

Lecture 05: Quantization Strategies for Efficient DNN Implementation

Notes

- Please send email to efficientaiaccelerator@gmail.com
- Lab1 has been released.
- Start considering the project topic, teaming.
- In-course quiz today, covering materials of DNN pruning.



Recap

- Why pruning?
 - Reduce running cost
 - Reduce storage
- 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

4-bit Fixed Point (INT4)



8-bit Fixed Point (INT8)

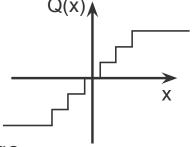


- Hyperparameter associated with the fixed-point format:
 - Clipping range (-L, L): usually symmetrical around 0
 - o Bitwidth (b)
- Quantization with Fixed-point format is called Fixed point quantization or INT quantization.



Fixed-Point Format (Symmetrical)

- How to convert a number x to INT representation?
 - Set the clipping range: (-L, L), bitwidth: b
 - Compute the scale: $s = 2L/(2^b 2)$
 - O Clip the input x: $x_c = Clip(x, L, -L)$
 - \circ Calculate the INT representation: $x_{int} = round(x_c/s)$
 - \circ Rescale: $x_q = sx_{int}$



- Have a uniform representation power within the clipping range.
- ullet All the computations can be performed using x_{int}

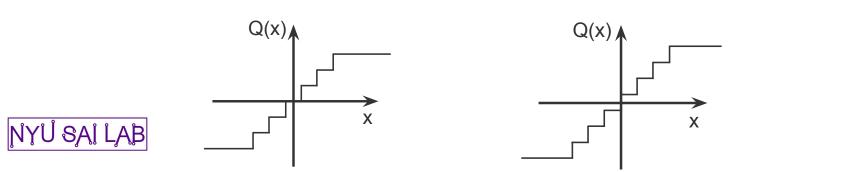


Fixed-Point Format (Symmetrical)

- Have a uniform representation power within the clipping range.
- ullet All the computations can be performed using x_{int}



• With s=2L/(2^b-2), zero can be represented using quantized number



Example

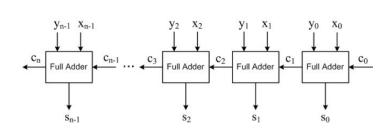
- X = [1.1, 2.4, -0.3, 0.8], bitwidth = 3, L = 2
- How to convert a number x to INT representation?
 - Set the clipping range: (-L, L), bitwidth: b b=3, L=2
 - \circ Compute the scale: $s=2L/(2^b-2)$ $\,$ s = 4/6 = 2/3 $\,$
 - o Clip the input x: $x_c = Clip(x, L, -L)$ xc = [1.1, 2, -0.3, 0.8]
 - \circ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



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

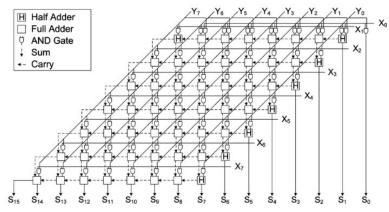
If the scales are the same

ullet Division: $x_q/y_q=x_{int}/y_{int}$



Fixed-point adder

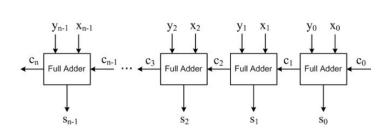




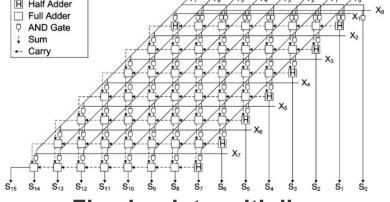
Fixed-point multiplier

- Addition/Subtraction: Hard to compute
- ullet Multiplication: $x_q imes y_q = s_x s_y (x_{int} imes y_{int})$
- ullet Division: $x_q/y_q = (s_x/s_y) imes (x_{int}/y_{int})$

If the scales are not the same



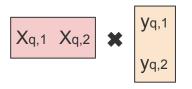
Fixed-point adder



Fixed-point multiplier

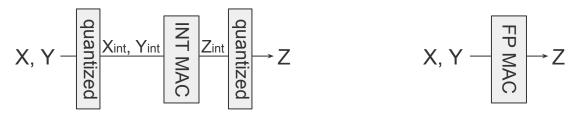


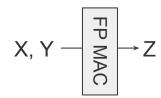
If we try to compute the dot product between X and Y:



 $X_{q,1}$ $X_{q,2}$ $X_{q,2}$ All elements within the tensors are quantized using the same scale, but the scale across the tensors can be different

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







 INT can be applied to a block of numbers, with the block size defined in a customizable manner.

1	2
10	11

1.1	2.4
10.5	11.8

Per tensor quantization

1.1	2.4
10.5	11.8

Row-wise quantization (low error)

1.1	2.4
10.5	11.8

Column-wise quantization (high error)

However, a higher quantization granularity will also incur a conversion overhead.





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: (Lmin, Lmax), bitwidth: b
 - \circ Compute the scale: $s=(L_{max}-L_{min})/(2^b-1)^{-1}$
 - \circ Clip the input x: $x_c = Clip(x, L_{min}, L_{max})$
 - Calculate the fixed-point representation:

$$x_{int} = round((x_c - L_{min})/s)$$

 \circ Rescale: $x_q = sx_{int} + L_{min}$





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: (Lmin, Lmax), bitwidth: b b=3, Lmax=2, Lmin=-0.5
 - \circ Compute the scale: $s=(L_{max}-L_{min})/(2^b-1)$ s = 0.357
 - Clip the input x: $x_c = Clip(x, L_{min}, L_{max})$ Xc = [1.1, 2, -0.3, 0.8]
 - Calculate the fixed-point representation:

$$x_{int} = round((x_c - L_{min})/s)$$
 Xint = [4,7,1,4]

 \circ Rescale: $x_q = sx_{int} + L_{min}$ Xq = [0.93, 2.0, -0.14, 0.93]



Addition/Subtraction:

$$x_q + y_q = s(x_{int} + y_{int}) + 2L_{min} \hspace{5mm} x_q - y_q = s(x_{int} - y_{int})$$

Multiplication (needs additional computation):

$$x_q imes y_q = s_x s_y (x_{int} imes y_{int}) + L_{min,x} y_q s_y + L_{min,y} x_q s_x + L_{min,x} L_{min,y}$$

Division: hard to implement



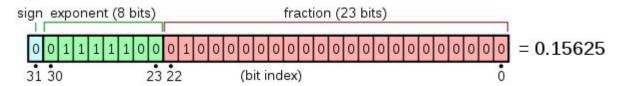
Floating-Point Arithmetic



- 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:

$$x = (-1)^s imes 2^{e-bias} imes (1+m) \ m = (0.b_0b_1b_2\dots b_{22})_2$$

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

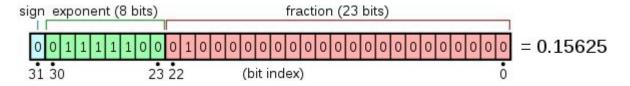
For example:

$$5.5 = (-1)^0 \times 2^{129-127} \times (1.011)_2$$
 $s = 0, e = 10000001, m = 0110000...0$

$$\circ$$
 -71 = (-1)¹ × 2¹³³⁻¹²⁷ × (1.000111)₂ s = 1, e = 10000101, m = 0001110...0

$$0.34375 = (-1)^0 \times 2^{125-127} \times (1.011)_2$$
 s = 0, e = 01111101, m = 0110000...0

Floating-Point Arithmetic



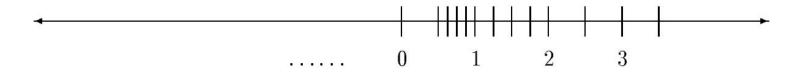
IEEE-754 standard:

$$x = (-1)^s imes 2^{e-bias} imes (1+m) \ m = (0.b_0b_1b_2\dots b_{22})_2$$

- 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?

$$egin{aligned} & \mathsf{x} = |\mathsf{x}| & \mathsf{s} = \mathsf{sign}(\mathsf{x}) \ & a = \lfloor log_2 x
floor & e = a + bias & m = rac{x}{2^a} - 1 \end{aligned}$$



Example

```
x = -13.24, bias=127 

x = |x| s = sign(x) a = \lfloor log_2 x \rfloor e = a + bias m = \frac{x}{2^a} - 1 

a = 3, e = 130, m = 0.655 

s = (1)_2, e = (100000010)_2, m = (10100111101011100001000)_2
```



Difference in Representation Power Between INT and FP

- FP provides relative precision that scales with magnitude. Small numbers near zero have finer granularity, while very large numbers have coarser steps.
 - For FP, the as the magnitude getting larger, the granularity will also decrease.

 Under the fixed exponent, mantissa

$$\stackrel{\circ}{} x=(-1)^s imes 2^{e-bias} imes \overline{(1+m)}$$

As exponent increase, the granularity will decrease exponentially

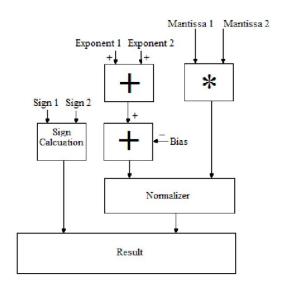
fills the gap evenly

 INT provides uniform precision. Each step between representable values is exactly the same.



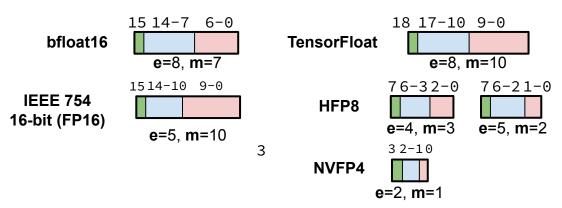
Computation with FP Representation

- Addition/Subtraction:
 - Need to align the exponent
 011010 + 001111 = 011010 + 011011 = 011101
 s₁e₁ m₁ s₂ e₂ m₂ Alignment
- Multiplication/Subtraction:
 - Sum the exponent, multiply the mantissa 011010×001111 $e = e_1 + e_2$ $s_1e_1 m_1$ $s_2 e_2 m_2$ $1+m = normalizer(1.m_1 x 1.m_2)$
- Addition and subtraction is expensive for FP.





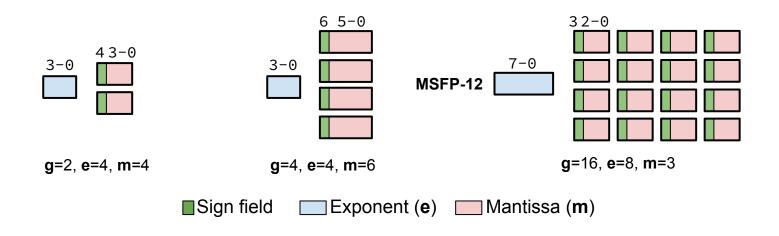
Customized FP Representation



- Numerous customized FP representations have been developed to facilitate DNN execution.
- FP can be applied to a block of numbers, with the block size defined in a customizable manner.



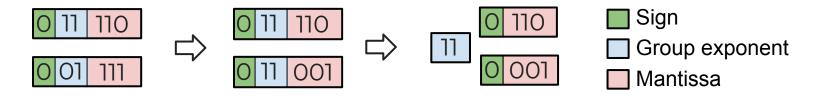
Block Floating Point (BFP)



• 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)



- Block floating point (BFP) is a numerical representation method that applies a shared exponent to a block of fixed-point values, balancing precision and dynamic range while reducing computational complexity compared to full floating-point arithmetic.
- There is no "leading 1".

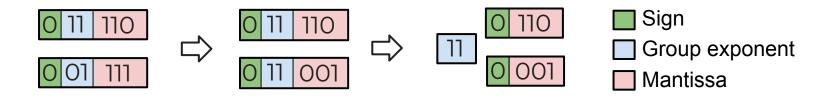
$$x=(-1)^s imes 2^{e-bias} imes (1+m) \ m=(0.b_0b_1b_2\ldots b_{22})_2$$

$$x = (-1)^s \times 2^{e-bias} \times m$$
 $m = (b_0.b_1b_2b_3...b_{22})_2$

BFP



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.

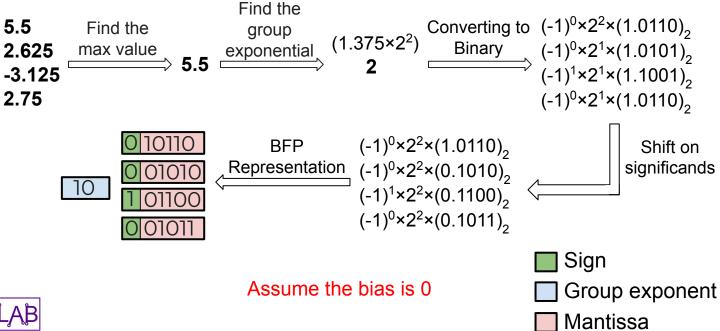
$$x=(-1)^s imes 2^{e-bias} imes (1+m) \ m=(0.b_0b_1b_2\ldots b_{22})_2$$

$$x = (-1)^s \times 2^{e-bias} \times m$$
 $m = (b_0.b_1b_2b_3...b_{22})_2$

BFP



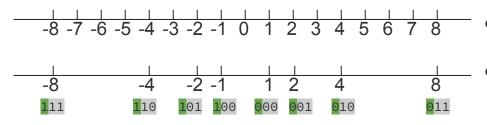
Example





Logarithm Arithmetics

- A specialized form of integer (INT) quantization
- Utilizes only power-of-two integer values, making hardware multiplication more efficient and cost-effective.



- Each INT number can be represented by its exponent = log(INT).
- A total of 8 numbers, 3 bits are needed to encode the bits.

$$a = (1100)_2$$

$$a \times 2 = (11000)_2$$

$$a = (1100)_2$$
 $a \times 2 = (11000)_2$ $a \times 8 = (1100000)_2$



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

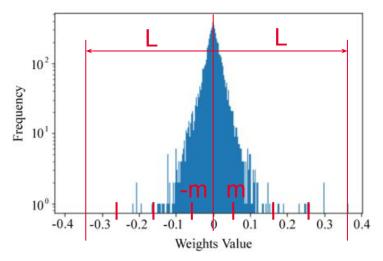


Taxonomy of Quantization

- 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

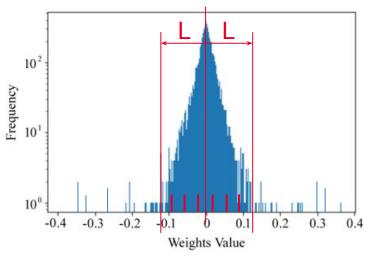


Weight distribution in ResNet

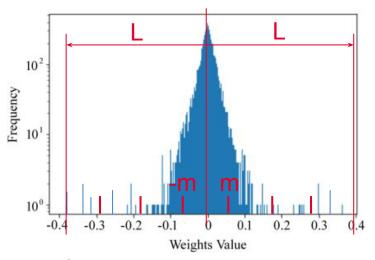
- The weight distribution follows a gaussian-like distribution.
- The outlier will lead to large quantization error.
- A good selection on the clip range L is critical for accuracy performance.



Weight Quantization



- Large truncation error
- Low quantization error for small values



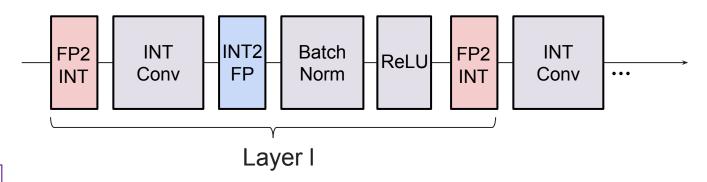
- Small truncation error
- Large quantization error for small values



 $L = 0.9 \times max(|W|), L = 0.95 \times max(|W|), 0.9 and 0.95 are chosen by experience.$

Activation Quantization

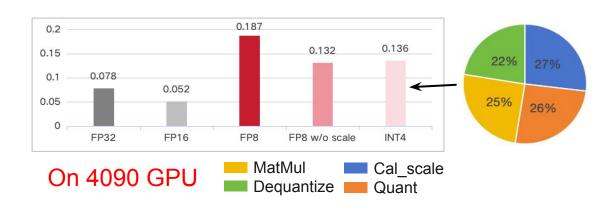
- Quantization on activation needs to be performed dynamically. This will introduce additional compute overhead.
- Also the activation will pass the nonlinear functions, which are usually very sensitive to quantization error, so dequantization is required to convert back to FP 16/32.





Activation Quantization

(577×1024)× (1024×1024) Projection Layer: Input: 577x1024 Weight: 4096x1024



- For low-precision quantization, the quantization process may cause more computation than the computational savings achieved by using low-precision quantization.
- Set the clipping range: (-L, L), bit width: b
- Compute the scale: $s = 2L/(2^b 2)$
- Clip the input x: $x_c = Clip(x, L, -L)$
- Calculate the INT representation: $x_{int} = round(x_c/s)$
- \circ Rescale: $x_q = sx_{int}$



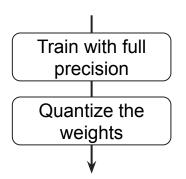
Taxonomy of Quantization

- 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

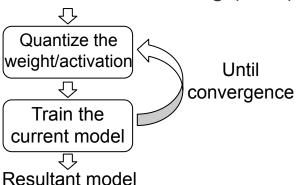


When to Quantize?

Post-training quantization (PTQ)



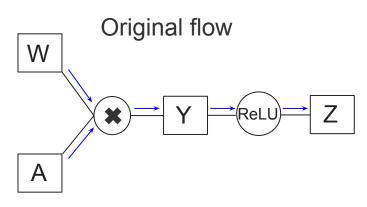
Quantization-aware Training (QAT)



- 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.



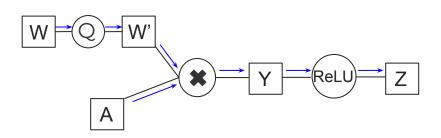
Another Way to Look at Quantization



$$Y = WA, Z = ReLU(Y)$$

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

Flow with quantization

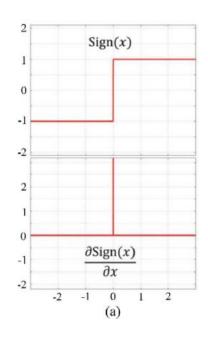


$$\frac{\partial L}{\partial W} = \frac{\partial L}{\partial Z} \frac{\partial Z}{\partial Y} \frac{\partial Y}{\partial W'} \frac{\partial W'}{\partial W}$$

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 makes 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.

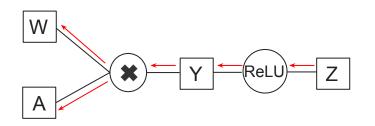


Straight Through Estimator (STE)

Forward pass

W Q W' A Q Y ReLU Z

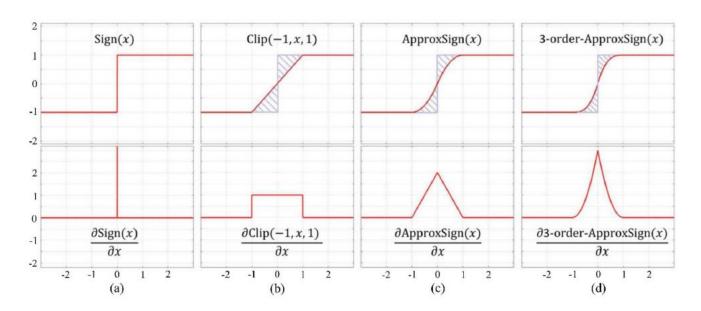
Backward pass



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



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 wa
```



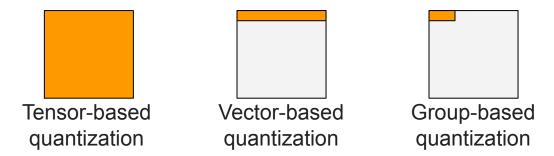
Taxonomy of Quantization

- 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.





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



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.



Data gradient Computation

$$\nabla \mathbf{Y} \times \mathbf{W}^{\mathsf{T}} = \nabla \mathbf{X}$$

Weight gradient Computation

X: input

∇X: input gradient

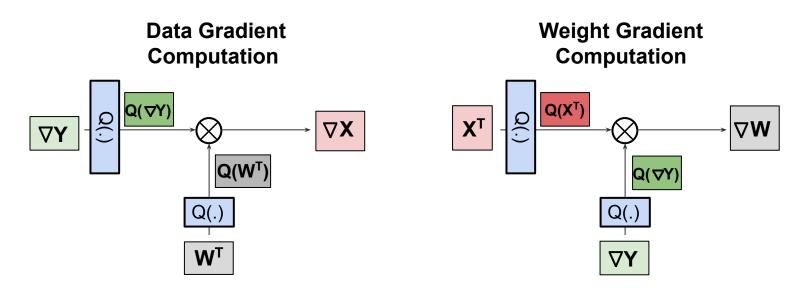
W: weight filters

∇W: weight gradient

Y: output

∀Y: output gradient

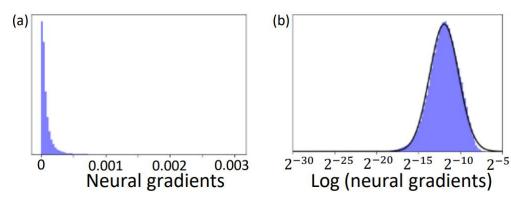




Gradient is much more sensitive to quantization error.



DNN Gradient Distribution



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



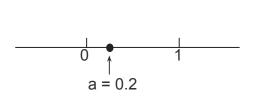
Log (neural gradients)

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



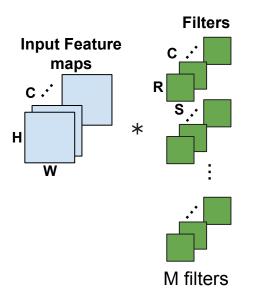
- To quantize a, conventional linear quantization will make q(a) = 0. However, this will cause a bias.
- With stochastic quantization:

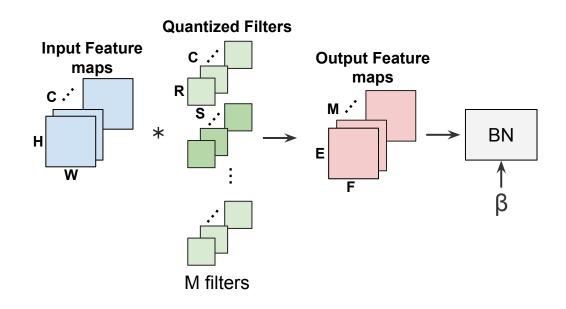
$$q(a) = egin{cases} 1 & ext{for } p = 0.2 \ 0 & ext{for } p = 0.8 \end{cases}$$

- For QAT, 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



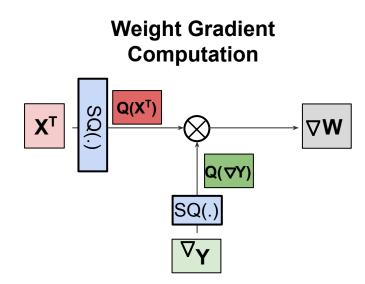






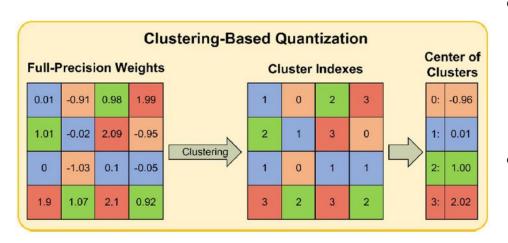


Data Gradient Computation VY Q(VY) Q(WT) SQ(.) WT





Clustering-Based Quantization



- Quantization and clustering share similarities. In clustering, each value is assigned to a centroid, while in quantization, each full-precision value is mapped to one of the predefined quantization levels.
- Clustering is usually used to compress the weight matrix and efficient storage, but it is hard for accelerating computations.
- However, due to the flexibility of selecting the centroid, clustering usually achieves a better accuracy.



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



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

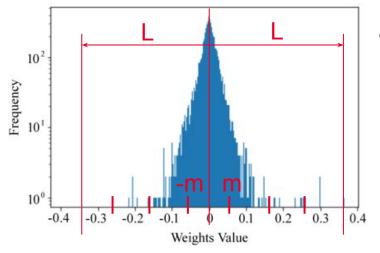


- How to convert a number to INT8 representation?
 - Set the clipping range: (-Lmin, Lmax), bitwidth: b
 - \circ Compute the scale: $s=(L_{max}-L_{min})/(2^b-1)^{-1}$
 - \circ Clip the input x: $x_c = Clip(x, L_{min}, L_{max})$
 - Calculate the fixed-point representation:

$$egin{aligned} x_{int} = round((x_c - L_{min})/s) \end{aligned}$$

 \circ Rescale: $x_q = sx_{int} + L_{min}$



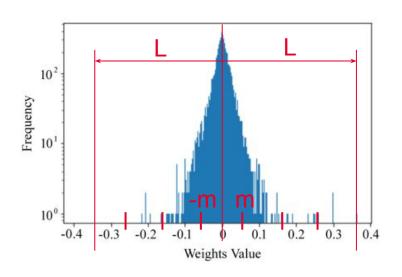


Weight distribution in ResNet

- How to convert a number to INT8 representation?
 - Set the clipping range: (-I, I), bitwidth: b
 - \circ Compute the scale: $s = (2l)/(2^b-1)$
 - Olip the input x: $x_c = Clip(x, l, -l)$
 - Calculate the fixed-point representation:x_{int} = round(x_c/s)
 - Rescale: xq = sxint

I = 0.9×max(|W|), I = 0.95×max(|W|)
Can learn by learnt during training?



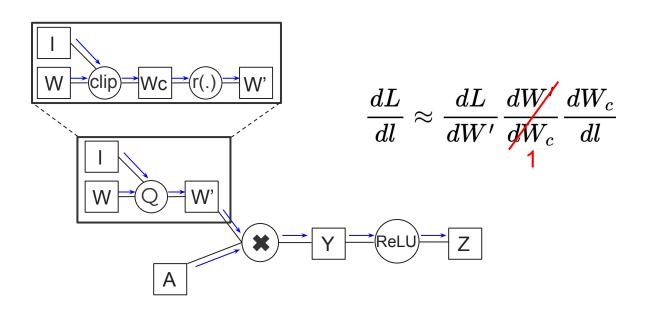


First we need to apply CLIP function to the input x, where the clip function has a range of (-I, I).

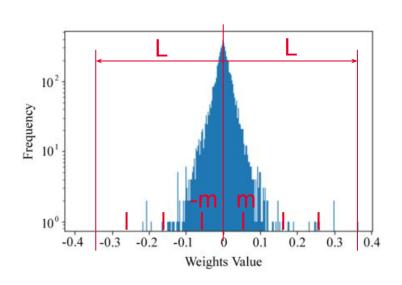
$$egin{aligned} oldsymbol{x}_c &= Clip(x,l) = egin{cases} l, & ext{if } x \geq l \ x, & -l \leq x \leq l \ -l, & x \leq l \end{cases} \ oldsymbol{x}_q &= round(rac{x_c}{s}) imes s \end{cases}$$

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







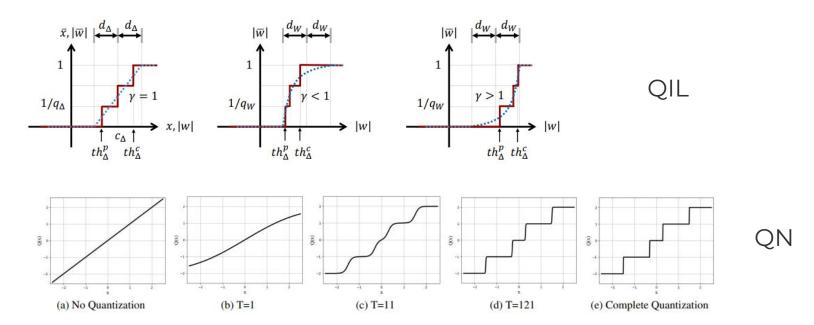


$$Clip(x,l) = egin{cases} l, & ext{if } x \geq l \ x, & -l \leq x \leq l \ -l, & x \leq l \end{cases}$$

$$rac{dClip(x,l)}{dx} = egin{cases} 0, & ext{if } x \geq l \ 1, & -l \leq x \leq l \ 0, & x \leq l \end{cases}$$

$$rac{dClip(x,l)}{dl} = egin{cases} 1, & ext{if } x \geq l \ 0, & -l \leq x \leq l \ -1, & x \leq l \end{cases}$$
 L can be learnable





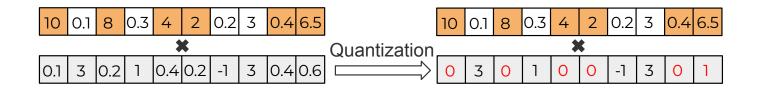


Jung, Sangil, et al. "Learning to quantize deep networks by optimizing quantization intervals with task loss." *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*. 2019.

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

Problem of Quantization

 The major drawback of quantization is that it does not consider the impact of the input when making the quantization decision.





Problem of Quantization

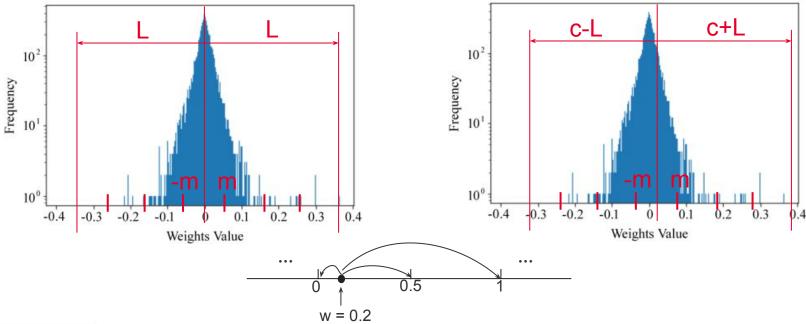


- We use a calibration dataset and profile some data xi, yi across each layer.
- After that, we use these data to train the optimal quantized weights.

$$\min_{W_l} ||Y_l - X_l Q(W_l)||^2$$
 For each I



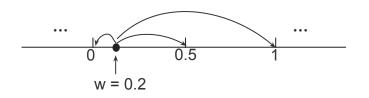
Quantization Interval Learning (QIL)





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.

Quantization Interval Learning (QIL)



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

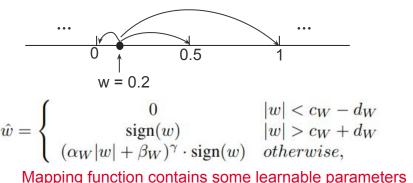
$$w_q = Q(w) \longrightarrow w_q = Q(F(w))$$

 F(.) is a function which contains learnable hyperparameters.

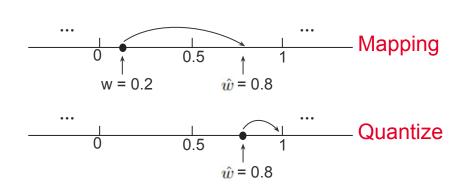
$$\hat{w} = \begin{cases} 0 & |w| < c_W - d_W \\ \operatorname{sign}(w) & |w| > c_W + d_W \\ (\alpha_W |w| + \beta_W)^{\gamma} \cdot \operatorname{sign}(w) & otherwise, \end{cases}$$

Quantization Interval Learning (QIL)

QIL offers flexibility to round the FP weights.



Mapping function contains some learnable parameters

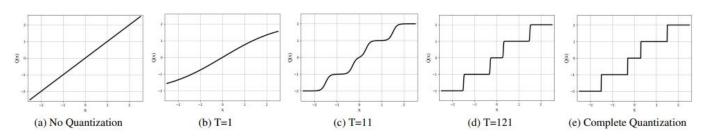


- $w_q = Q(F(w))$ are stored for inference after the training process finished.
- We can not apply this techniques over the activation, due to its large computational overhead.



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.



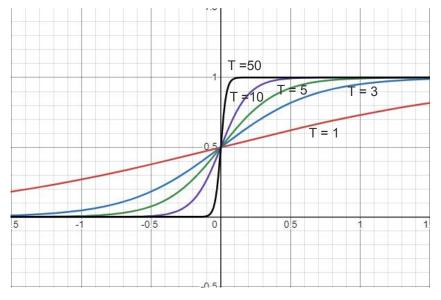
- $y = \alpha(\sum_{i=1}^{n} s_i \mathcal{A}(\beta x b_i) o)$
 - $\mathcal{A}(x) = \begin{cases} 1 & x \ge 0, \\ 0 & x < 0. \end{cases}$
- n + 1 is the number of quantization intervals
- β is the scale factor of inputs
- si and bi are the scales and biases for the unit step functions



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

Quantization Networks

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



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



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

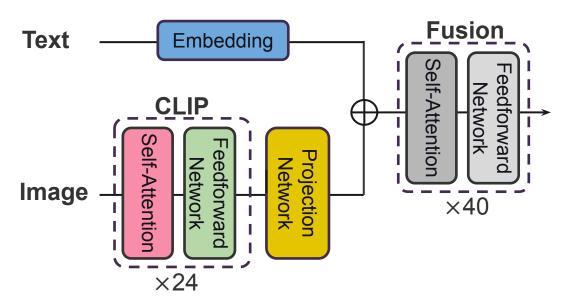


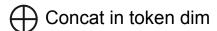
Post Training Quantization

- 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.



Case Study: CLIP in Llava

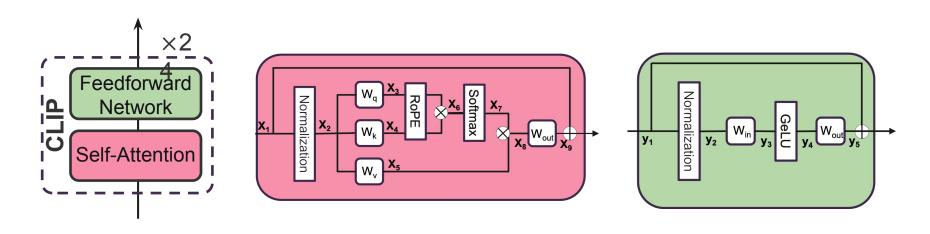




Visual input example, Tricky Question:	
User	If there are factual errors in the questions, point it out if not, proceed answering the question. What's happening in the desert?
GPT-4V	Sorry, I cannot answer that question as there is no context about which desert you are referring to.
LLaVA	In the image, 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.



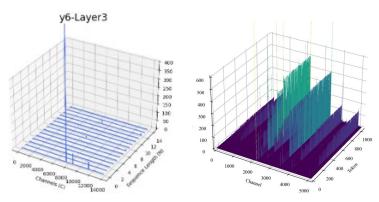
CLIP Architecture





Types of Outlier

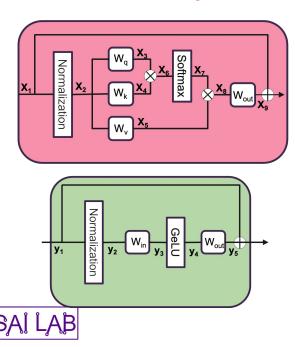
- Massive Activation:
 - For an activation matrix A, an massive activation is an element Aij within it that satisfies:
 - \circ Aij > $\eta \times mean(|A|)$
 - Aij > γ
 - o η=300, γ=50
- Channelwise Outlier:
 - \circ mean(Ai) > $\eta \times std(A) + mean(|A|)$
 - \circ std(Ai) < β
 - \circ $\eta = 3, \beta = 0.6$

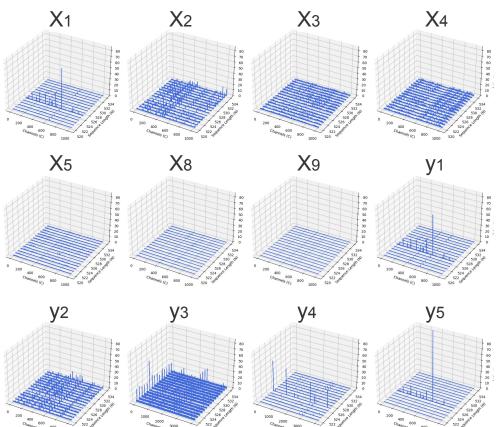




Outlier Study: CLIP Activations

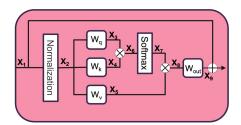
3D activation within layer 12



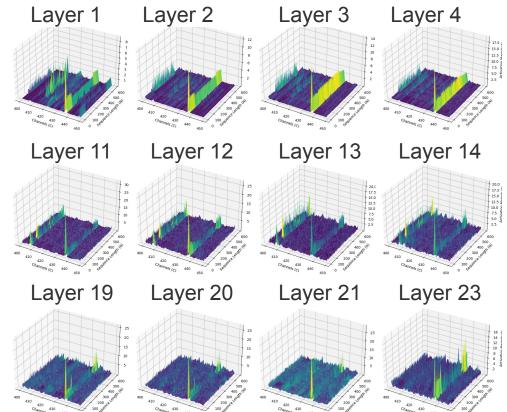


Outlier Study: CLIP Activations

 3D plots of x2 across layers.



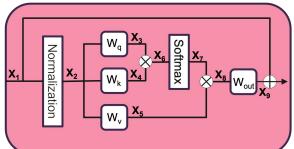
 x2 exhibits channel wise outlier



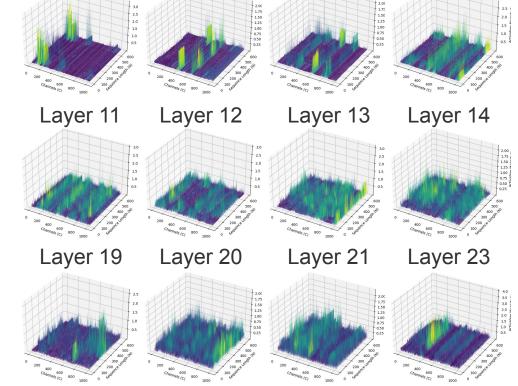


Outlier Study: CLIP Activations Layer 1 Layer 2 Layer 3

 3D plots of x8 across layers.



 x8 exhibits channel wise outlier

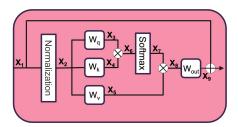


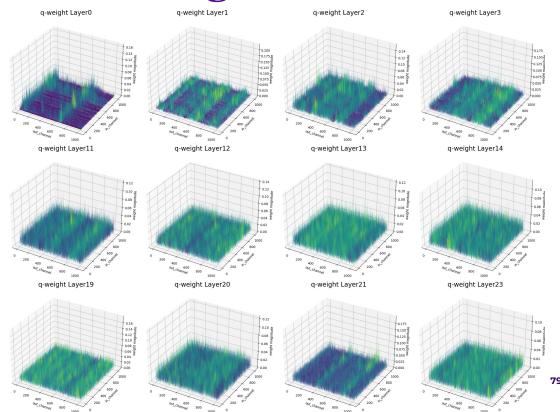


Layer 4

Outlier Study: CLIP Weights

Wq across CLIP layers.

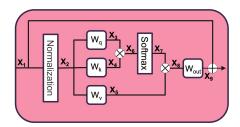


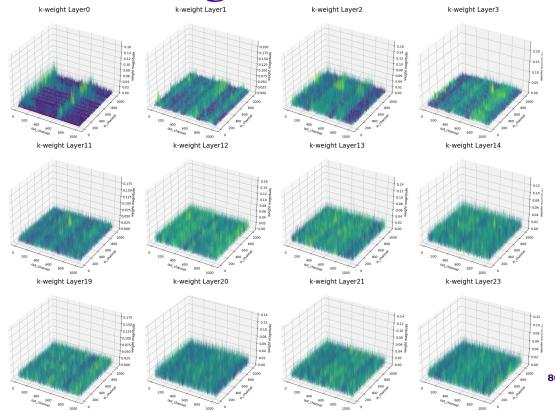




Outlier Study: CLIP Weights

Wk across CLIP layers.

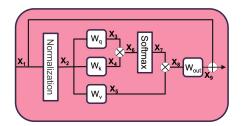


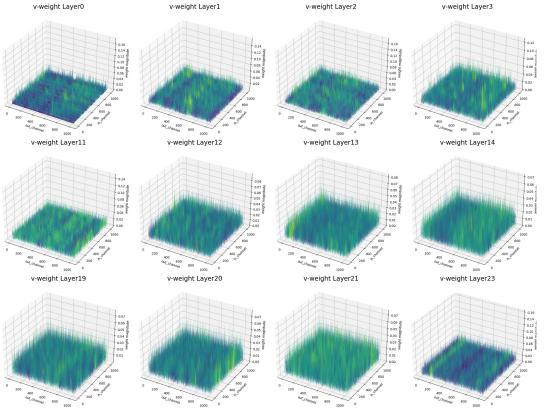




Outlier Study: CLIP Weights

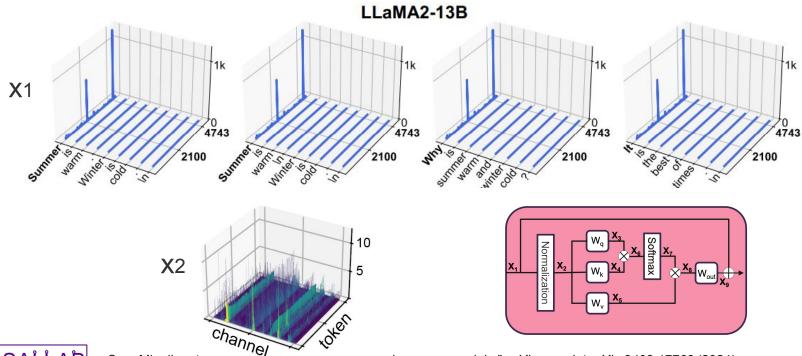
W_V across CLIP layers.







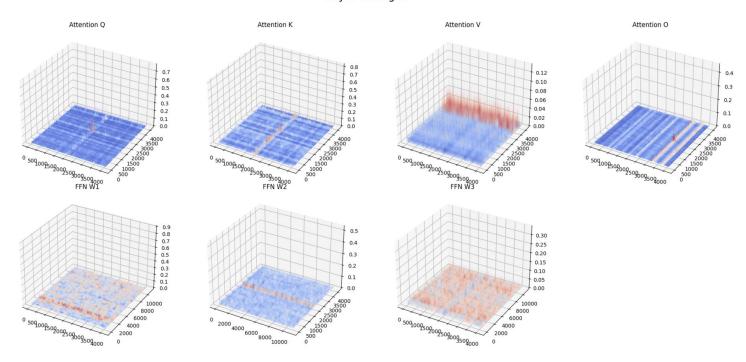
Outlier Study: LLaMA Activations



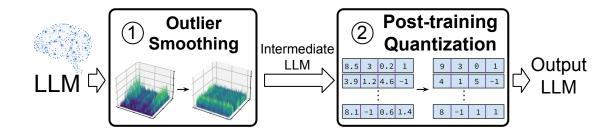


Outlier Study: LLaMA Weights

Layer 0 Weights



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.

