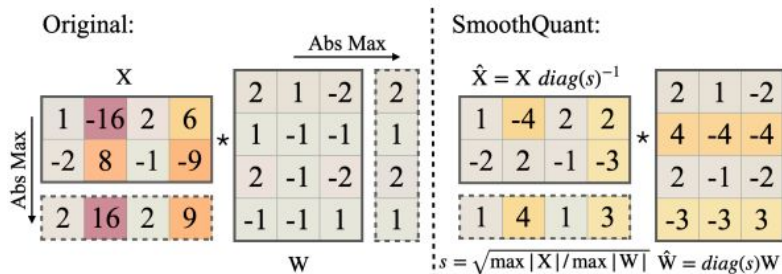


- The intermediate results within LLM usually have a lot of outliers.
- SmoothQuant smooths the activation outliers by offline migrating the quantization difficulty from activations to weights with a mathematically equivalent transformation.

$$\mathbf{Y} = (\mathbf{X}\text{diag}(\mathbf{s})^{-1}) \cdot (\text{diag}(\mathbf{s})\mathbf{W}) = \hat{\mathbf{X}}\hat{\mathbf{W}}$$

- \mathbf{s} depends on the square root of the magnitude of the largest channel

SmoothQuant

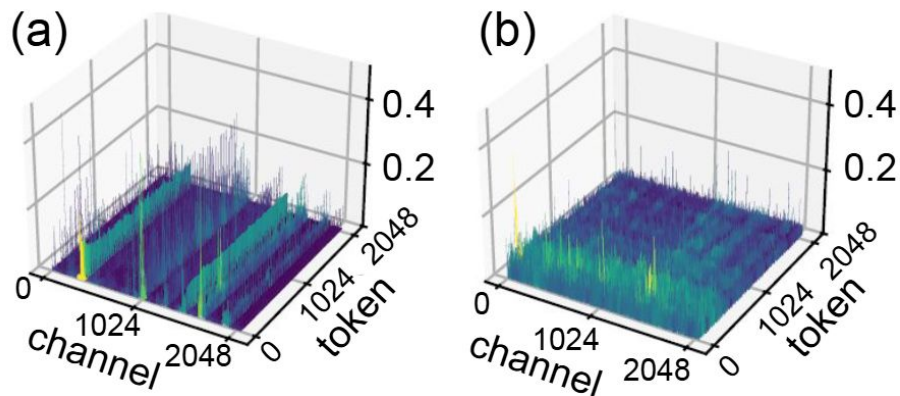


$$Y = (X \text{diag}(s)^{-1}) \cdot (\text{diag}(s)W) = \hat{X} \hat{W}$$

$$s_j = \max(|X_j|)^\alpha / \max(|W_j|)^{1-\alpha}$$

<i>OPT-175B</i>	LAMBADA	HellaSwag	PIQA	WinoGrande
FP16	74.7%	59.3%	79.7%	72.6%
W8A8	0.0%	25.6%	53.4%	50.3%
ZeroQuant	0.0%*	26.0%	51.7%	49.3%
LLM.int8()	74.7%	59.2%	79.7%	72.1%
Outlier Suppression	0.00%	25.8%	52.5%	48.6%
SmoothQuant-O1	74.7%	59.2%	79.7%	71.2%
SmoothQuant-O2	75.0%	59.0%	79.2%	71.2%
SmoothQuant-O3	74.6%	58.9%	79.7%	71.2%

QuaRot



- QuaRot introduces a novel methods to convert the weights and activation of LLM.
- After conversion, most of the outliers within the activation and weights are removed.
- This conversion introduces almost no additional cost during the inference.

QuaRot

- Assume $Y = AW$, where A may have outliers, quantizing A and W as $Q(A)$ and $Q(W)$ could result in increased quantization error. Consequently, $Q(A)Q(W)$ may differ significantly from AW .
- With QuaRot, a orthogonal matrix is applied to eliminate the outliers within A .

$$A \longrightarrow \boxed{W} \longrightarrow AW$$

$$AR \longrightarrow \boxed{R^T W} \longrightarrow AW$$

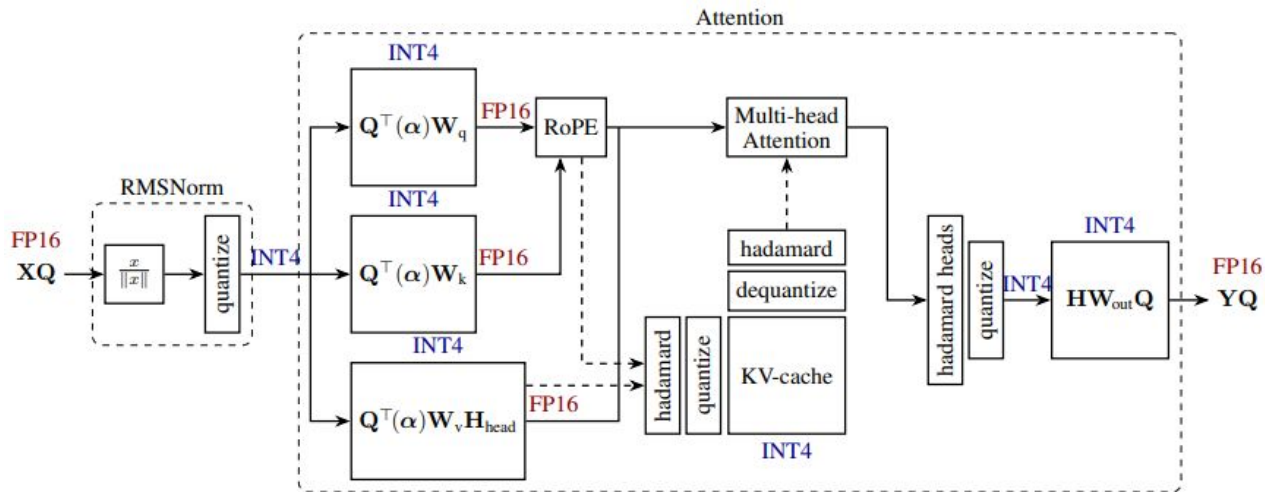
$$R^T R = R R^T = I$$

$$Q(A) \longrightarrow \boxed{Q(W)} \longrightarrow Q(A)Q(W)$$

$$Q(AR) \longrightarrow \boxed{Q(R^T W)} \longrightarrow Q(AR)Q(R^T W)$$

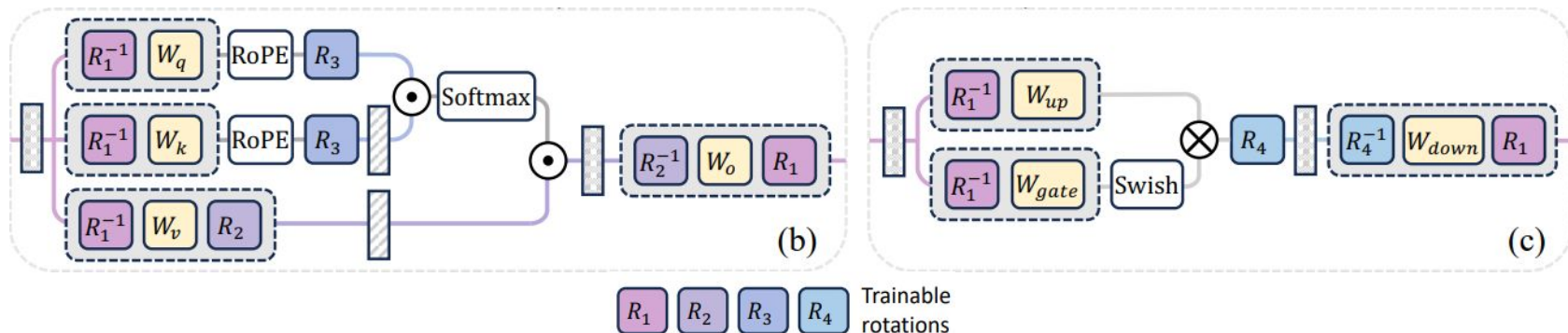
- $R^T W$ can be computed offline, AR can be generated by modifying the weight matrices of the last layer.

QuaRot



- For some of the layers, the conversion needs to be performed online
- We can use Hadamard matrix, which consists of only 1 and -1 to facilitate the matrix multiplications.

SpinQuant



$$\arg \min_{R \in \mathcal{M}} \mathcal{L}_Q(R_1, R_2 \mid W, X)$$

- SpinQuant optimizes (or learns) the rotation matrices to obtain the minimal changes on the training loss.
- We have to ensure the rotational matrix still satisfies the orthogonal property \rightarrow Cayley Optimization.