

# **CS 260D: Large-scale Machine Learning**

**Lecture #14 & #15: Neural Network Quantization**

Baharan Mirzasoleiman  
UCLA Computer Science

# Neural Network Quantization

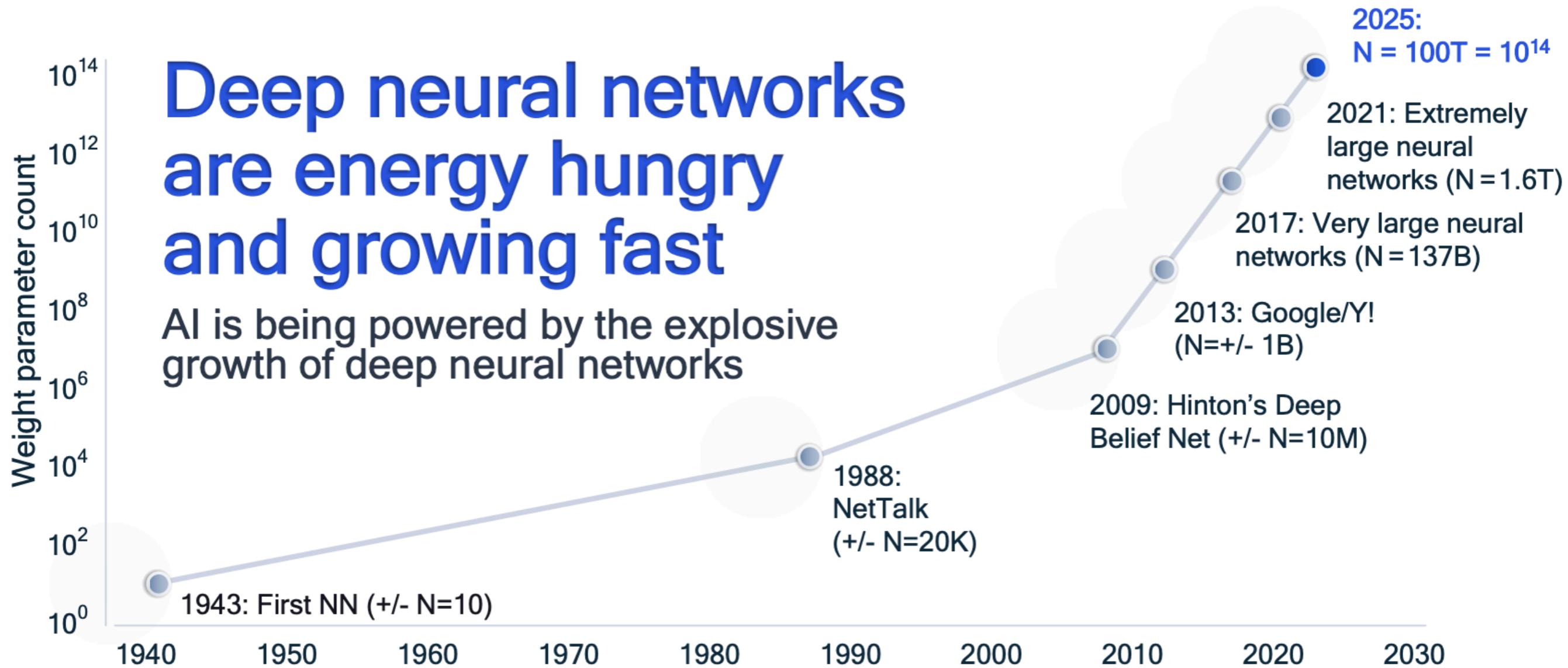
# Outline

- Energy-efficient machine learning and the need for quantization
- Introduction to neural network quantization
- Simulating quantization in neural networks
- Post-training quantization (PTQ)
- Quantization-aware training (QAT)

From Marios Fournarakis

# Deep neural networks are energy hungry and growing fast

AI is being powered by the explosive growth of deep neural networks



2025

Increasingly large and complex neural networks for Natural Language Processing, Image and Video Processing

# The AI power and thermal ceiling

## The challenge of AI workloads



- Very compute intensive
- Complex concurrencies
- Real-time
- Always-on

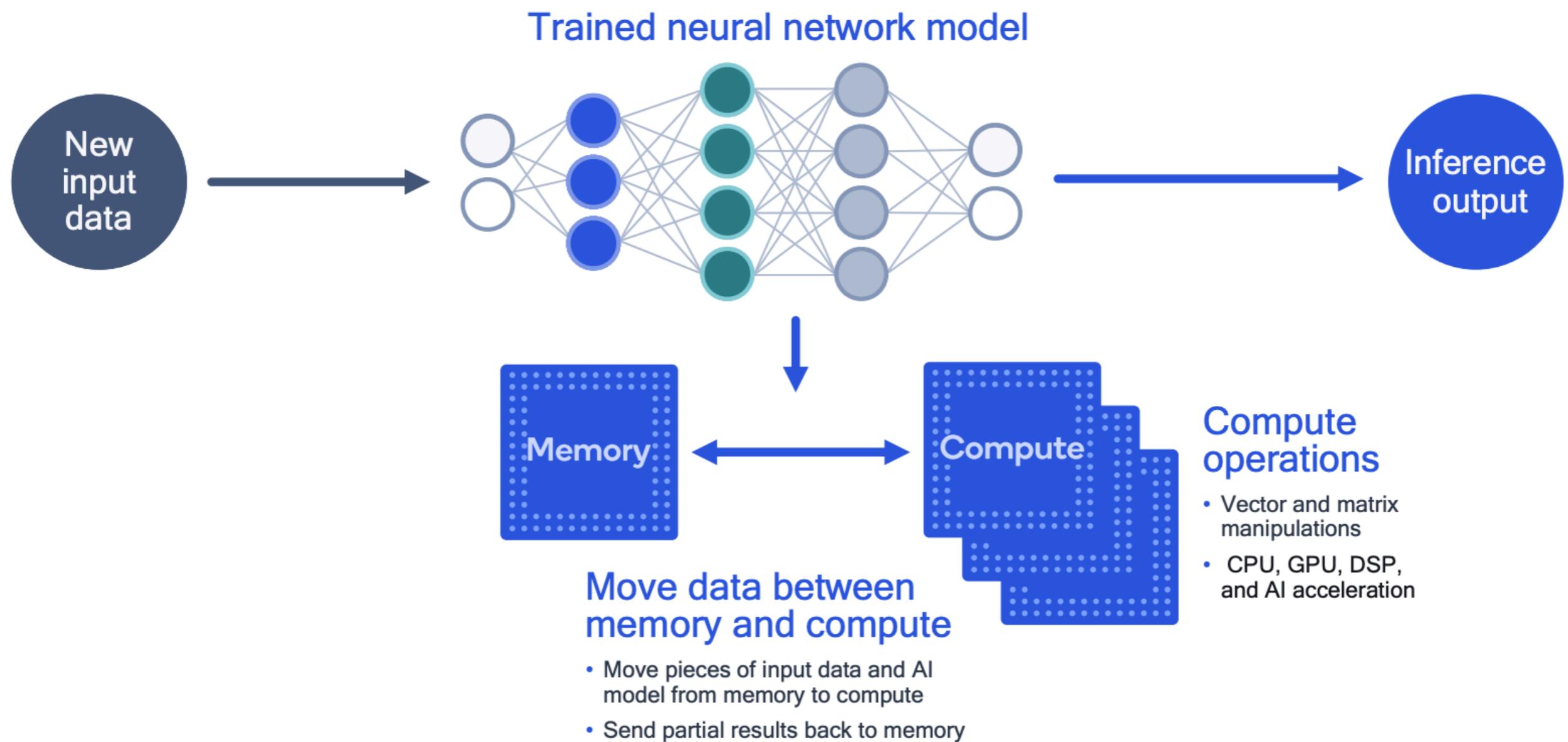
## Constrained mobile environment

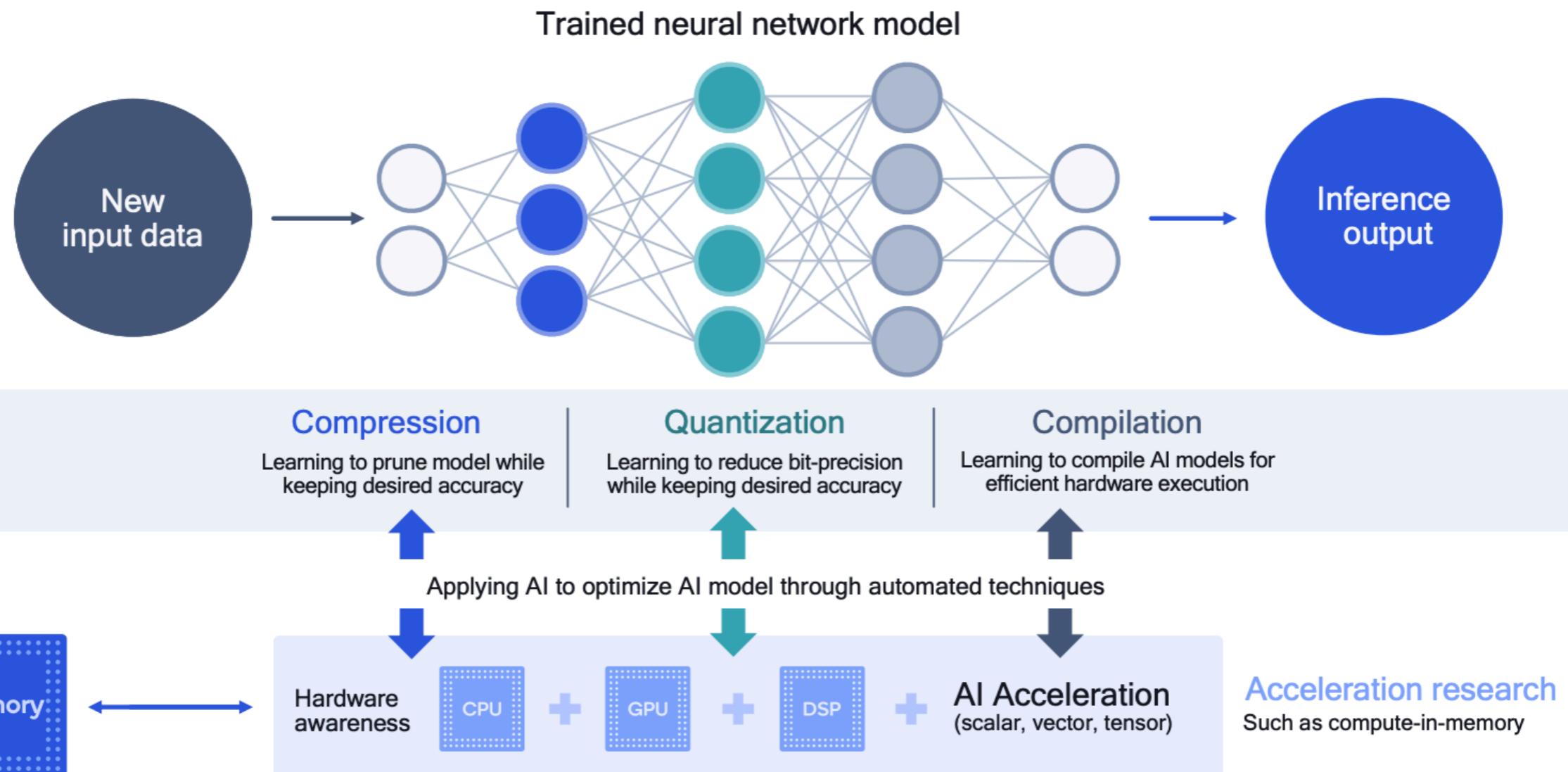


- Must be thermally efficient for sleek, ultra-light designs
- Requires long battery life for all-day use
- Storage/memory bandwidth limitations



# Advancing AI research to increase power efficiency





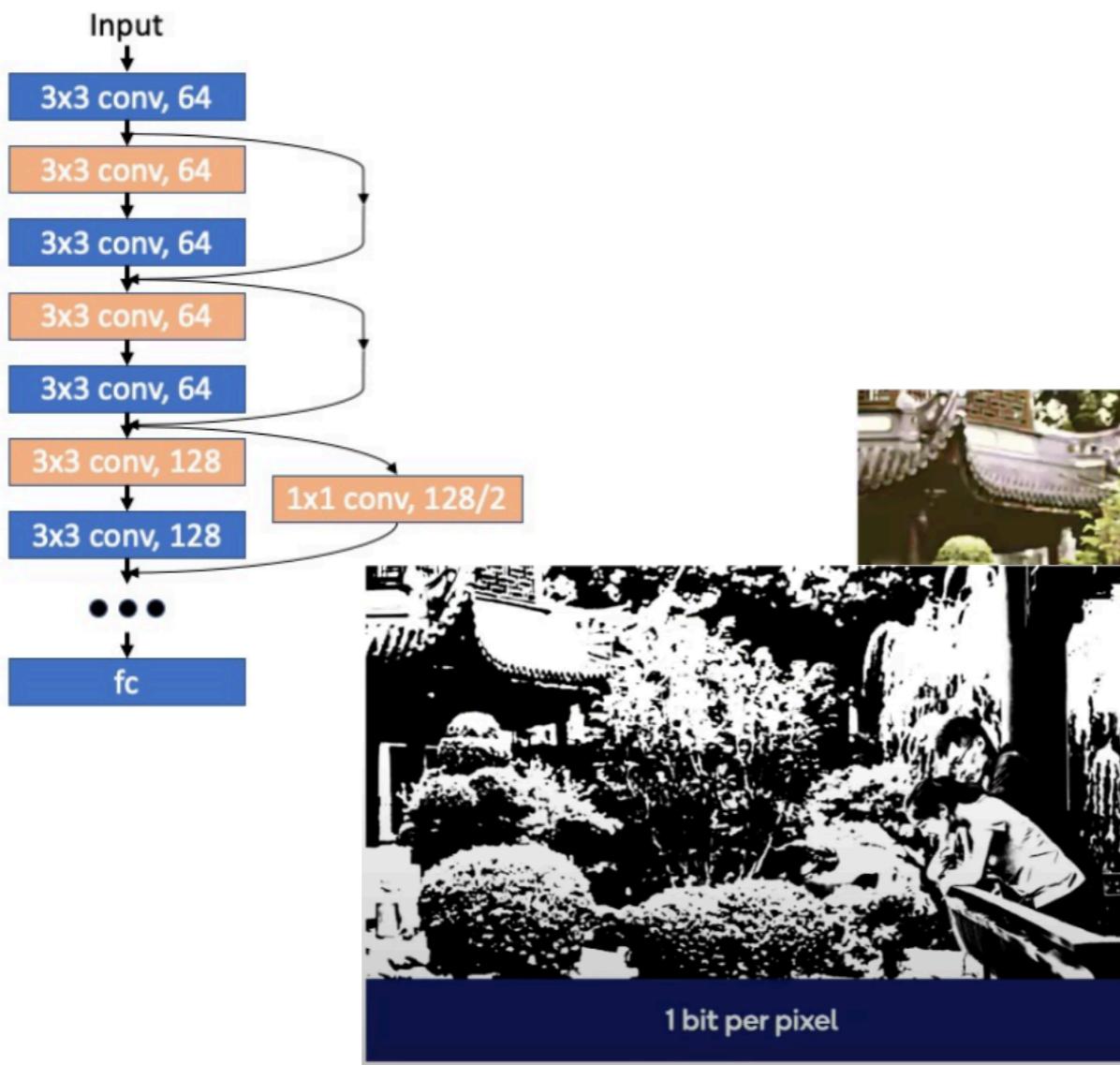
Advancing AI research to increase power efficiency

# What is Neural Network Quantization?

# What is neural network quantization?

For any given trained neural network:

- Store weights in low bits (INT8)
- Compute calculations in low bits



Quantization Analogy

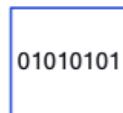
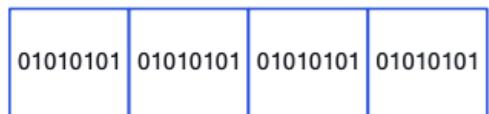
Use fewer bits to represent each pixel in an image



# Quantizing AI models offers significant benefits

## Memory usage

8-bit versus 32-bit weights and activations stored in memory



## Power consumption

Significant reduction in energy for both computations and memory access

Add energy (pJ)		Mem access energy (pJ)	
INT8	FP32	Cache (64-bit)	
0.03	0.9	8KB	10
<b>30X energy reduction</b>		32KB	20
		1MB	100
Mult energy (pJ)		DRAM	
INT8	FP32	1300-2600	
0.2	3.7	<b>Up to 4X energy reduction</b>	
<b>18.5X energy reduction</b>			

## Latency

With less memory access and simpler computations, latency can be reduced



## Silicon area

Integer math or less bits require less silicon area compared to floating point math and more bits

Add area ( $\mu\text{m}^2$ )	
INT8	FP32
36	4184
<b>116X area reduction</b>	
Mult area ( $\mu\text{m}^2$ )	
INT8	FP32
282	7700
<b>27X area reduction</b>	

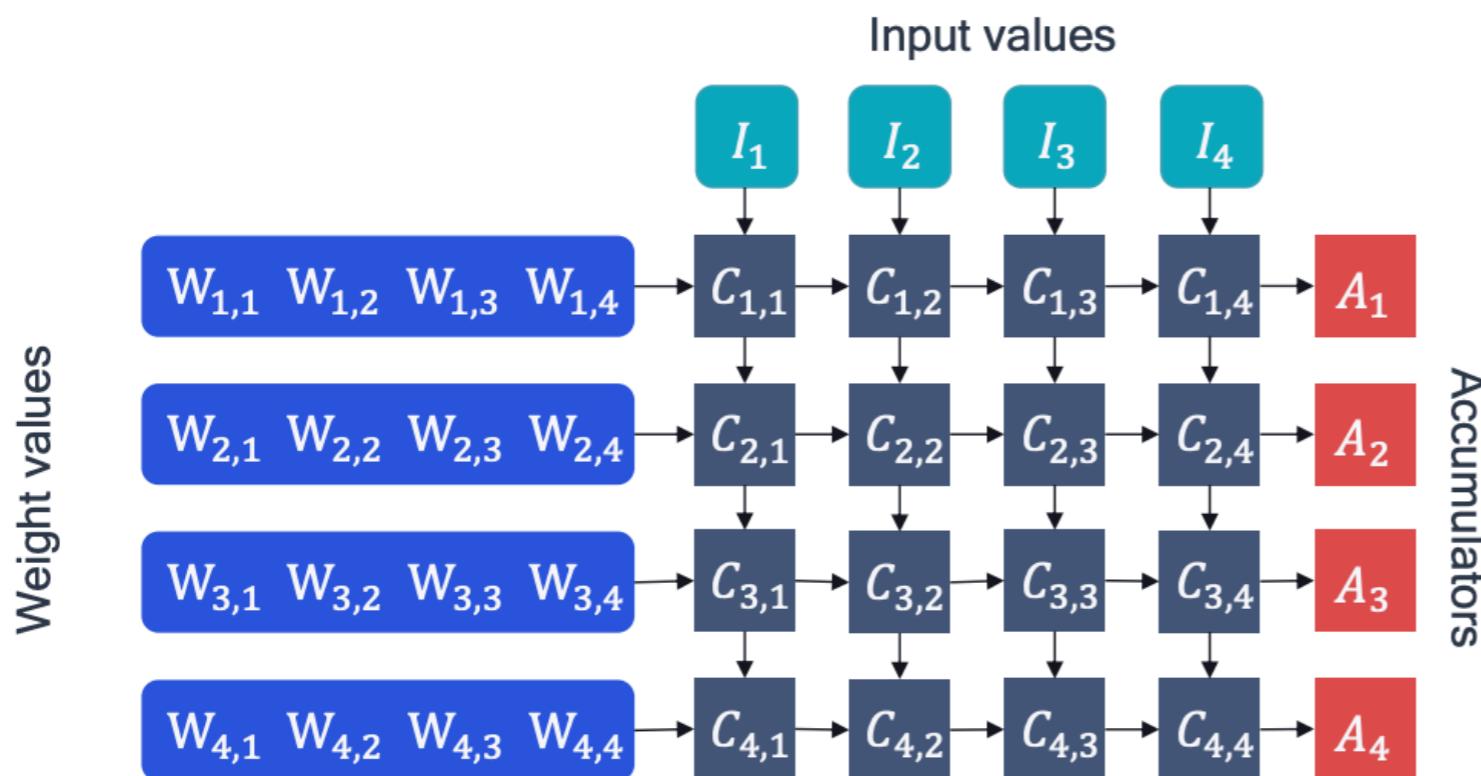
# Matrix operations are the backbone of neural networks

A running example to showcase how to make these operations more efficient

$$\mathbf{W} = \begin{pmatrix} 0.97 & 0.64 & 0.74 & 1.00 \\ 0.58 & 0.84 & 0.84 & 0.81 \\ 0.00 & 0.18 & 0.90 & 0.28 \\ 0.57 & 0.96 & 0.80 & 0.81 \end{pmatrix} \quad \mathbf{X} = \begin{pmatrix} 0.41 & 0.25 & 0.73 & 0.66 \\ 0.00 & 0.41 & 0.41 & 0.57 \\ 0.42 & 0.24 & 0.71 & 1.00 \\ 0.39 & 0.82 & 0.17 & 0.35 \end{pmatrix} \quad \mathbf{b} = \begin{pmatrix} 0.1 \\ 0.2 \\ 0.3 \\ 0.4 \end{pmatrix}$$

How to most efficiently calculate  $\mathbf{WX} + \mathbf{b}$ ?

# A schematic MAC array for efficient computation

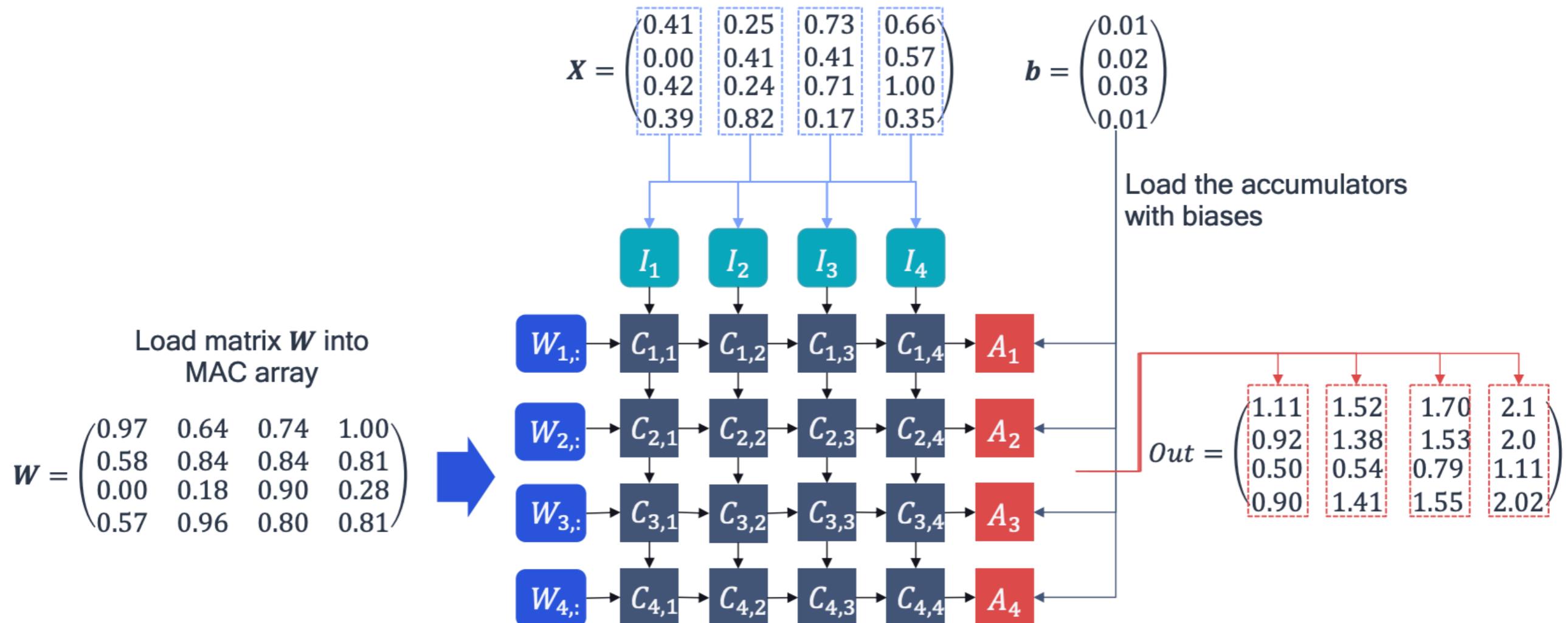


The array efficiently calculates the dot product between multiple vectors

$$A_i = \sum_j C_{i,j} + b_i$$

$$A_i = W_i \cdot I_1 + W_i \cdot I_2 + W_i \cdot I_3 + W_i \cdot I_4$$

# Step-by-step matrix multiplication in MAC array



# Quantization comes at a cost of lost precision

- We can approximate an FP tensor with an integer tensor multiplied by a scale-factor,  $s_X$ :

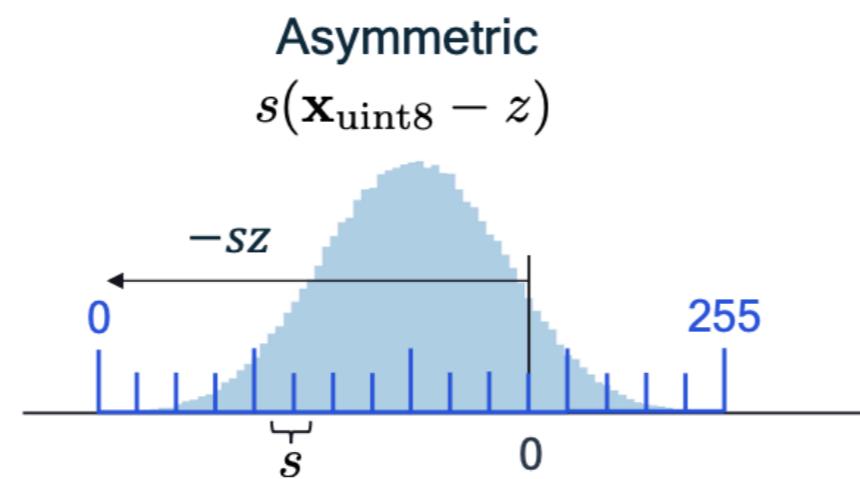
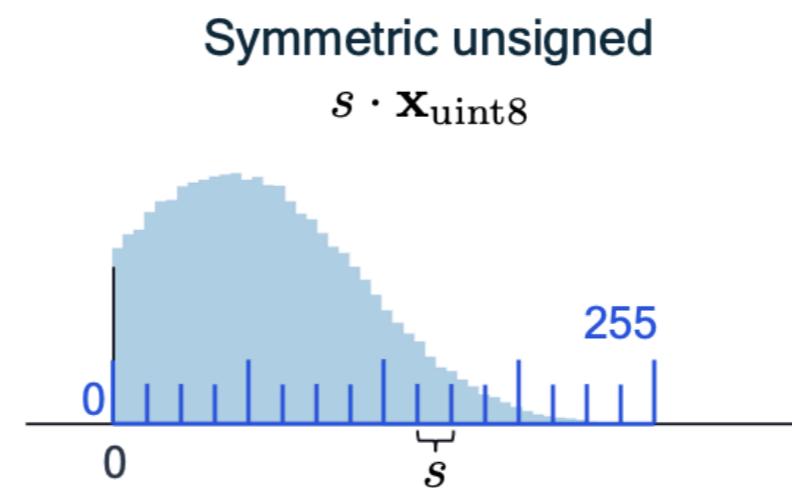
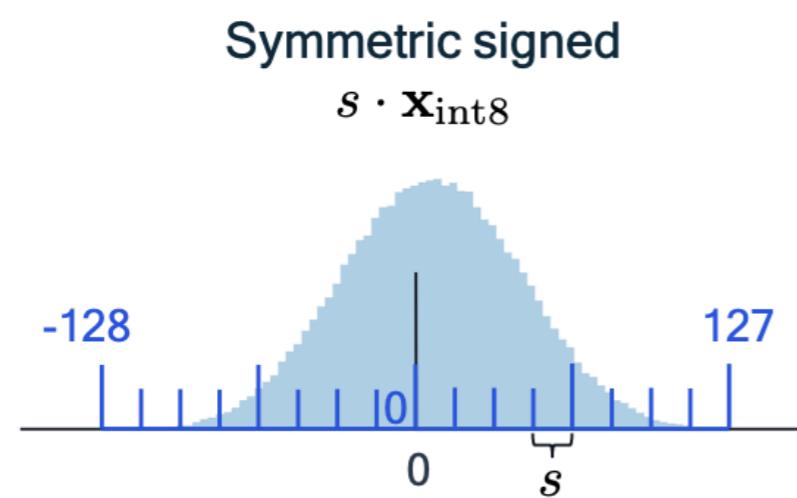
$$\text{FP32 tensor} \xrightarrow{\quad} X \approx s_X X_{\text{int}} = \hat{X} \xleftarrow{\quad} \text{scaled quantized tensor}$$
$$W = \begin{pmatrix} 0.97 & 0.64 & 0.74 & 1.00 \\ 0.58 & 0.84 & 0.84 & 0.81 \\ 0.00 & 0.18 & 0.90 & 0.28 \\ 0.57 & 0.96 & 0.80 & 0.81 \end{pmatrix} \approx \frac{1}{255} \begin{pmatrix} 247 & 163 & 189 & 255 \\ 148 & 214 & 214 & 207 \\ 0 & 46 & 229 & 71 \\ 145 & 245 & 204 & 207 \end{pmatrix} = s_W W_{\text{uint8}}$$

- Quantization is not free:

$$\epsilon = W - s_W W_{\text{int}} = \frac{1}{255} \begin{pmatrix} 0.35 & 0.20 & -0.3 & 0 \\ -0.1 & 0.20 & 0.20 & -0.45 \\ 0.00 & -0.1 & -0.5 & 0.40 \\ 0.35 & -0.2 & 0 & -0.45 \end{pmatrix}$$

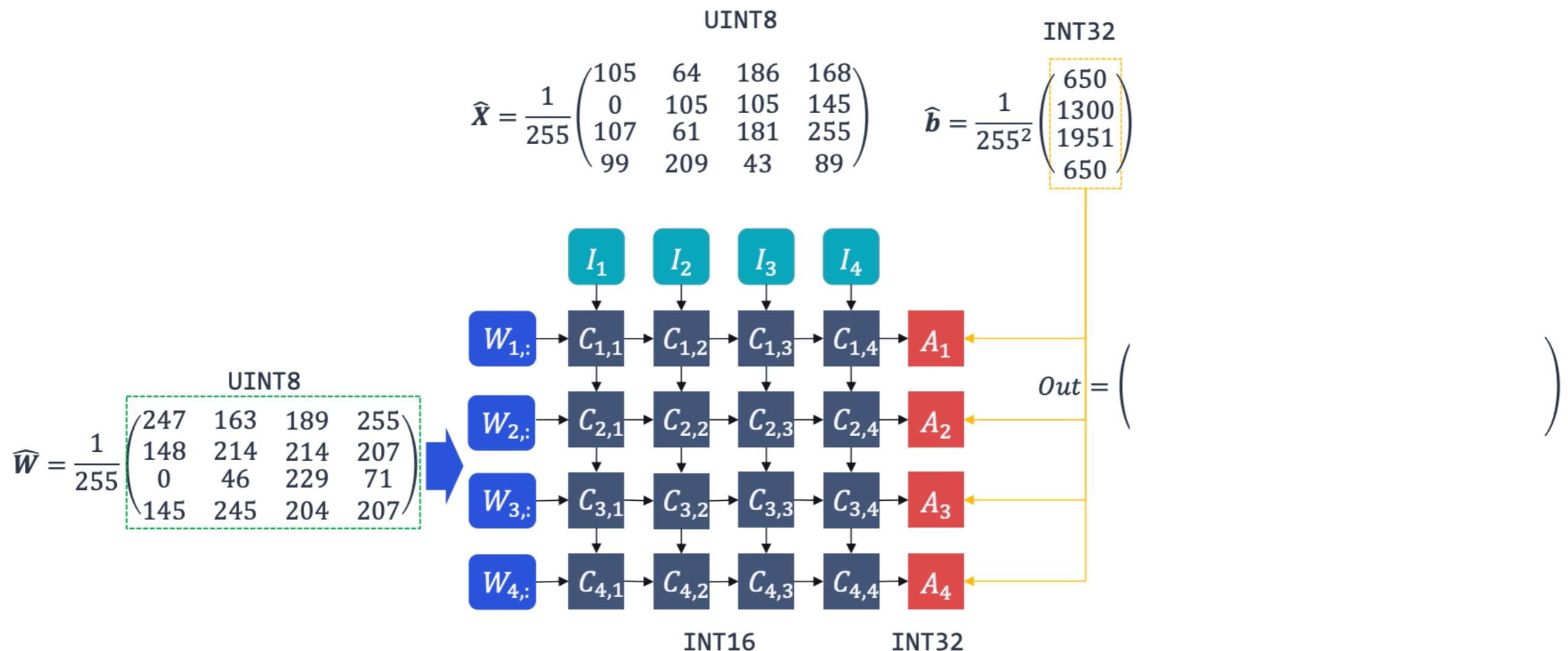
# Different types of quantization have pros and cons

Symmetric, asymmetric, signed, and unsigned quantization

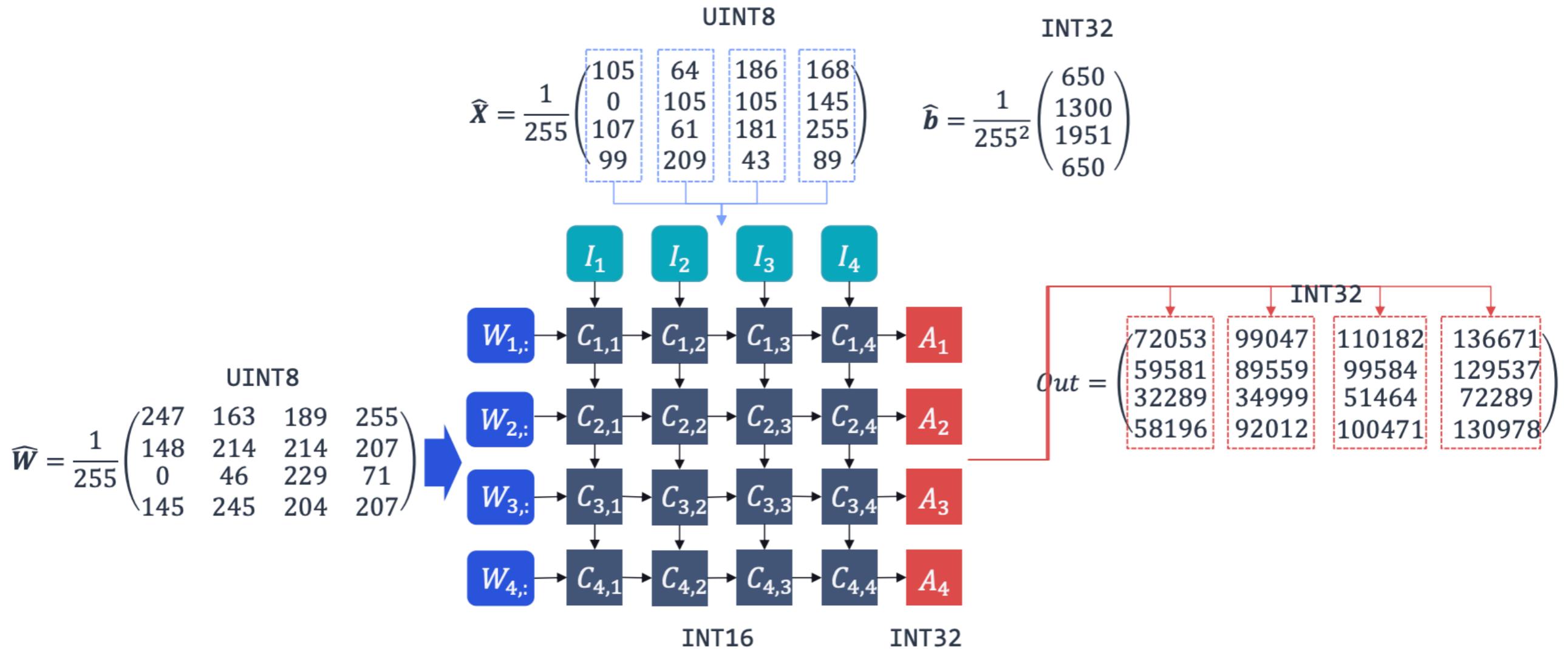


Fixed point grid  
Floating point grid  
 $s$ : scale factor  
 $z$ : zero-point

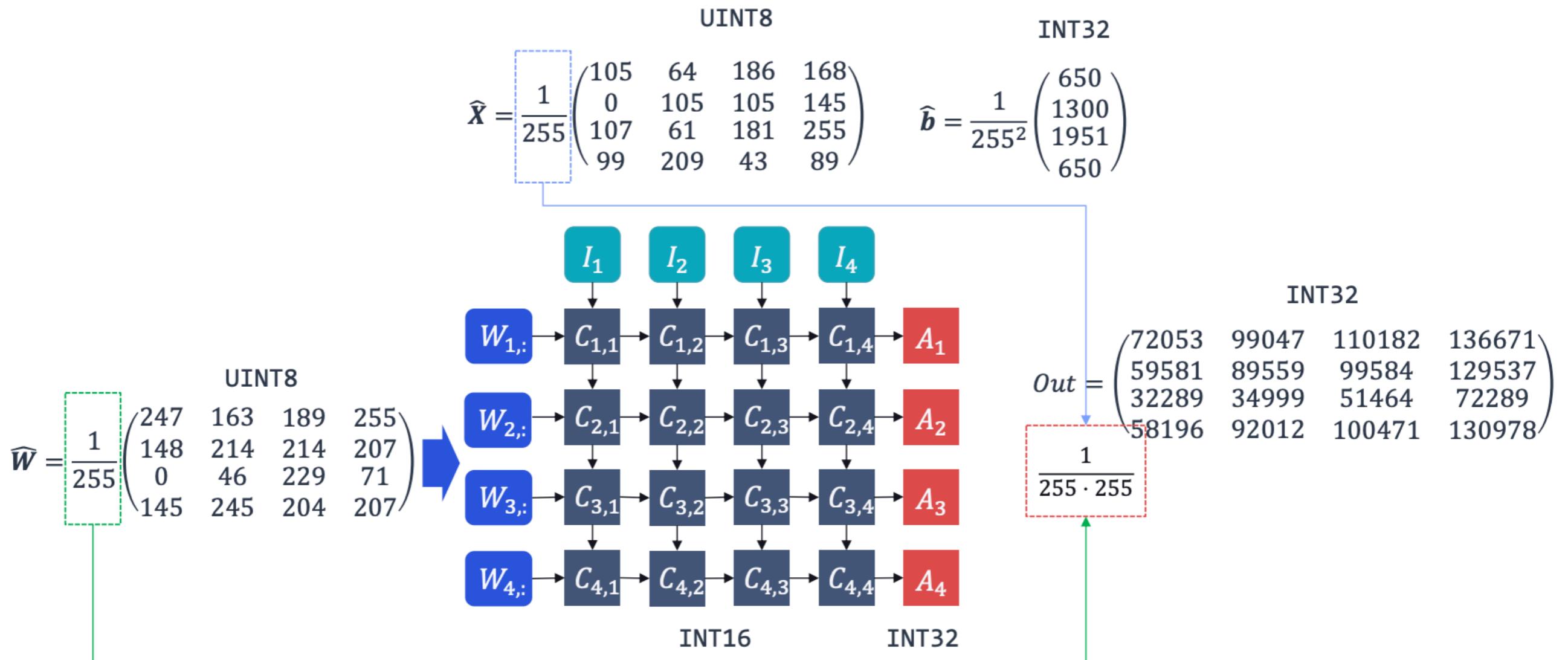
# Quantized inference using symmetric quantization



# Quantized inference using symmetric quantization

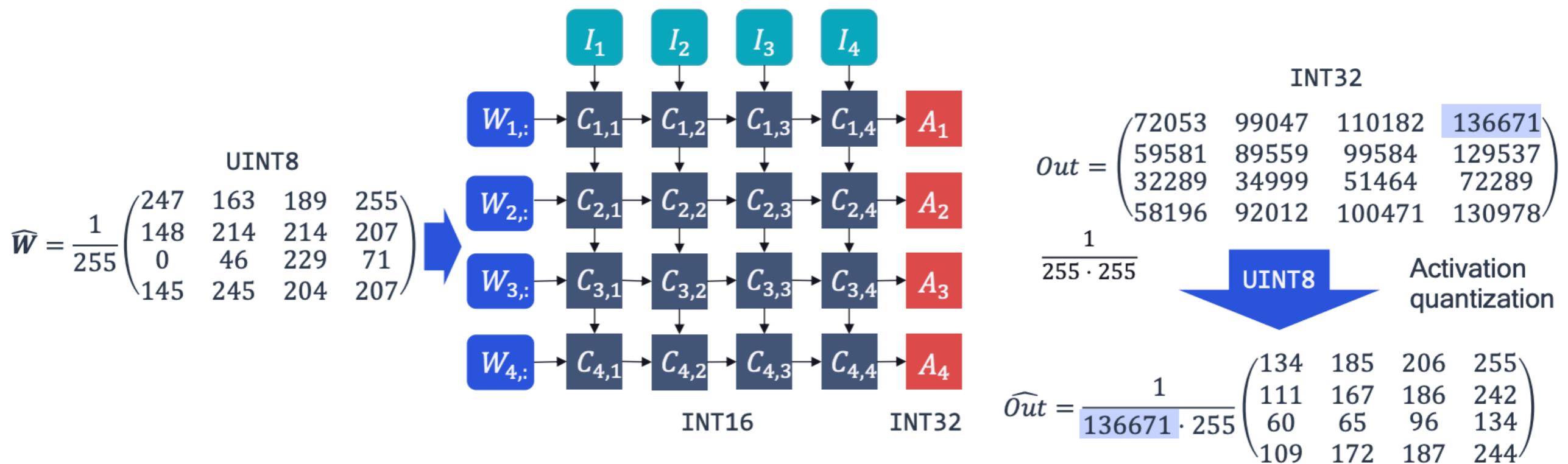


# Quantized inference using symmetric quantization



# Quantized inference using symmetric quantization

$$\widehat{\mathbf{X}} = \frac{1}{255} \begin{pmatrix} 105 & 64 & 186 & 168 \\ 0 & 105 & 105 & 145 \\ 107 & 61 & 181 & 255 \\ 99 & 209 & 43 & 89 \end{pmatrix} \quad \widehat{\mathbf{b}} = \frac{1}{255^2} \begin{pmatrix} 650 \\ 1300 \\ 1951 \\ 650 \end{pmatrix}$$



# What type of quantization should you use?

$W$  : weight matrix

$X$  : input of a layer

## Symmetric quantization

$$WX \approx s_W(W_{\text{int}}) s_X(X_{\text{int}})$$
$$= s_W s_X (W_{\text{int}} X_{\text{int}})$$

Same calculation

## Asymmetric quantization

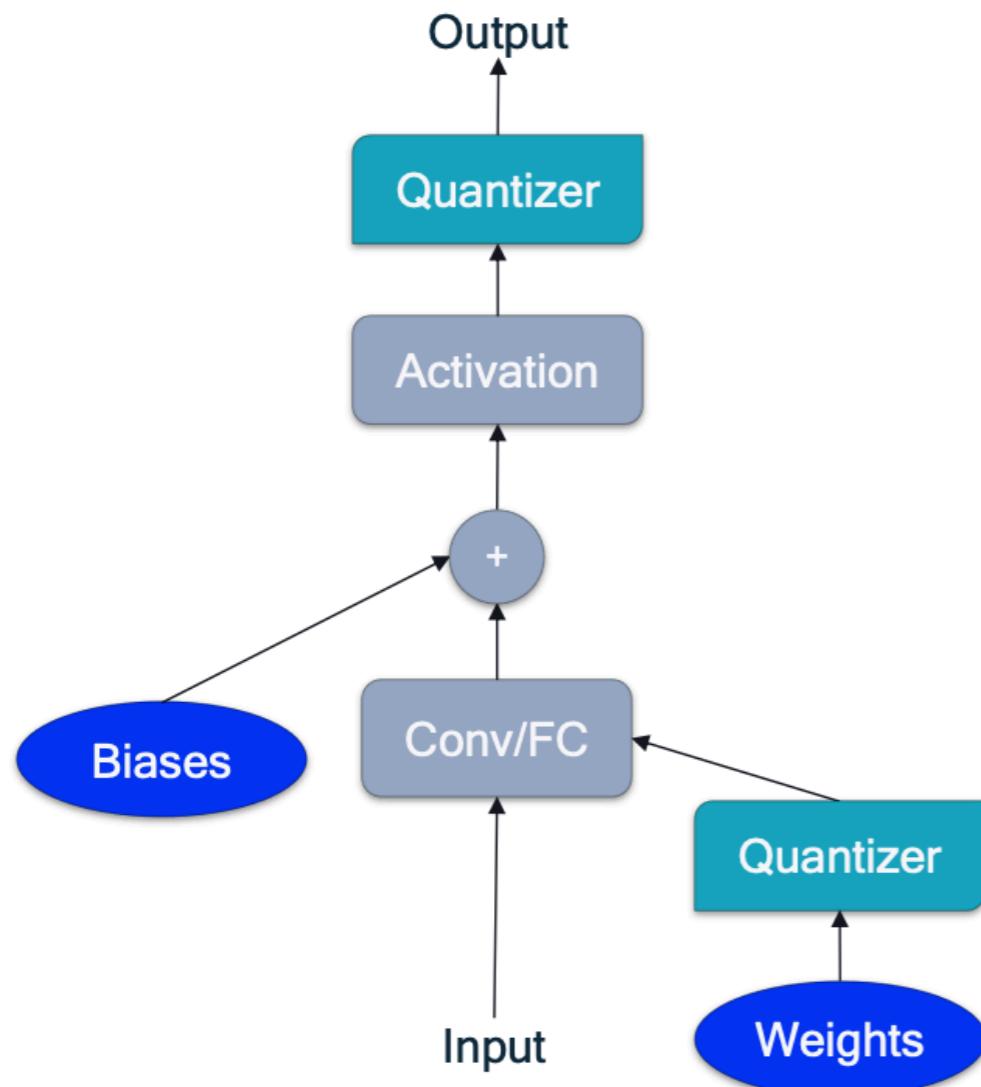
$$WX \approx s_W(W_{\text{int}} - z_W) s_X(X_{\text{int}} - z_X)$$
$$= s_W s_X (W_{\text{int}} X_{\text{int}}) + \underbrace{s_W s_X z_X W_{\text{int}} + s_W z_W s_X z_X}_{\text{Precompute, add to layer bias}} + \underbrace{s_W s_X z_W X_{\text{int}}}_{\text{Data-dependent overhead}}$$

Asymmetric weight quantization is equivalent to adding an input channel

Symmetric weights and asymmetric activations more hardware efficient

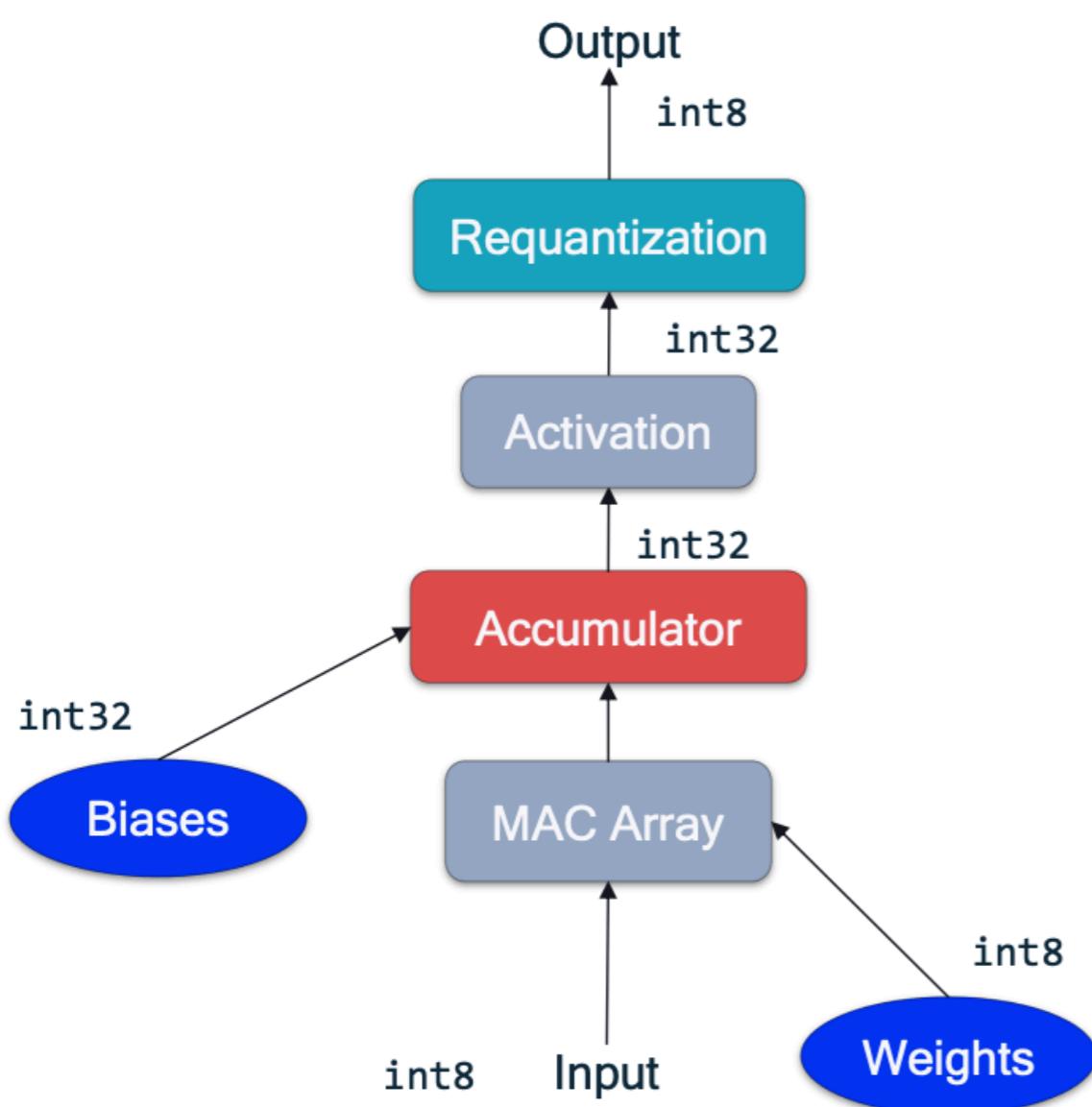
# Simulating Quantization

# Why simulate quantization?

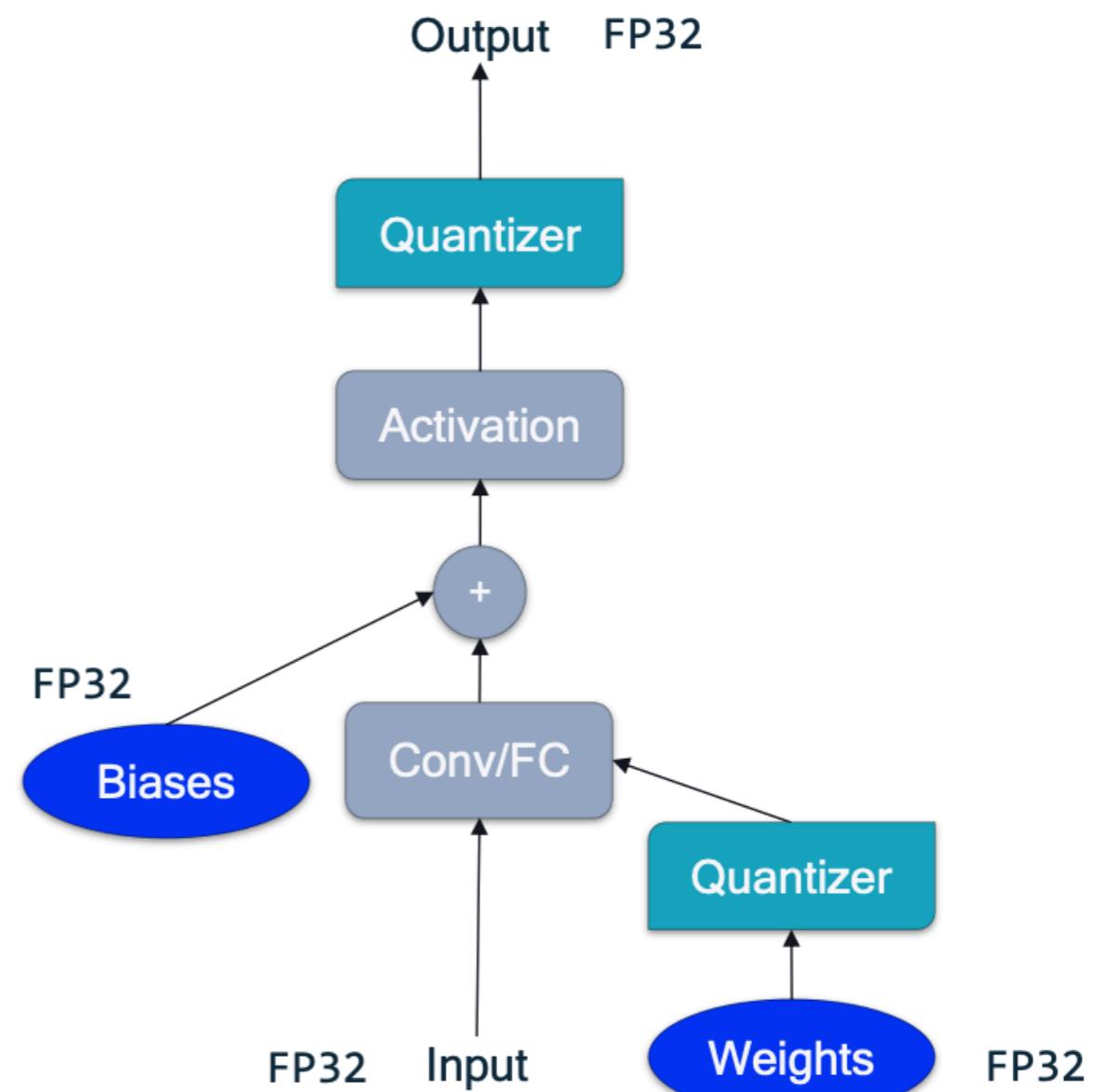


- We simulate fixed-point operations with floating-point numbers using general purpose hardware (e.g. CPU, GPU)
- This simulation is achieved by introducing simulated **quantization operations** (quantizers) to the compute graph.
- Quantization simulation benefits:
  - Enables GPUs acceleration
  - No need for dedicated kernels
  - Test various quantization option and bit-widths

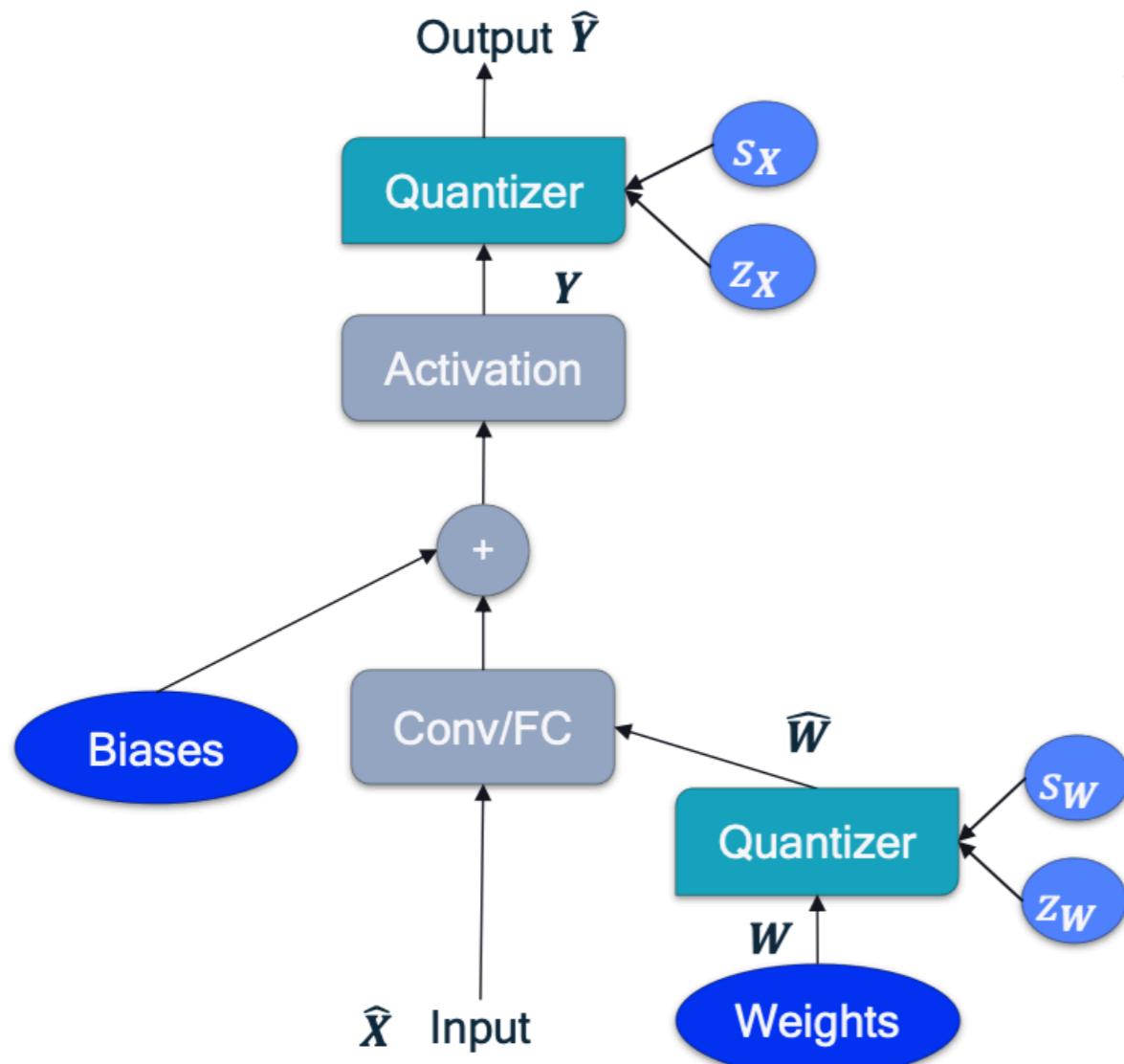
## On-device fixed-point inference



## Simulated quantized inference



# What operations do the quantizer perform?



Assuming asymmetric quantization the quantization operation applied to input tensor  $X$ :

$$X_{\text{int}} = \text{clip} \left( \text{round} \left( \frac{X}{s} \right) + z, \min = 0, \max = 2^b - 1 \right)$$

$$\hat{X} = s (X_{\text{int}} - z)$$

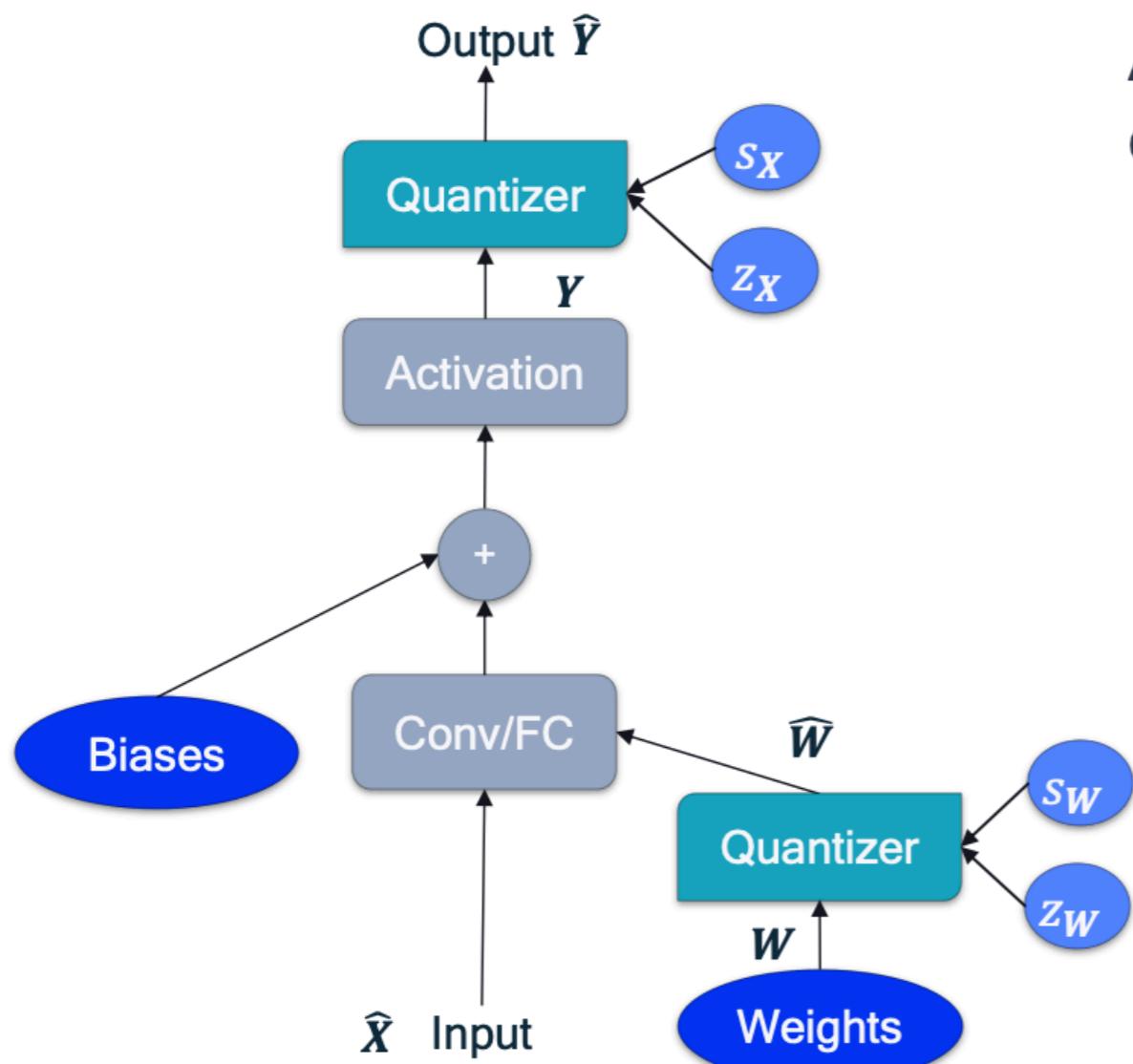
Example using  $b = 4$ :

$$X = \begin{pmatrix} 0.41 & 0.0 \\ 0.8 & -0.5 \end{pmatrix}$$

$$s = \frac{1}{15} = 0.067$$

$$z = \text{round} \left( \frac{0.5}{0.067} \right) = 8$$

# What operations do the quantizer perform?



Assuming asymmetric quantization the quantization operation applied to input tensor  $X$ :

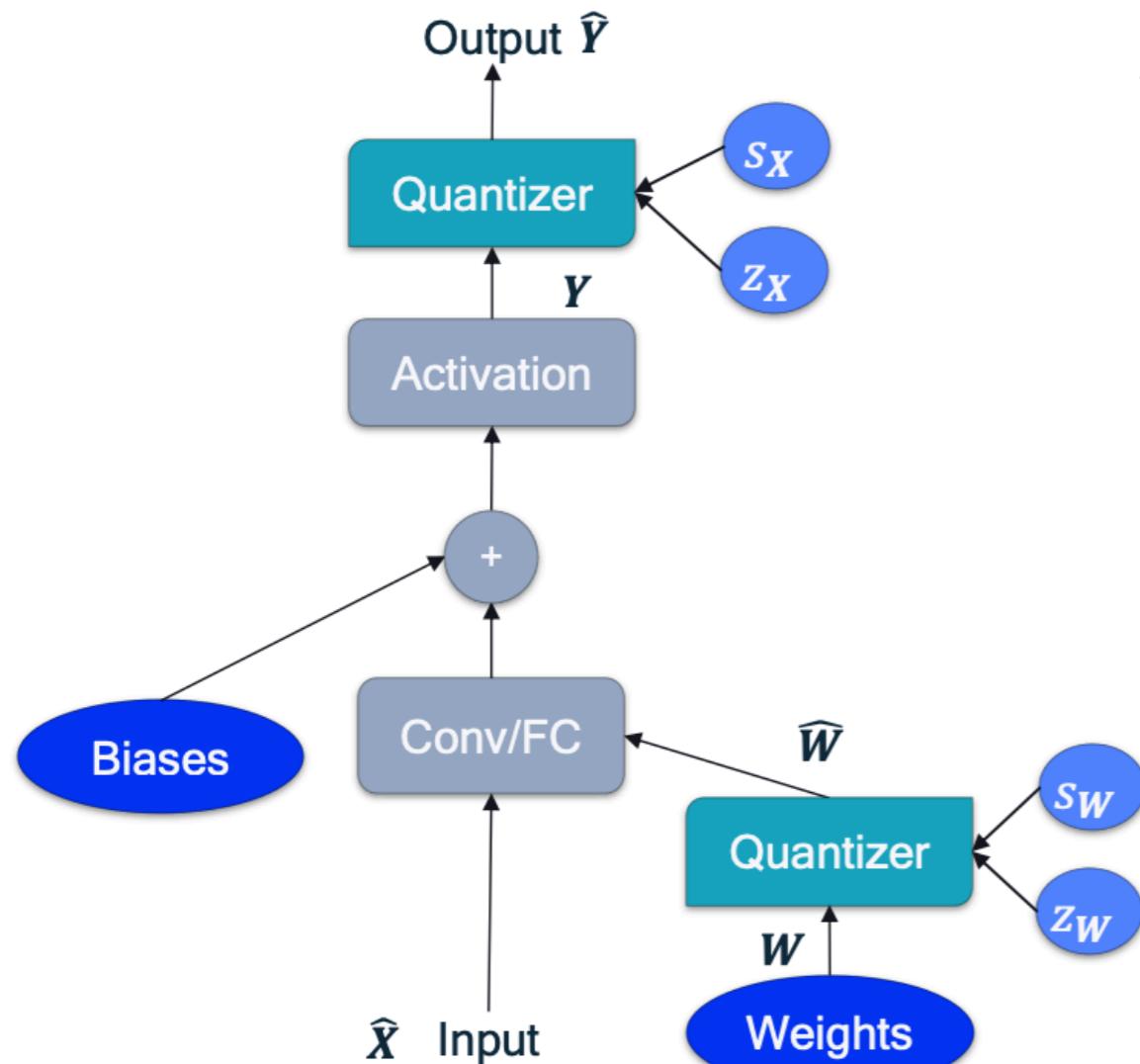
$$X_{\text{int}} = \text{clip} \left( \text{round} \left( \frac{X}{s} \right) + z, \min = 0, \max = 2^b - 1 \right)$$

$$\hat{X} = s (X_{\text{int}} - z) \quad \text{round} \left( \frac{X}{s} \right) + z = \begin{pmatrix} 14 & 8 \\ 20 & 0 \end{pmatrix}$$

Example using  $b = 4$ :  $s = 0.067$   $z = 8$

$$\frac{X}{s} = \begin{pmatrix} 6.15 & 0.0 \\ 12 & -7.5 \end{pmatrix}$$

# What operations do the quantizer perform?



Assuming asymmetric quantization the quantization operation applied to input tensor  $X$ :

$$X_{\text{int}} = \text{clip}\left(\text{round}\left(\frac{X}{s}\right) + z, \min = 0, \max = 2^b - 1\right)$$

$$\hat{X} = s(X_{\text{int}} - z) \quad \text{round}\left(\frac{X}{s}\right) + z = \begin{pmatrix} 14 & 8 \\ 20 & 0 \end{pmatrix}$$

Example using  $b = 4$ :  $s = 0.067$   $z = 8$

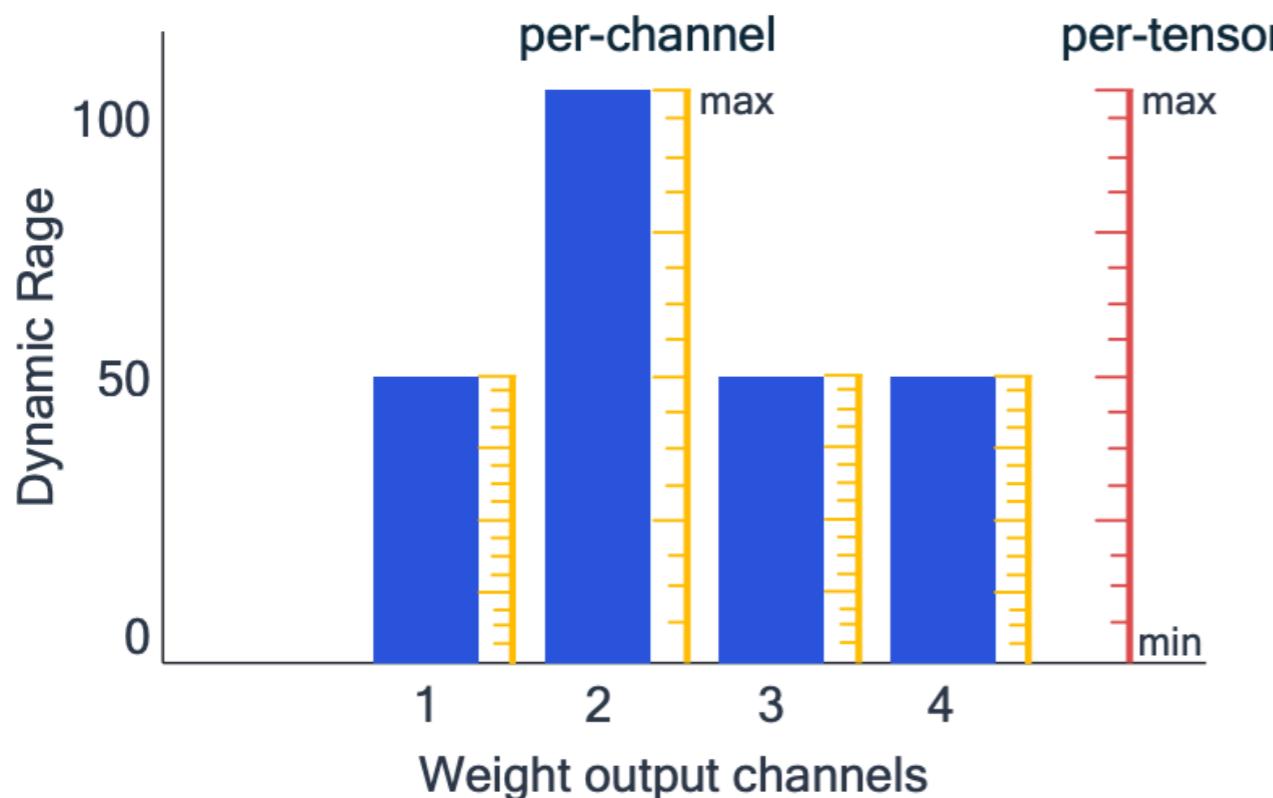
$$\text{round}\left(\frac{X}{s}\right) + z = \begin{pmatrix} 14 & 8 \\ 20 & 0 \end{pmatrix} \xrightarrow{\text{clip}} \begin{pmatrix} 14 & 8 \\ 15 & 0 \end{pmatrix}$$

de-quantize

$$X = \begin{pmatrix} 0.41 & 0.0 \\ 0.8 & -0.5 \end{pmatrix} \quad \hat{X} = \begin{pmatrix} 0.4 & 0.0 \\ 0.47 & -0.53 \end{pmatrix}$$

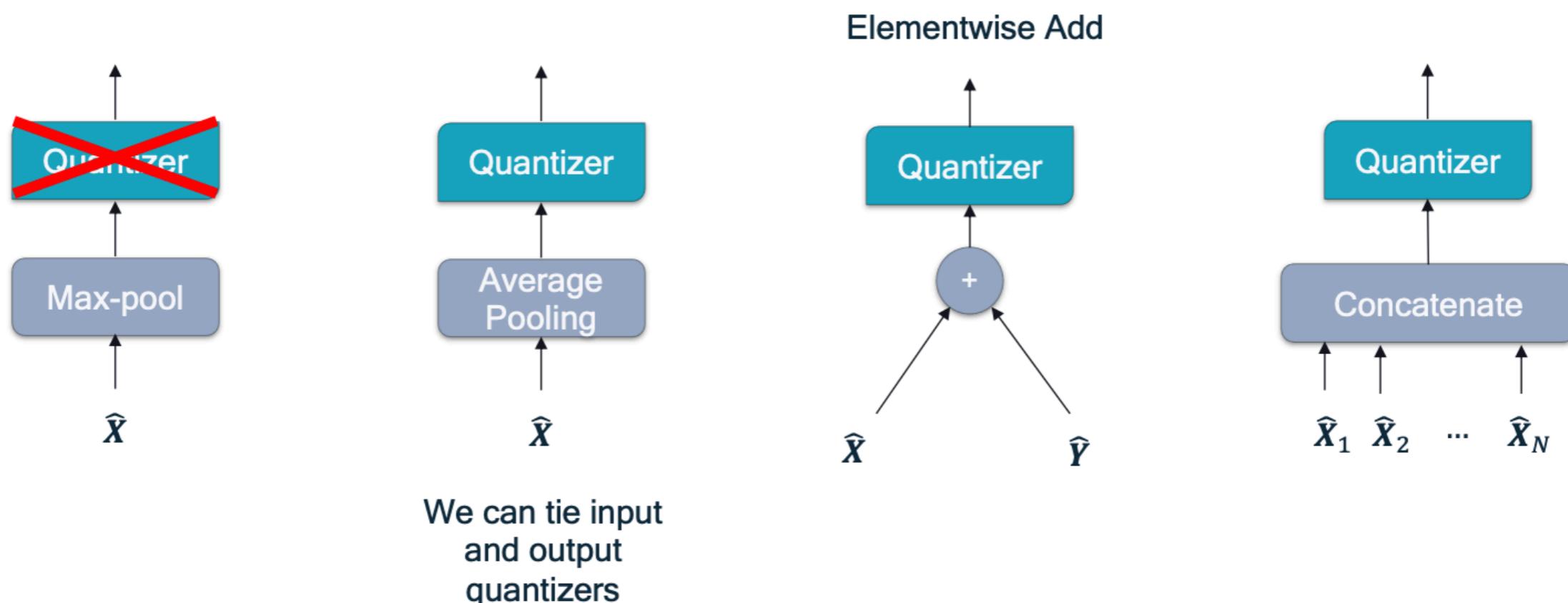
# Per-channel vs Per-tensor quantization of weights

Schematic histogram of weight ranges for layer



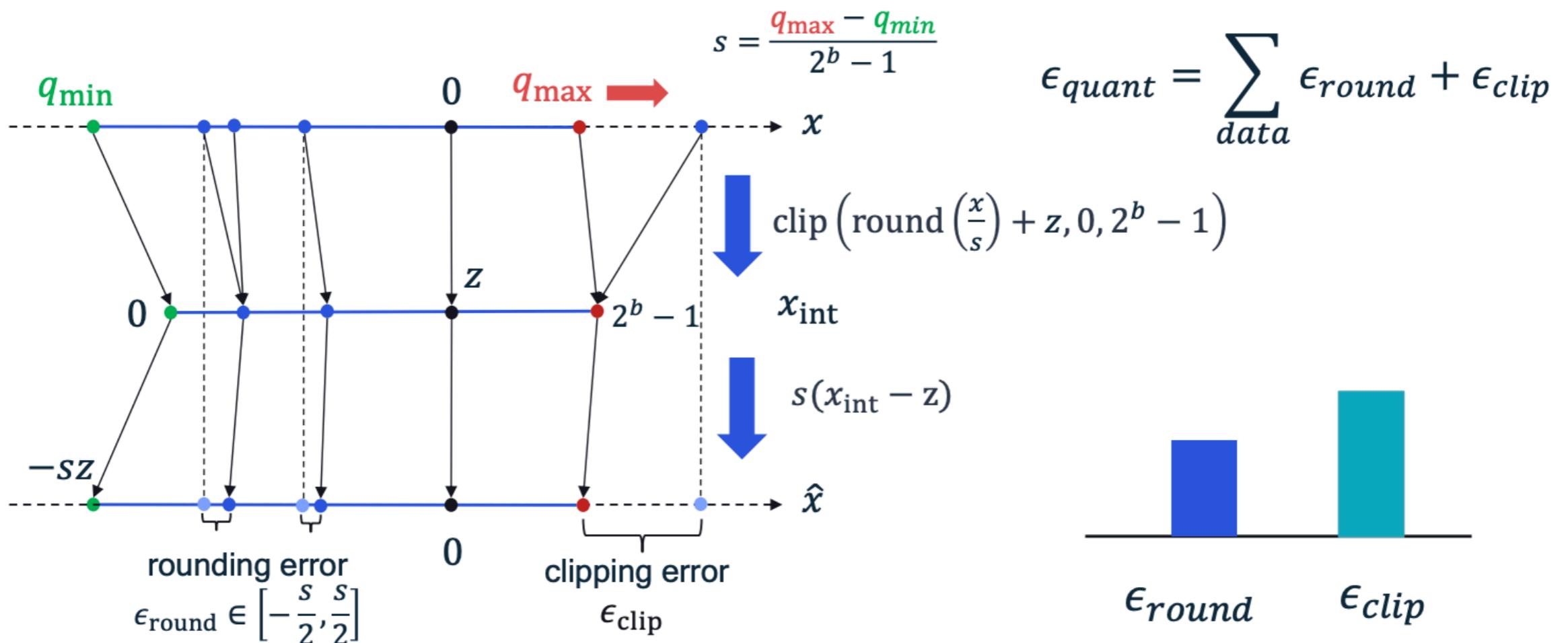
- **Per-tensor quantization** most supported by fixed-point accelerators
- **Per-channel quantization** better utilizes the quantization grid
- Per-channel quantization increasingly popular for weights
- Check for HW support

# How to simulate quantization in common DL layers

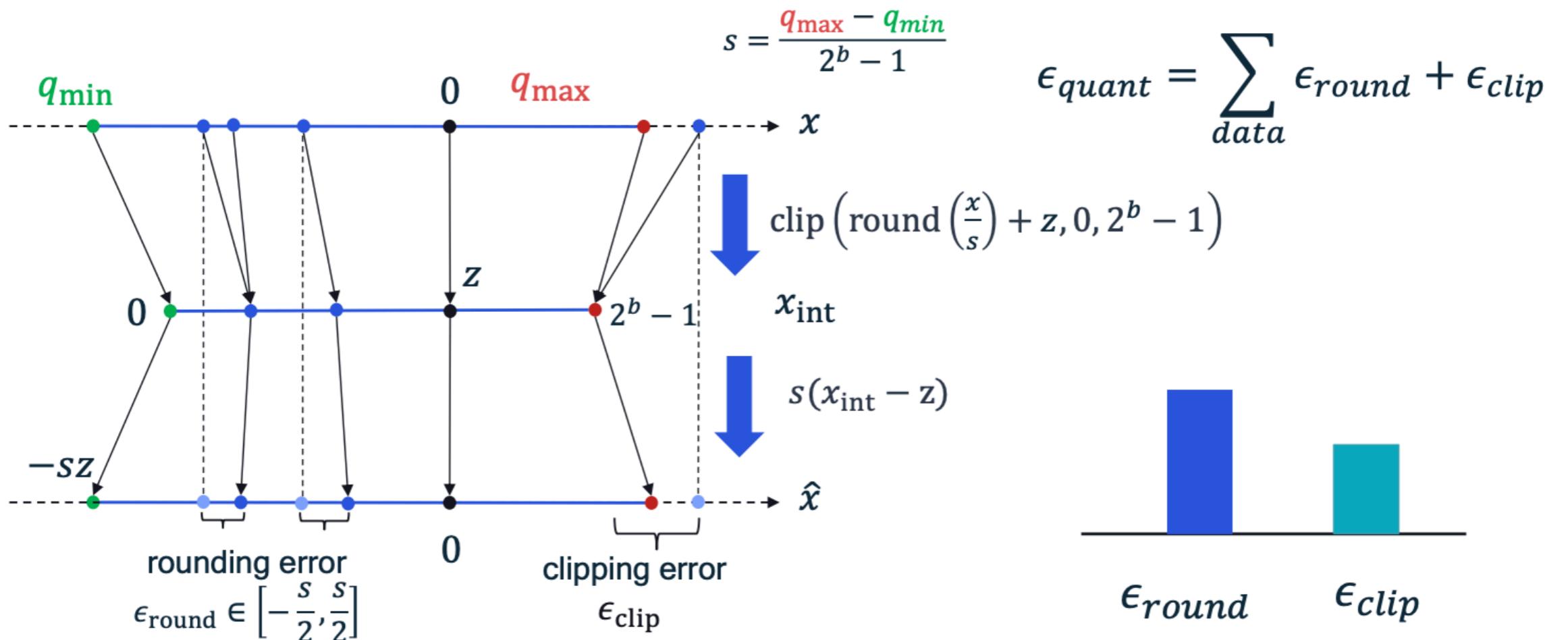


# Choosing the quantization parameters

# Sources of quantization error



# Sources of quantization error



# Quantization range setting methods

- Min-max range:

$$q_{\min} = \min \mathbf{X}$$

$$q_{\max} = \max \mathbf{X}$$

- Optimization-based methods:

$$\operatorname{argmin}_{q_{\min}, q_{\max}} \ell(\mathbf{X}, \hat{\mathbf{X}}(q_{\min}, q_{\max}))$$

MSE                          Cross-entropy

- Batch-Norm Based [1]:

$$q_{\min} = \min (\boldsymbol{\beta} - \alpha \boldsymbol{\gamma})$$

$$q_{\max} = \max (\boldsymbol{\beta} + \alpha \boldsymbol{\gamma})$$

$$\begin{aligned} & \text{BatchNorm}(\mathbf{z}_k) \\ &= \boldsymbol{\gamma}_k \frac{\mathbf{z}_k - \boldsymbol{\mu}_k}{\sqrt{\boldsymbol{\sigma}_k + \epsilon}} + \boldsymbol{\beta}_k \end{aligned}$$

[1] Nagel et al, 2019, Data-Free Quantization Through Weight Equalization and Bias Correction

# Quantization setting methods ablation study

Model (FP32 Accuracy)	ResNet18 (69.68)		MobileNetV2 (71.72)		
	Bit-width	A8	A8	A6	
Min-Max		69.60	68.19	70.96	64.58
MSE		69.59	67.84	71.35	67.55
MSE & X-entropy		69.60	68.91	71.36	68.85
BN ( $\alpha = 6$ )		69.54	68.73	71.32	71.32

Average ImageNet validation accuracy (%) over 5 seeds

## Post-Training Quantization (PTQ)

- ✓ Takes a pre-trained network and converts it to a fixed-point network without access to the training pipeline
- ✓ Data-free or small calibration set needed
- ✓ Use though single API call
- ✗ Lower accuracy at lower bit-widths

## Quantization-Aware Training (QAT)

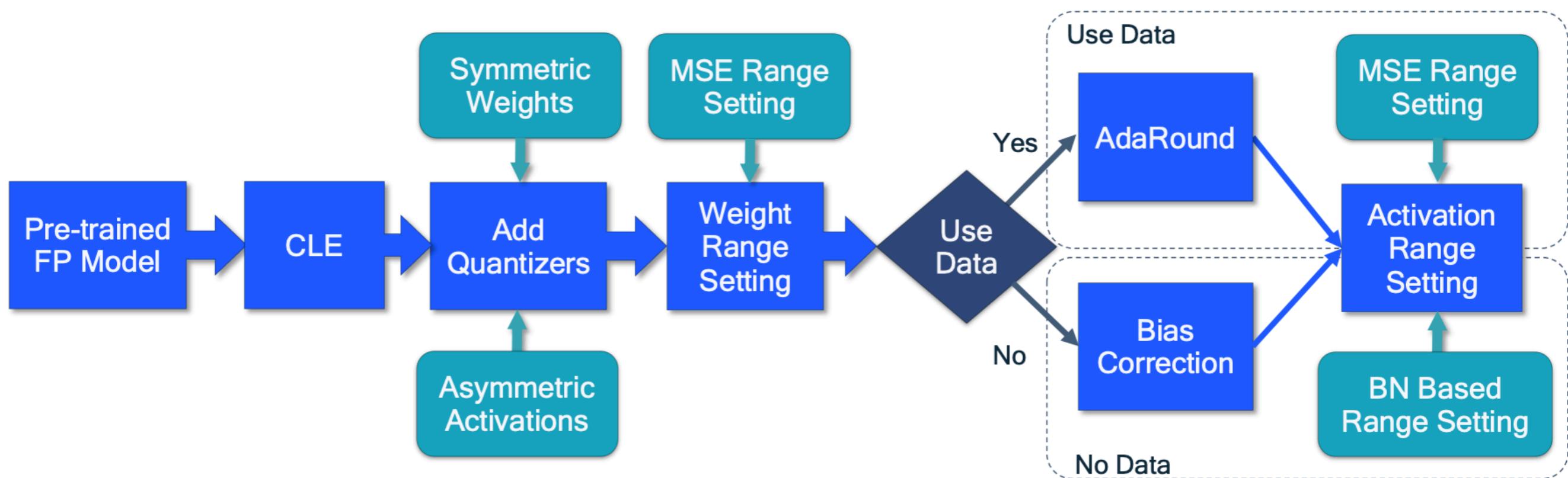
- ✗ Requires access to training pipeline and labelled data
- ✗ Longer training times
- ✗ Hyper-parameter tuning
- ✓ Achieves higher accuracy

Source sample text

What algorithm to choose to improve accuracy?

# Post training quantization

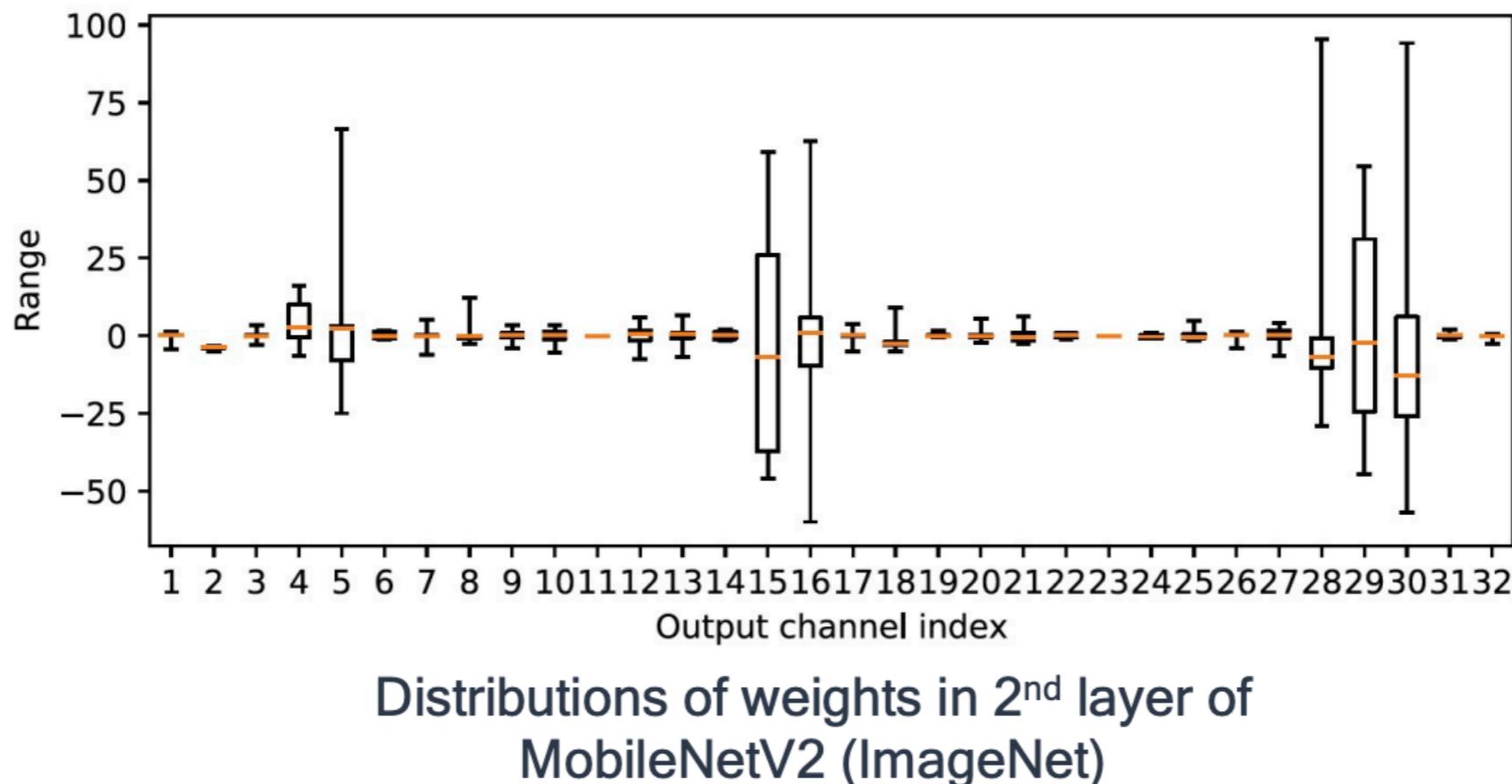
# Post-training quantization pipeline



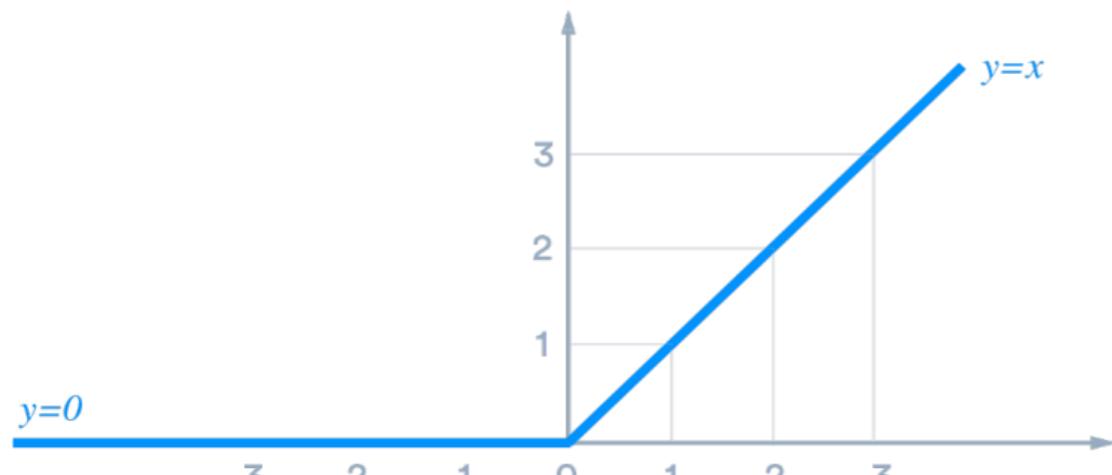


## Cross-Layer Equalization

# Imbalanced weights is a common problem in practice



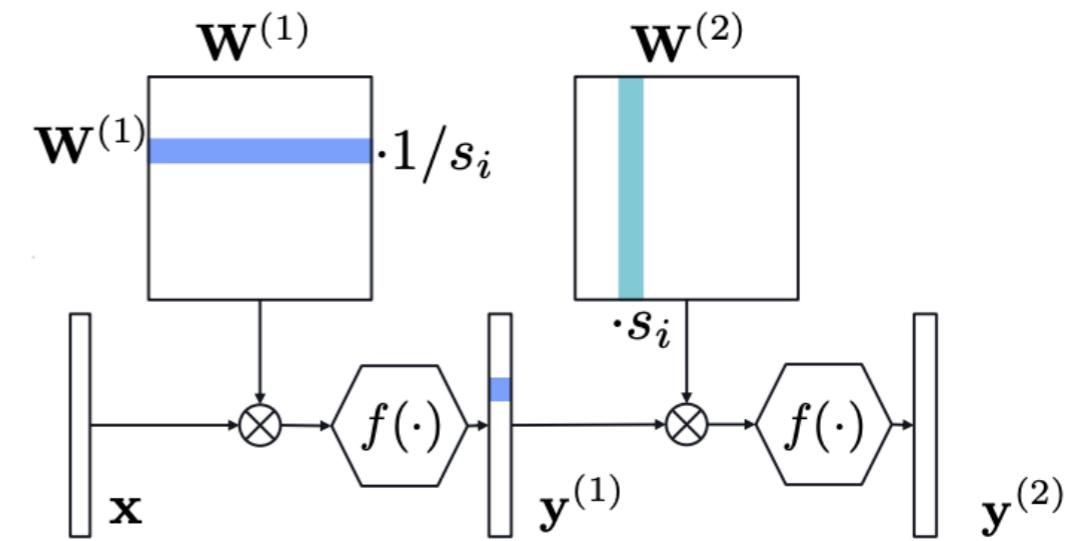
# Cross-layer equalization scales weights in neighboring layers for better quantization



$$\text{ReLU}(x) = \max(0, x)$$

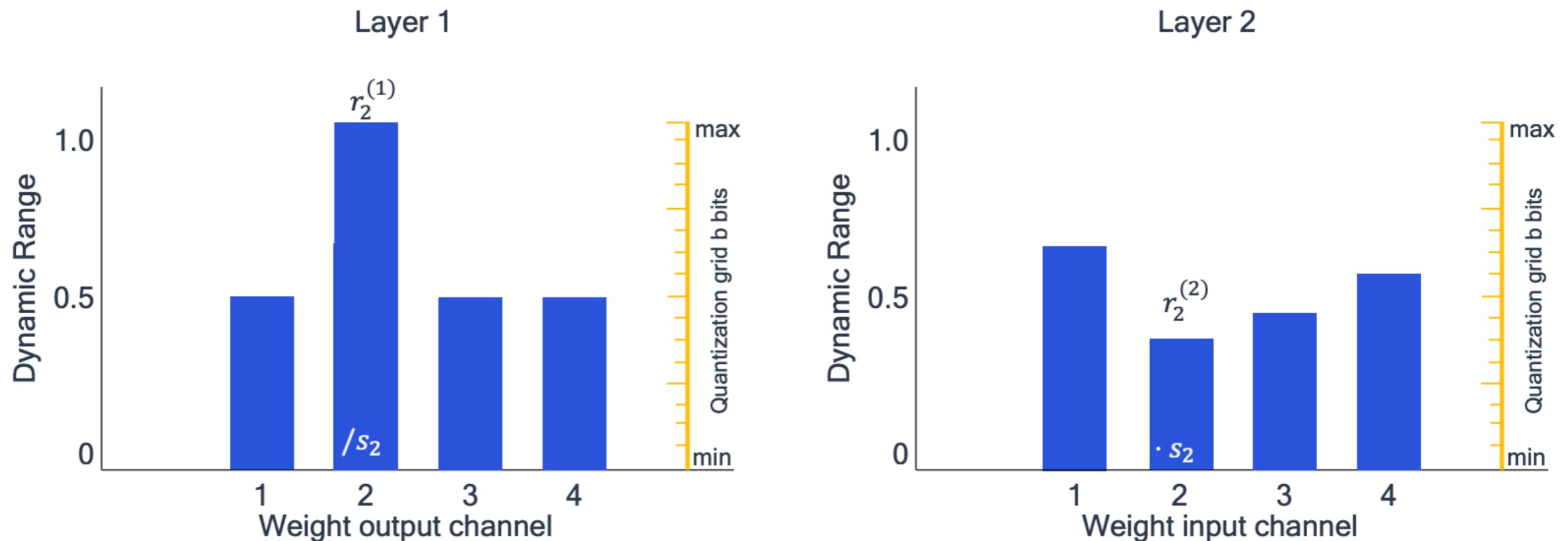
ReLU is scale-equivariant

$$\text{ReLU}(\mathbf{s}x) = \mathbf{s} \cdot \text{ReLU}(x)$$



We can scale two neighboring layers together to optimize it for quantization

# Finding the scaling factors for cross-layer equalization



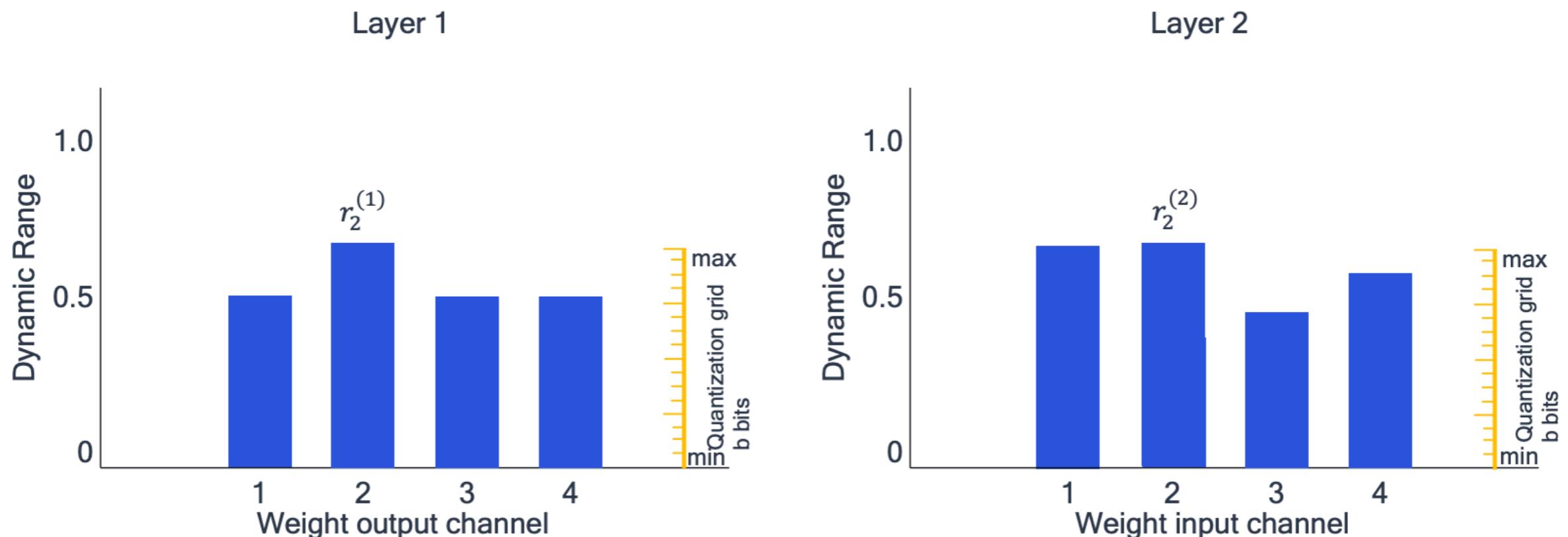
Equalize the weight channels of layer 1 with weight channel of layer 2

$$\text{by setting } s_i = \frac{1}{r_i^{(2)}} \sqrt{r_i^{(1)} r_i^{(2)}}$$

Proved to be optimal

$r_i$  :Range of layer  $i$

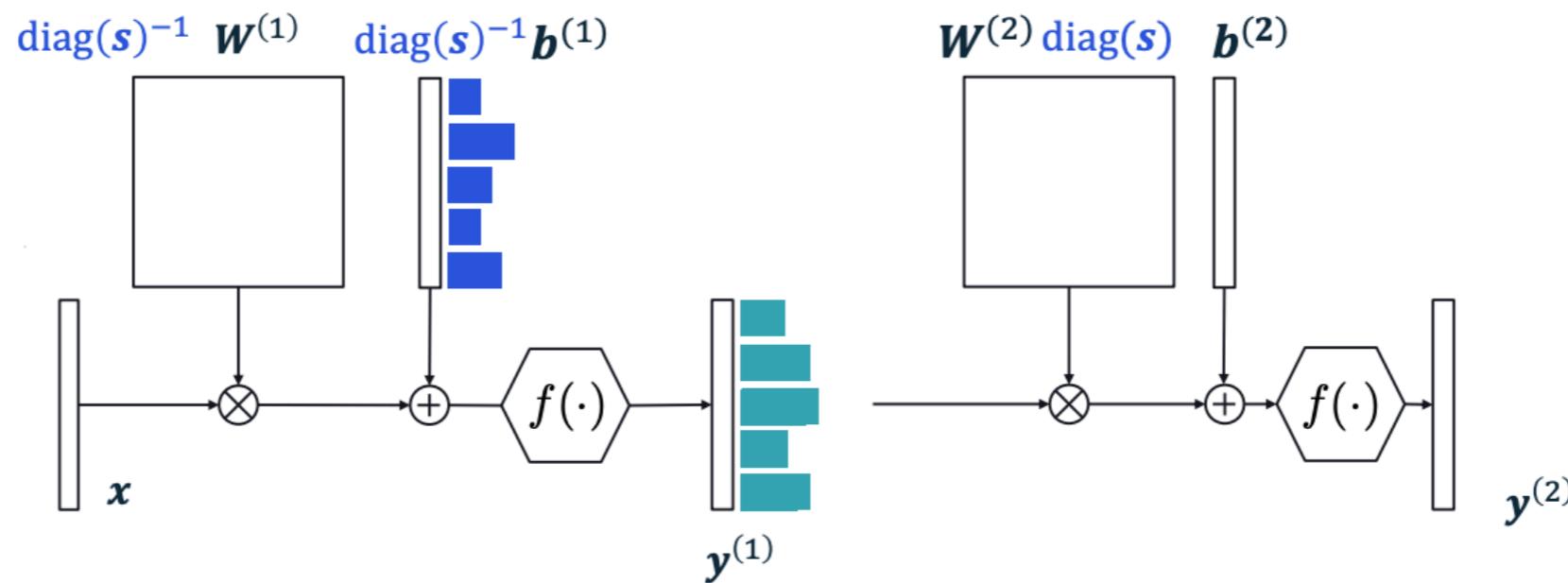
# Finding the scaling factors for cross-layer equalization



Equalize the weight channels of layer 1 with weight channel of layer 2

$$\text{by setting } s_i = \frac{1}{r_i^{(2)}} \sqrt{r_i^{(1)} r_i^{(2)}}$$

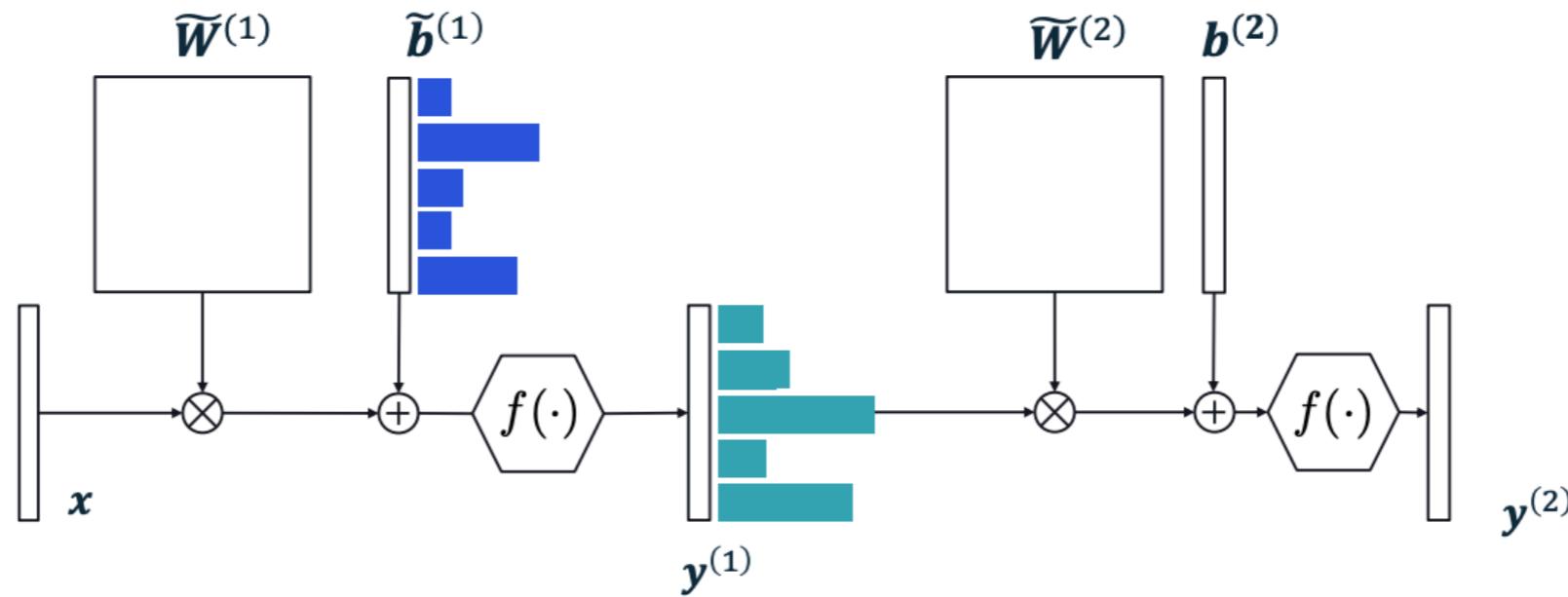
# Absorbing large biases to the next layer equalizes activation ranges



Source sample text

Equalize activation ranges by absorbing  $c$  from layer 1 into layer 2

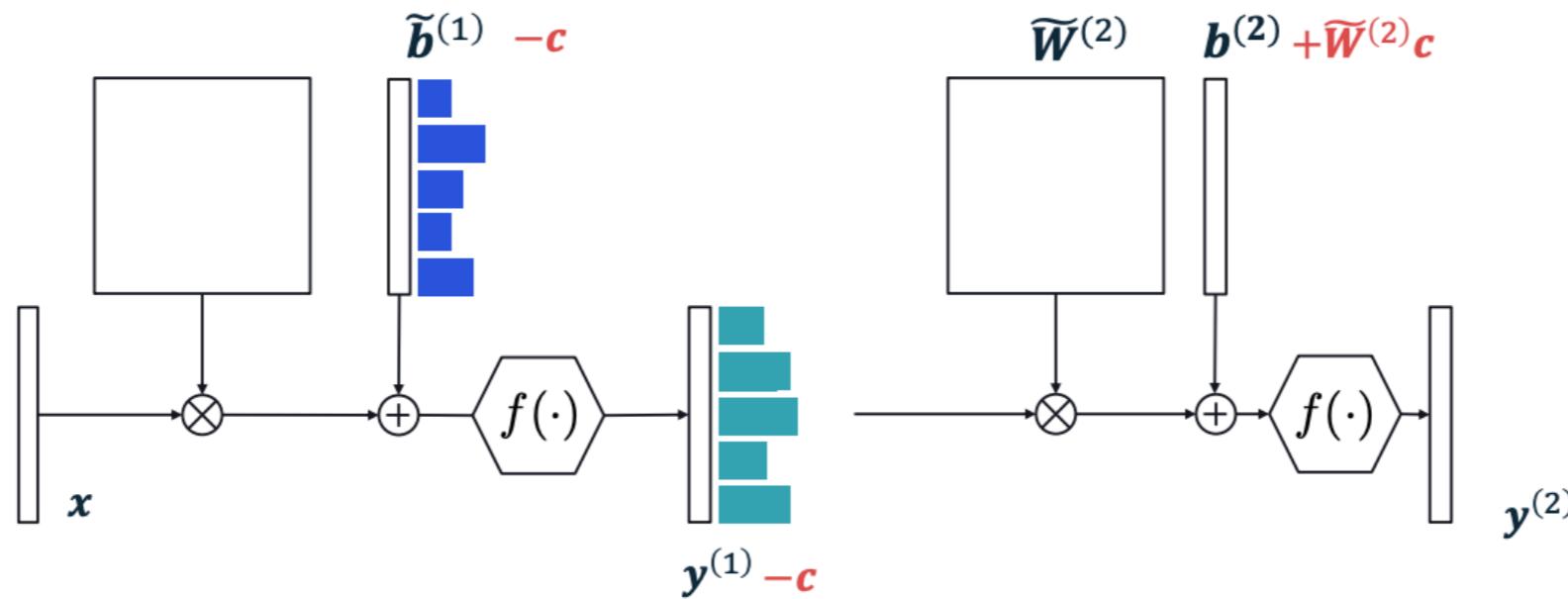
# Absorbing large biases to the next layer equalizes activation ranges



When  $s_i < 1$ , range of biases may increase

Equalize activation ranges by absorbing  $c$  from layer 1 into layer 2

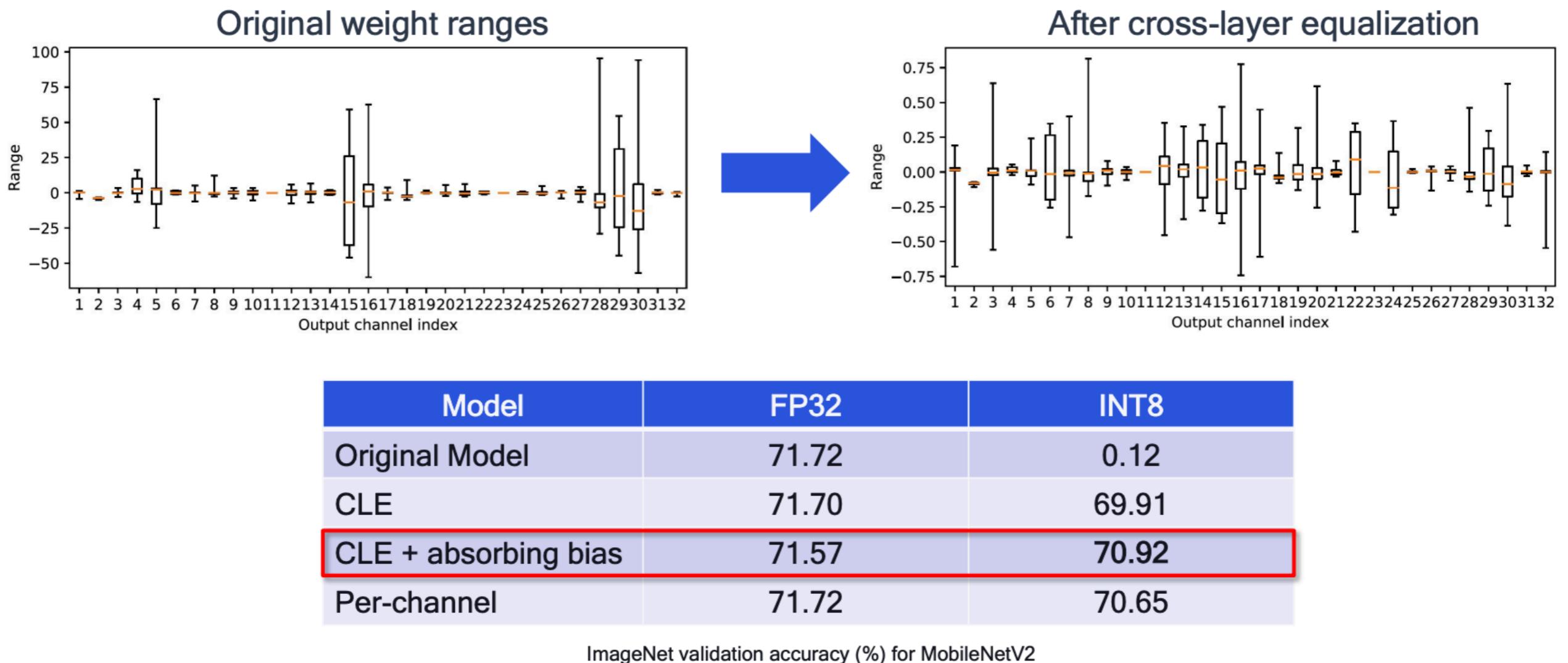
# Absorbing large biases to the next layer equalizes activation ranges



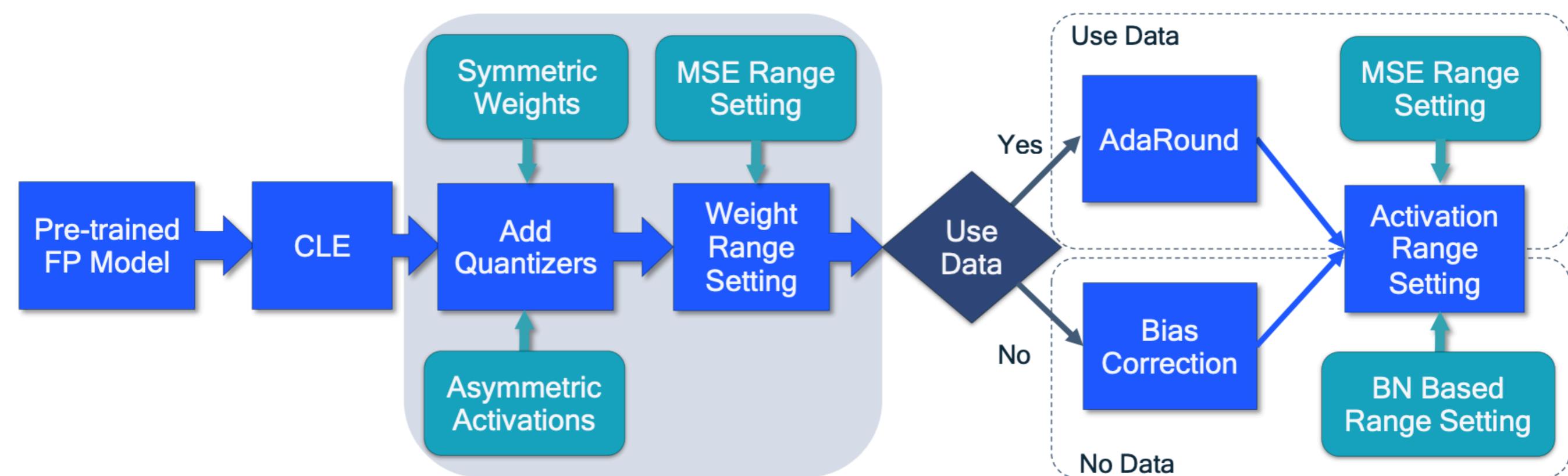
When  $s_i < 1$ , range of biases may increase

Equalize activation ranges by absorbing  $c$  from layer 1 into layer 2

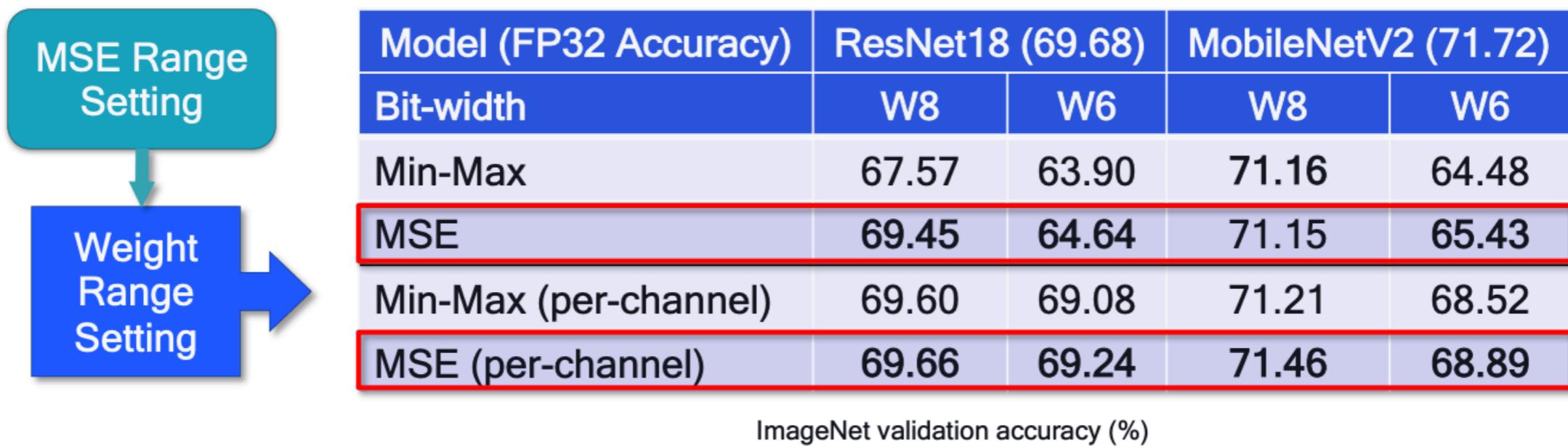
# Cross-layer equalization significantly improves accuracy



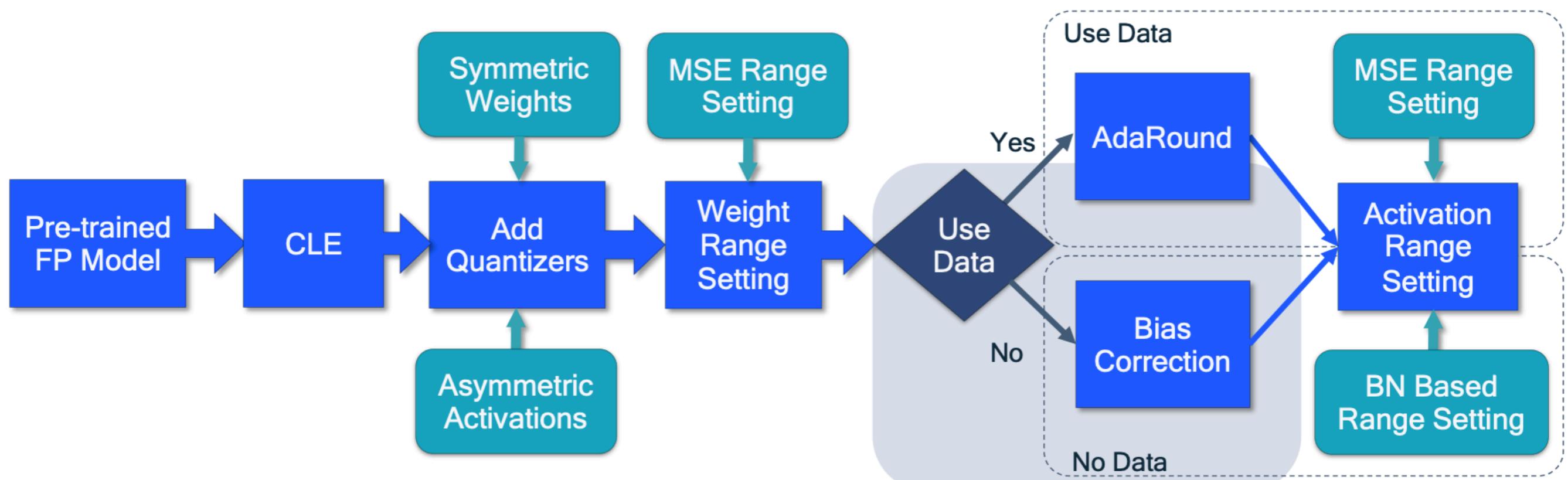
# Quantizer and range setting



# Quantizer and range setting



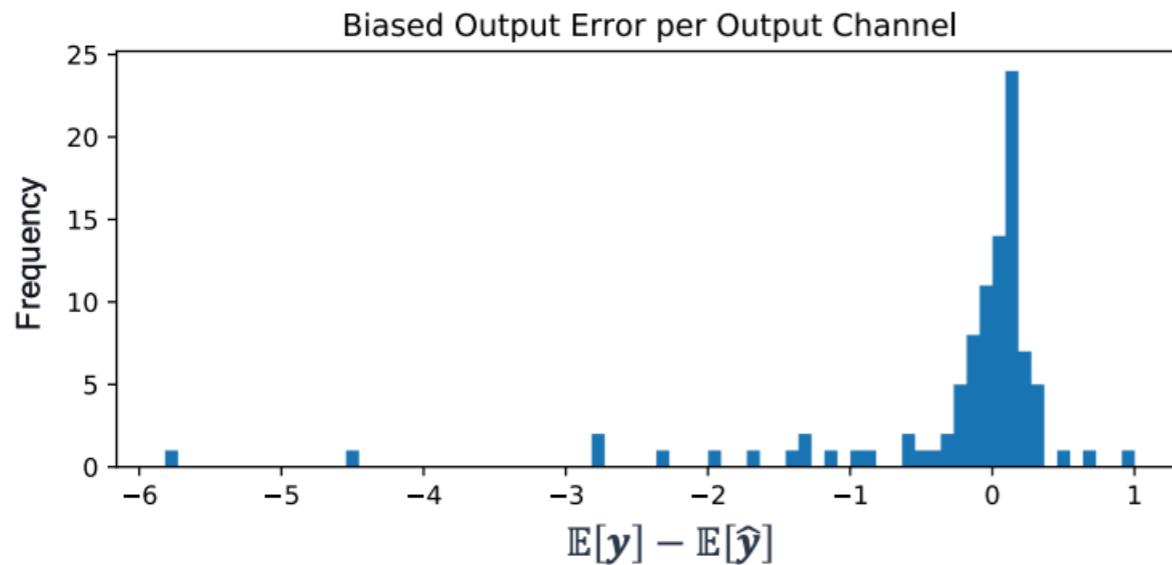
# Bias Correction



# Biased quantization error leads to accuracy drop

Quantization error is often biased

$$\begin{aligned}\mathbb{E}[y] - \mathbb{E}[\hat{y}] &= \mathbb{E}[Wx] - \mathbb{E}[\hat{W}x] \\ &= W\mathbb{E}[x] - \hat{W}\mathbb{E}[x] \\ &= \Delta W \mathbb{E}[x]\end{aligned}\quad \text{Empirical bias correction}$$



Per-channel biased output error introduced by weight quantization of the second depth-wise separable layer in MobileNetV2

**Key idea: Bias correction**



$$\begin{aligned}\mathbb{E}[x] &= \mathbb{E}[\text{ReLU}(x^{\text{pre}})] \\ &= \gamma \mathcal{N}\left(\frac{-\beta}{\gamma}\right) + \beta \left[1 - \Phi\left(\frac{-\beta}{\gamma}\right)\right]\end{aligned}$$

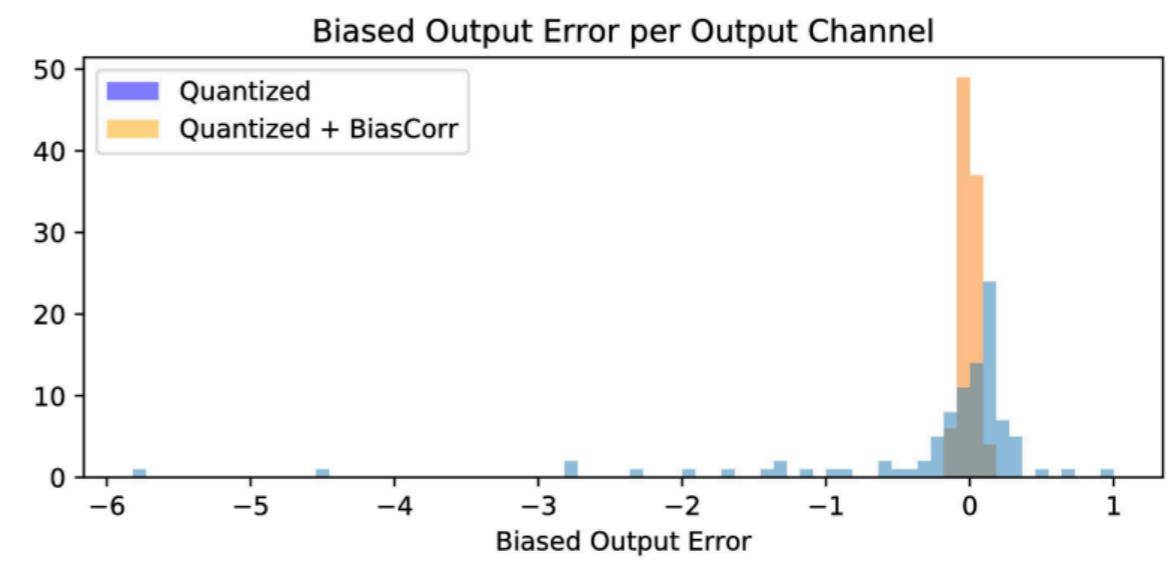
Analytic bias correction

$\phi$ :normal cdf  
 $\beta, \gamma$ :batch norm parameters

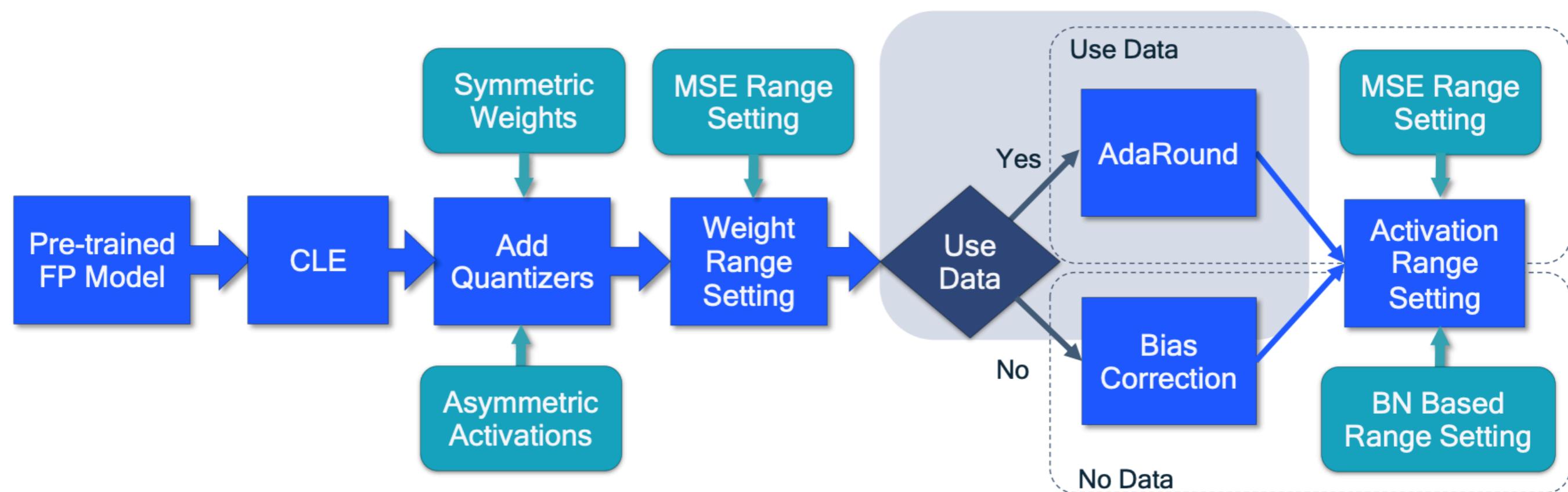
# Bias correction

Model	W8A8	FP32
Original Model	0.12	71.72
+bias correction	52.02	71.72
CLE + bias absorption	70.92	71.57
+bias correction	71.79	71.57

ImageNet val. accuracy for MobileNetV2



# AdaRound



# AdaRound

- Traditionally, in PTQ we use rounding-to-nearest operator

$$x_{\text{int}} = \text{clip} \left( \text{round} \left( \frac{x}{s} \right) + z, \min = 0, \max = 2^b - 1 \right)$$

- However, rounding-to-nearest is not optimal?

Rounding Method	Accuracy (%)
Nearest	52.29
Floor / Ceil	00.10
Stochastic	$52.06^{\pm 5.52}$
Stochastic (best)	63.06

4-bit weight quantization of 1<sup>st</sup> layer of Resnet18,  
validation accuracy on ImageNet.

Up or Down? Adaptive Rounding for Post-Training Quantization (Nagel, Amjad, et al., ICML 2020)

## Up or Down?

How can we systematically find the best rounding choice?

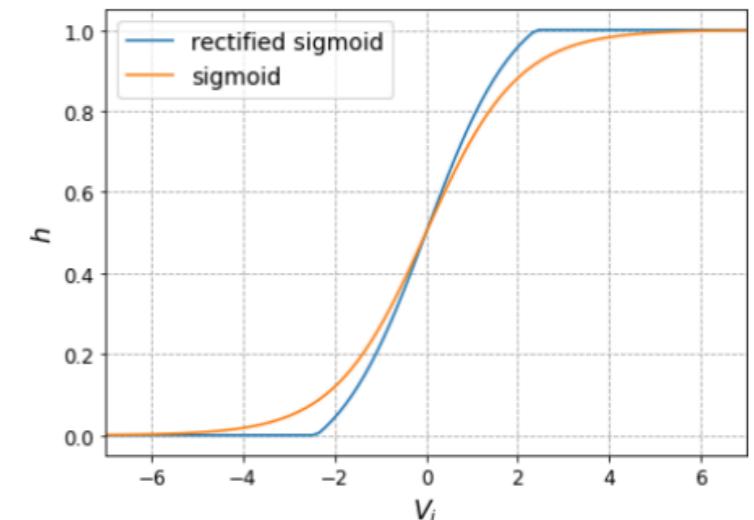
# AdaRound: learning to round

## AdaRound: learning to round

- Minimize local  $L_2$  loss per-layer rather than task loss:

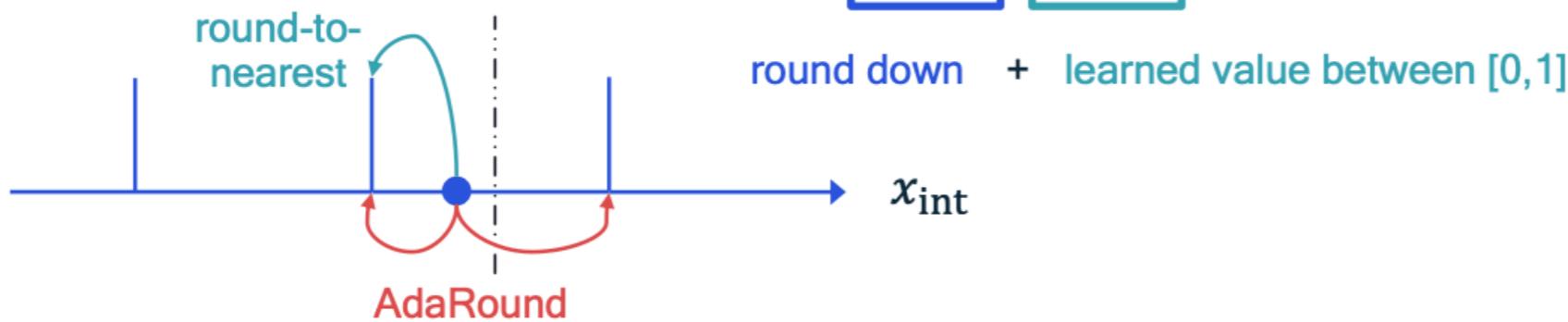
Surrogate for task loss

$$\arg \min_{\mathbf{V}} \|\mathbf{Wx} - \widetilde{\mathbf{W}}\mathbf{x}\|_F^2$$



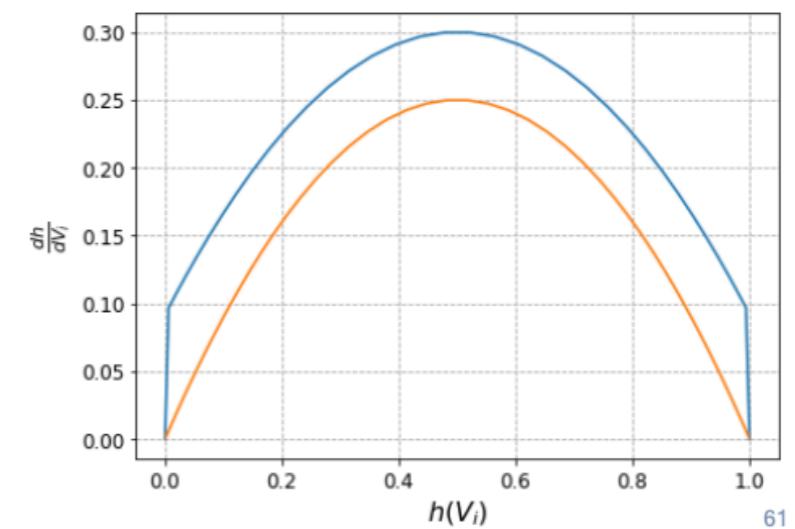
- where  $\widetilde{\mathbf{W}}$  are soft-quantized weights:

$$\widetilde{\mathbf{W}} = s \cdot \text{clip} \left( \left\lfloor \frac{\mathbf{W}}{s} \right\rfloor + h(\mathbf{V}), n, p \right)$$



$$h(\mathbf{V}) = \text{clip} (\sigma(\mathbf{V})(\zeta - \gamma) + \gamma, 0, 1)$$

rectified sigmoid



# AdaRound: learning to round

- Minimize local  $L_2$  loss per-layer rather than task loss:

$$\arg \min_{\mathbf{V}} \left\| \mathbf{Wx} - \widetilde{\mathbf{W}}\mathbf{x} \right\|_F^2 + \lambda f_{reg}(\mathbf{V})$$

regularizer forces  $h(\mathbf{V})$  to be 0 or 1

- where  $\widetilde{\mathbf{W}}$  are soft-quantized weights:

$$\widetilde{\mathbf{W}} = s \cdot clip \left( \left\lfloor \frac{\mathbf{W}}{s} \right\rfloor + h(\mathbf{V}), n, p \right)$$

$h(\mathbf{V}) = \text{clip} (\sigma(\mathbf{V})(\zeta - \gamma) + \gamma, 0, 1)$   
rectified sigmoid

round down + learned value between [0, 1]

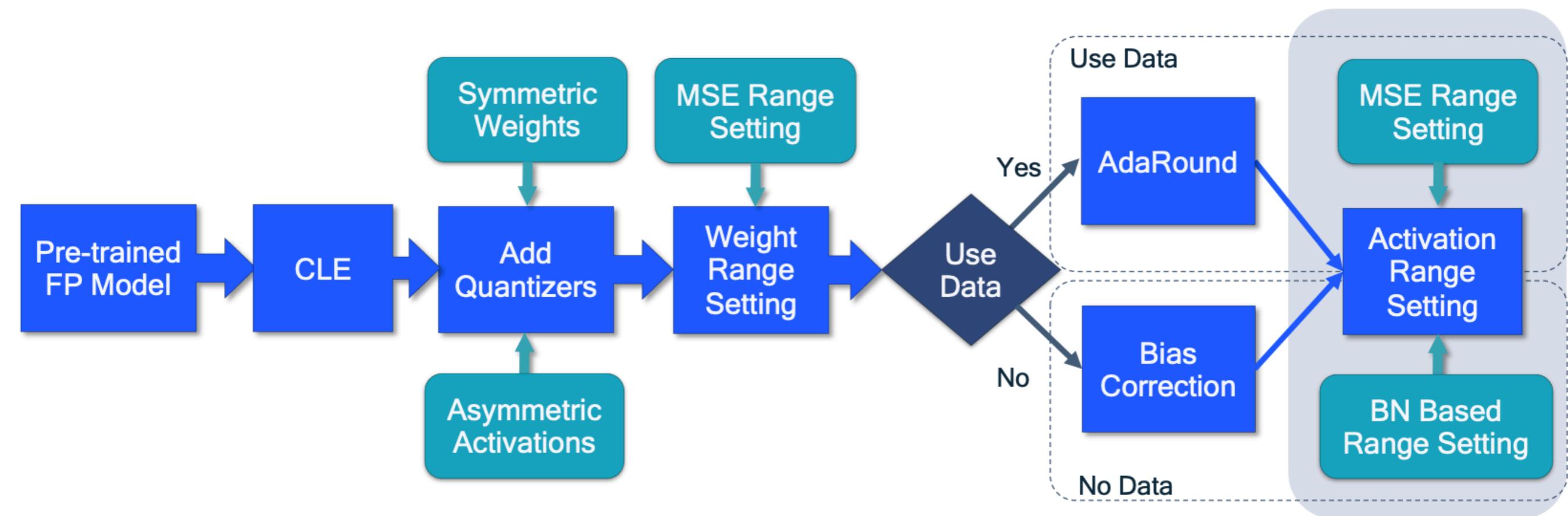
- Regularization:

$$f_{reg}(\mathbf{V}) = \sum_{i,j} 1 - |2h(\mathbf{V}_{i,j}) - 1|^\beta$$

# AdaRound results

Quantization method	#bits W/A	ResNet18	ResNet50	InceptionV3	MobileNetV2
Full precision	32/32	69.68	76.07	77.40	71.72
CLE + BC	4/8	38.98	52.84	-	46.67
Per channel bias corr*	4*/8	67.4	74.8	59.5	-
AdaRound	4/8	68.55	75.01	75.72	69.25

# Activation range setting



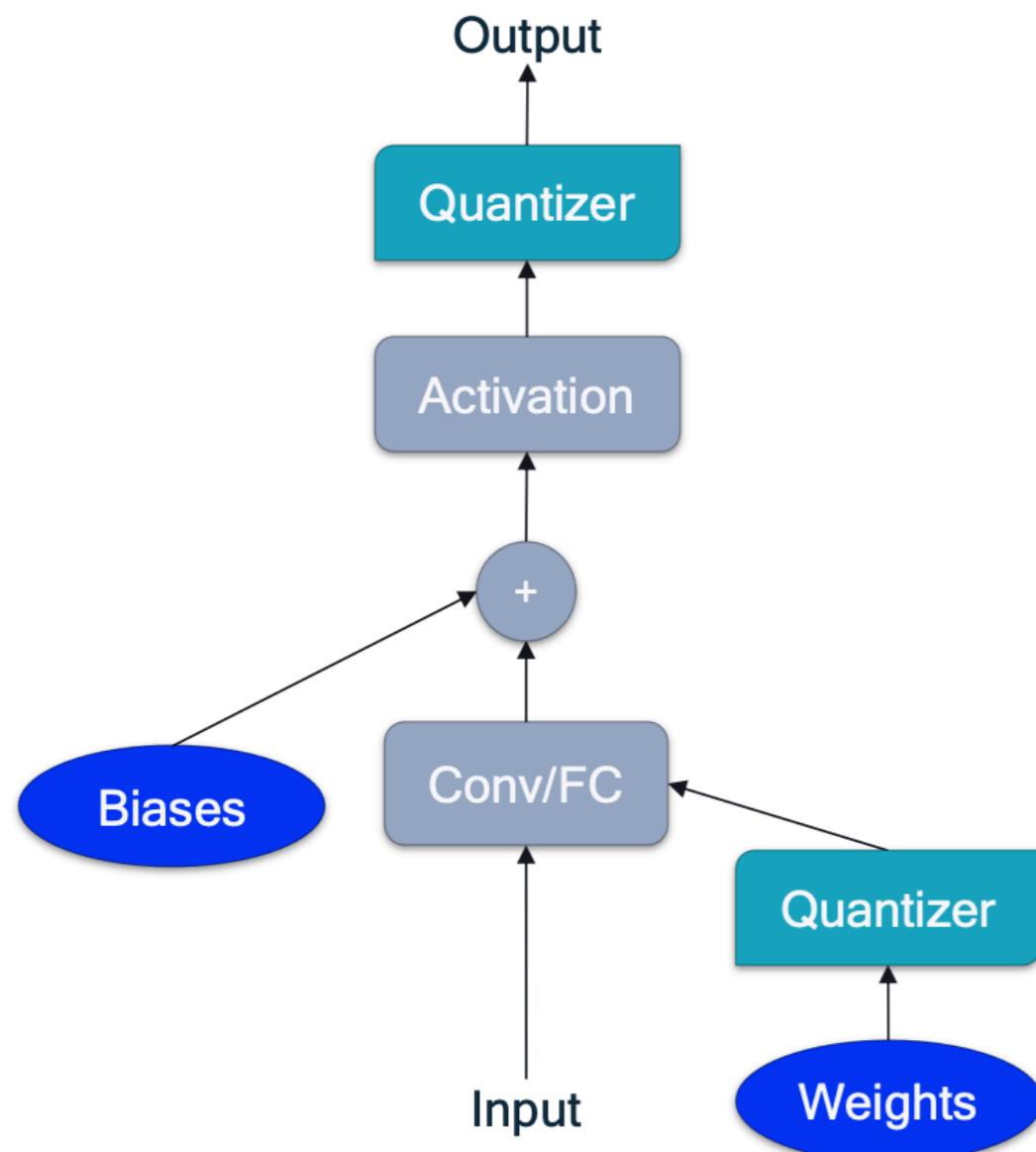
# PTQ results using previous pipeline

- drop  $\leq 1.0\%$
- $1.0\% < \text{drop} \leq 1.5\%$
- drop  $> 1.5\%$

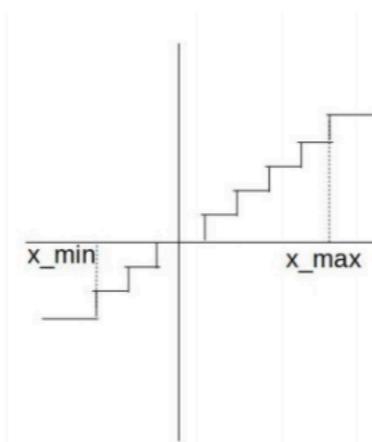
Models	FP32	Per-tensor				Per-channel			
		W8A8	diff	W4A8	diff	W8A8	diff	W4A8	diff
ResNet18	69.68	69.60	-0.08	68.62	-1.06	69.56	-0.12	68.91	-0.77
ResNet50	76.07	75.87	-0.20	75.15	-0.92	75.88	-0.19	75.43	-0.64
MobileNetV2	71.72	70.99	-0.73	69.21	-2.51	71.16	-0.56	69.79	-1.93
InceptionV3	77.40	77.68	+0.28	76.48	-0.92	77.71	-0.31	76.82	-0.58
EfficientNet lite	75.42	75.25	-0.17	71.24	-4.18	75.39	-0.03	74.01	-1.41
DeepLabV3	72.94	72.44	-0.50	70.80	-2.14	72.27	-0.67	71.67	-1.27
EfficientDet-D1	40.08	38.29	-1.79	0.31	-39.77	38.67	-1.41	35.08	-5.00
BERT-base	83.06	82.43	-0.63	81.76	-1.30	82.77	-0.29	82.02	-1.04

# Quantization aware training

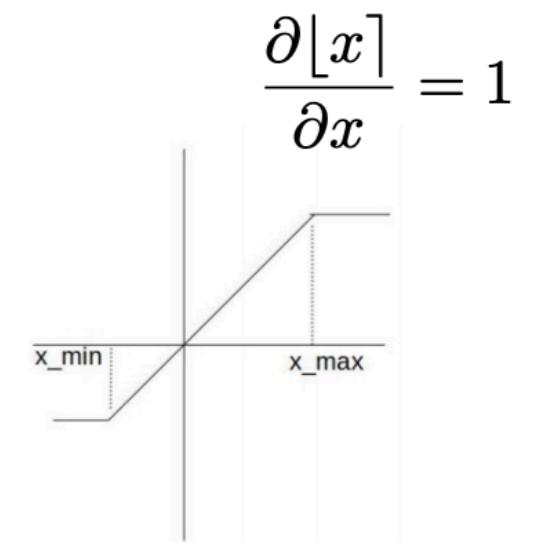
# Simulating quantization for backward path



- The round-to-nearest operation does not have meaningful gradients
- Gradient-based training impossible
- Solution: Redefine gradient with the “straight-through estimator” (STE)\*



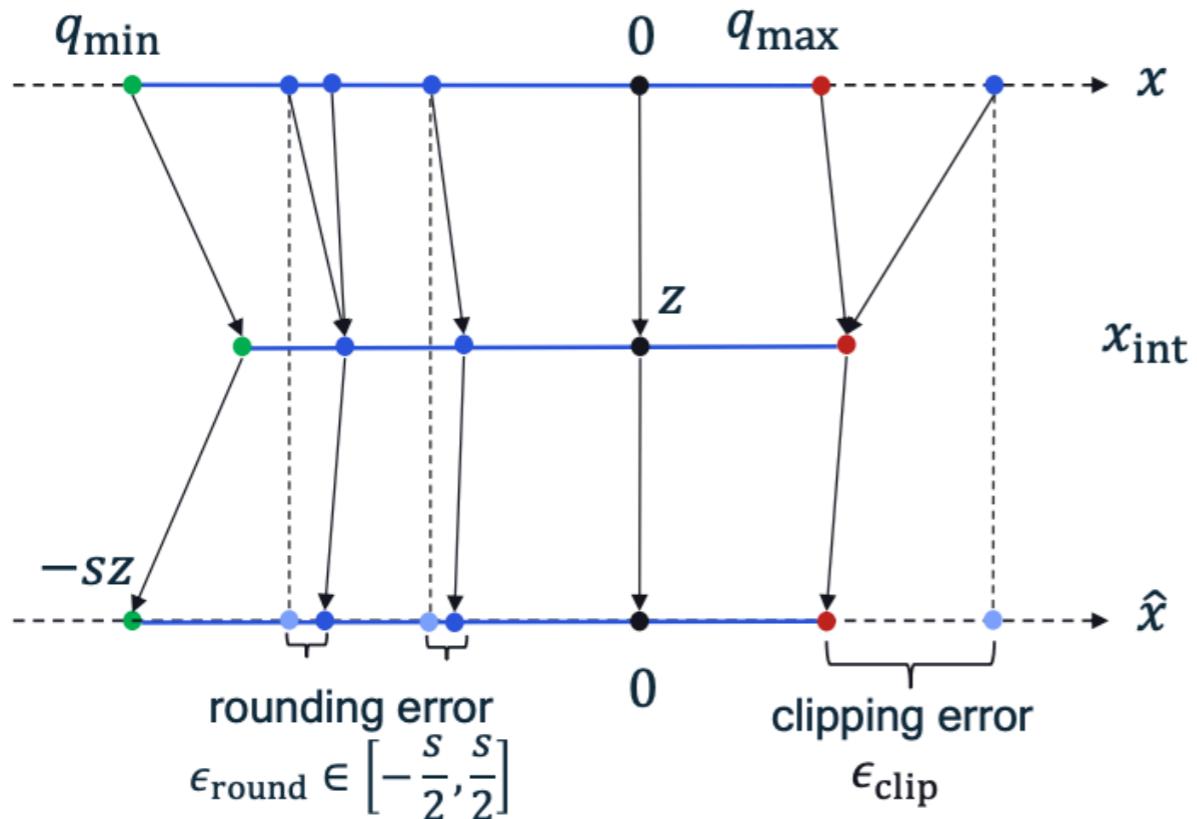
Real Forward pass



Simulated forward pass

$$\frac{\partial \lfloor x \rfloor}{\partial x} = 1$$

# Learning the quantization parameters



Learn quantization parameters during training using STE

$$x_{\text{int}} = \text{clamp} \left( \text{round} \left( \frac{x}{s} \right) + z, \min = 0, \max = 2^b - 1 \right)$$

$$\hat{x} = s (x_{\text{int}} - z)$$

Through task loss gradients, we find the optimal trade-off between  $\epsilon_{\text{clip}}$  &  $\epsilon_{\text{round}}$

[1] Esser, S. K., McKinstry, J. L., Bablani, D., Appuswamy, R., and Modha, D. S. Learned step size quantization, 2020

[2] Jain, S. R., Gural, A., Wu, M., and Dick, C. Trained uniform quantization for accurate and efficient neural network inference on fixed-point hardware.

[3] Bhalgat, Y., Lee, J., Nagel, M., Blankevoort, T., and Kwak, N. Lsq+: Improving low-bit quantization through learnable offsets and better initialization.

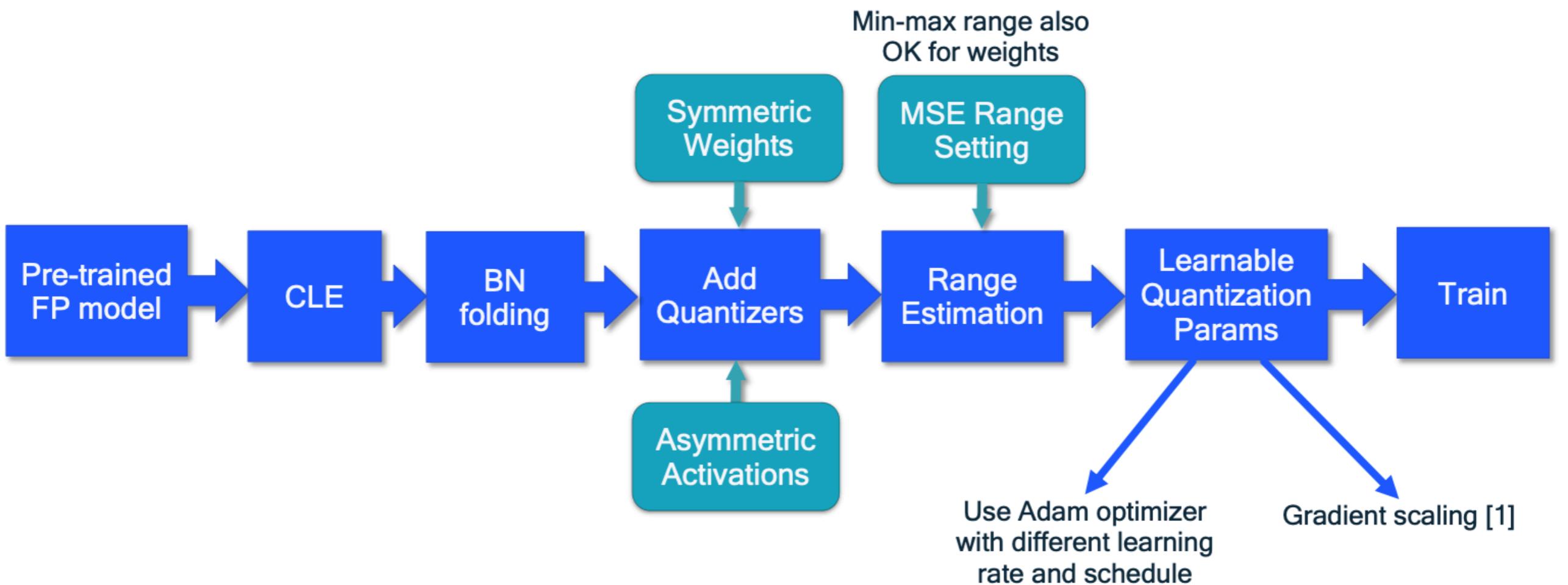
# Good initialization matters for QAT



Quantization setting	FP32	PTQ	QAT
W4A8 baseline	71.72	0.10	0.10
W4A8 w/ CLE	71.57	12.99	70.13
W4A8 w/ CLE + BC	71.57	46.90	70.07

Val. accuracy for MobileNetV2 for pet-tensor quantization

# Proposed QAT pipeline



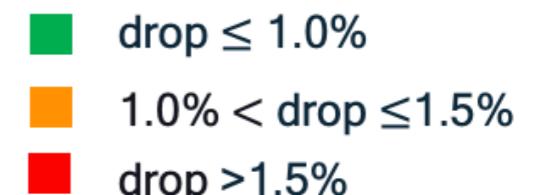
[1] Esser, S. K., McKinstry, J. L., Bablani, D., Appuswamy, R., and Modha, D. S. Learned step size quantization, 2020

# QAT results

█ drop  $\leq 1.0\%$   
█  $1.0\% < \text{drop} \leq 1.5\%$   
█ drop  $> 1.5\%$

Models	FP32	Per-tensor				Per-channel			
		W8A8	diff	W4A8	diff	W8A8	diff	W4A8	diff
ResNet18	69.68	70.38	+0.70	69.76	+0.08	70.43	+0.75	70.01	+0.33
ResNet50	76.07	76.21	+0.14	75.89	-0.18	76.58	+0.51	76.52	+0.45
MobileNetV2	71.72	71.76	+0.04	70.17	-1.55	71.82	+0.10	70.48	-1.24
InceptionV3	77.40	78.33	+0.93	77.84	+0.44	78.45	+1.05	78.12	+0.72
EfficientNet lite	75.42	75.17	-0.25	71.55	-3.87	74.75	-0.67	73.92	-1.50
DeepLabV3	72.94	73.99	+1.05	70.90	-2.04	72.87	-0.07	73.01	+0.07
EfficientDet-D1	40.08	38.94	-1.14	35.34	-4.74	38.97	-1.11	36.75	-3.33
BERT-base	83.06	83.26	+0.20	82.64	-0.42	82.44	-0.62	82.39	-0.67

# QAT and PTQ comparison



Difference from FP accuracy for W4A8 quantization

Models	FP32	Per-tensor		Per-channel	
		PTQ	QAT	PTQ	QAT
ResNet18	69.68	-1.06	+0.08	-0.77	+0.33
ResNet50	76.07	-0.92	-0.18	-0.64	+0.45
MobileNetV2	71.72	-2.51	-1.55	-1.93	-1.24
InceptionV3	77.40	-0.92	+0.44	-0.58	+0.72
EfficientNet lite	75.42	-4.18	-3.87	-1.41	-1.50
DeepLabV3	72.94	-2.14	-2.04	-1.27	+0.07
EfficientDet-D1	40.08	-39.77	-4.74	-5.00	-3.33
BERT-base	83.06	-1.30	-0.42	-1.04	-0.67

Relaxed Quantization for Discretized Neural Networks (Louizos, et al.)	ICLR 2019
Data-Free Quantization Through Weight Equalization and Bias Correction (Nagel, van Baalen, et al.)	ICCV 2019
Up or Down? Adaptive Rounding for Post-Training Quantization (Nagel, Amjad, et al.)	ICML 2020
Bayesian Bits: Unifying Quantization and Pruning (van Baalen, Louizos, et al.)	NeurIPS 2021
In-Hindsight Quantization Range Estimation for Quantized Training (Fournarakis, et al.)	CVPR 2021
A White Paper on Neural Network Quantization (Nagel, Fournarakis, et al.)	ArXiv 2021
Understanding and Overcoming the Challenges of Efficient Transformer Quantization (Bondarenko, et al.)	EMNLP 2021

Source sample text

Leading research in quantization