

Quantization - 2

EECE695D: Efficient ML Systems

Spring 2025

Agenda

- Linear quantization
- Various issues
 - Granularity
 - Rescaling
 - Clipping
 - Rounding
 - QAT

Linear Quantization

Idea

- Represent each weight as a **scaled-and-shifted integers**
- **Advantage.** Less compute, easy encoding/decoding
 - Matmuls can be done in integer, and scaled-and-shifted back

$$\begin{array}{c} [1.2, 2.4] \begin{bmatrix} -1.1 \\ 3.3 \end{bmatrix} \\ 4 \text{ FLOPs} \end{array} = \begin{array}{c} (1.2) \cdot (1.1) \cdot [1, 2] \begin{bmatrix} -1 \\ 3 \end{bmatrix} \\ 1 \text{ FLOP} + 4 \text{ Integer Op} \end{array}$$

Formalization

- Represent each weight as **scaled-and-shifted integers**
 - That is, our decoder is:

$$\hat{\mathbf{w}} = s \cdot (\mathbf{c} - z\mathbf{1})$$

- $\hat{\mathbf{w}}$: reconstructed weight
- \mathbf{c} : code (e.g., INT8 $\in \{-128, \dots, 127\}$)
- s : scaling factor (e.g., FP32)
- z : zero point (e.g., INT)

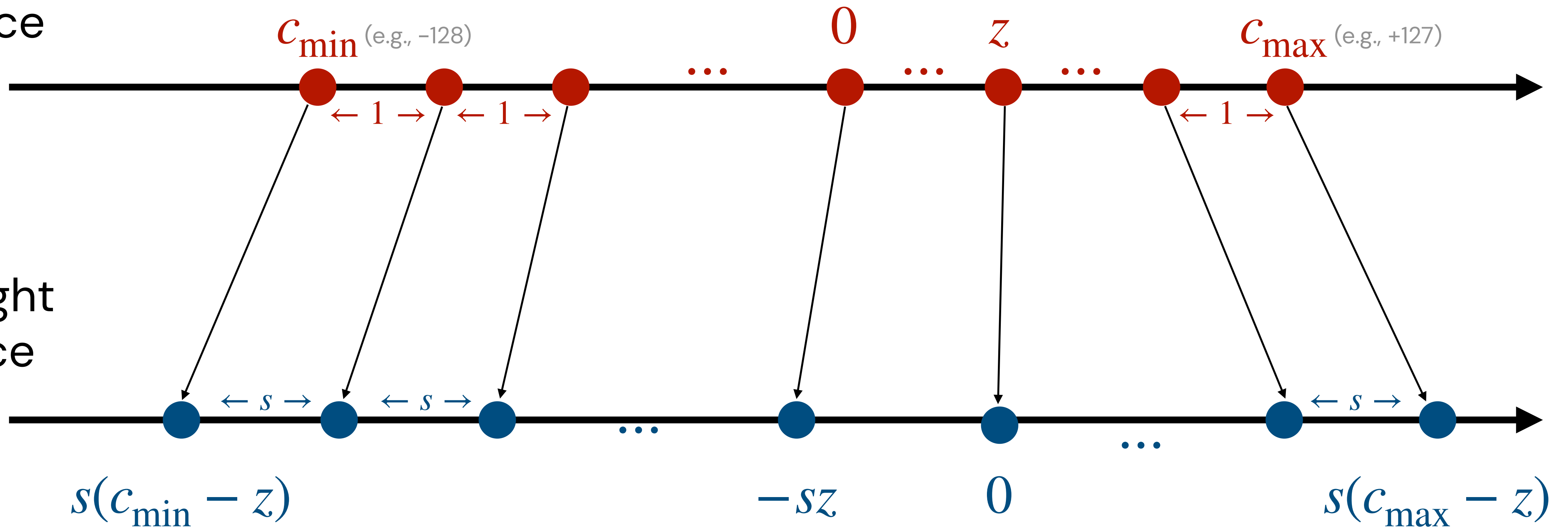
Formalization

$$\hat{\mathbf{w}} = s \cdot (\mathbf{c} - z\mathbf{1})$$

- Visually, what this decoder does is as follows:

Code space

Weight space



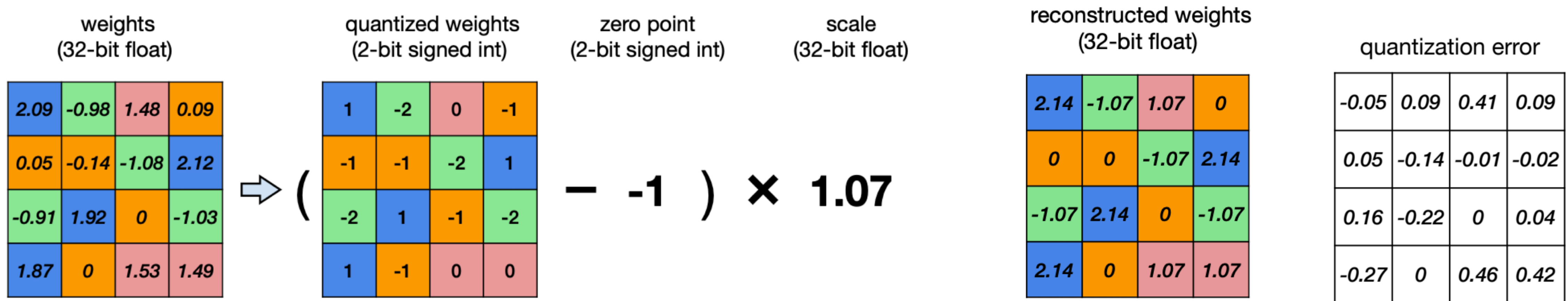
Formalization

$$\hat{\mathbf{w}} = s \cdot (\mathbf{c} - z\mathbf{1})$$

- Given this decoder, we encode to the nearest neighbor

$$\mathbf{c} = \text{round}(\mathbf{w}/s + z\mathbf{1})$$

- $\text{round}(\cdot)$: mapping to the nearest integer inside the range



Formalization

- Now we have the encoder and decoder
 - Encoder: $f(\mathbf{w}) = \text{round}(\mathbf{w}/s + z\mathbf{1})$
 - Decoder: $q(\mathbf{c}) = s \cdot (\mathbf{c} - z\mathbf{1})$
- **Want-to-do.** Given the weights \mathbf{w} , select the parameters s, z so that it solves

$$\min_{s,z} \hat{L}(g(f(\mathbf{w})))$$

- Of course, this is difficult; thus we use heuristic methods

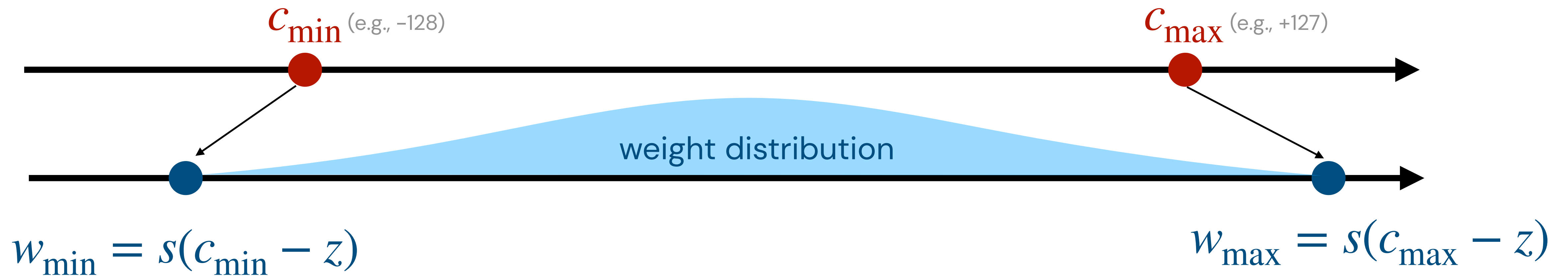
Minmax Quantization

- A crude but working way: **match the range!**

- That is, determine s, z so that

$$\min_i w_i = s(c_{\min} - z), \quad \max_i w_i = s(c_{\max} - z)$$

- Philosophy. Range that is just enough to capture the largest weights



Minmax Quantization

$$\min_i w_i = s(c_{\min} - z), \quad \max_i w_i = s(c_{\max} - z)$$

- Solving this, we get:

$$s = \frac{w_{\max} - w_{\min}}{c_{\max} - c_{\min}}, \quad z = c_{\min} - \text{round}(w_{\min}/s)$$

- One can do a similar thing to **quantize activations**:
 - Requires some “calibration data” to compute x_{\max}, x_{\min}
 - Quite brittle; often needs some clipping
- Note. It is also popular (and often better) to simply do “grid search”

Minmax Quantization

- **Brain teaser.** Suppose that we want to choose x_{\max}, x_{\min}
 - If data and model parameters are sharded over many servers, how much communication cost would we need?

(will be problematic for “dynamic quantization”)

Linear Quantization: Matmuls

Matmuls

- Consider the matmul:

$$\mathbf{Y} = \mathbf{W}\mathbf{X}$$

- Suppose that we have good quantizers for \mathbf{W} , \mathbf{X} , \mathbf{Y} :
 - That is, we have $s_{\mathbf{W}}$, $z_{\mathbf{W}}$, $s_{\mathbf{X}}$, $z_{\mathbf{X}}$, $s_{\mathbf{Y}}$, $w_{\mathbf{Y}}$

- Then, we get:

$$\begin{aligned} s_{\mathbf{Y}}(\mathbf{c}_{\mathbf{Y}} - z_{\mathbf{Y}}) &= s_{\mathbf{W}}(\mathbf{c}_{\mathbf{W}} - z_{\mathbf{W}}) \cdot s_{\mathbf{X}}(\mathbf{c}_{\mathbf{X}} - z_{\mathbf{X}}) \\ &= s_{\mathbf{W}}s_{\mathbf{X}}(\mathbf{c}_{\mathbf{W}}\mathbf{c}_{\mathbf{X}} - z_{\mathbf{W}}\mathbf{c}_{\mathbf{X}} - z_{\mathbf{X}}\mathbf{c}_{\mathbf{W}} + z_{\mathbf{W}}z_{\mathbf{X}}) \end{aligned}$$

Matmuls

- Rewriting, we have a formula for computing the **codes of Y**

$$\mathbf{c}_Y = \frac{s_W s_X}{s_Y} (\mathbf{c}_W \mathbf{c}_X - z_W \mathbf{c}_X - z_X \mathbf{c}_W + z_W z_X) + z_Y$$

- We have separated out FP ops from INT ops
- **Problem.** To compute \mathbf{c}_Y , we need (FP) * (INT) operation
 - Empirically, $s_W s_X / s_Y \in (0,1)$
 - \Rightarrow Write it as $2^{-n} \times M_0$, with M_0 being an INT
(bit shift)

Matmuls

- Some INT ops can be pre-computed, reducing inference-time compute

$$\mathbf{c}_Y = \frac{s_W s_X}{s_Y} (\mathbf{c}_W \mathbf{c}_X - z_W \mathbf{c}_X - z_X \mathbf{c}_W + z_W z_X) + z_Y$$

- Also, many weight distributions are nearly symmetric:

- If we let $z_W = 0$,

$$\mathbf{c}_Y = \frac{s_W s_X}{s_Y} (\mathbf{c}_W \mathbf{c}_X - z_X \mathbf{c}_W) + z_Y$$

- Not doable for activations, usually.

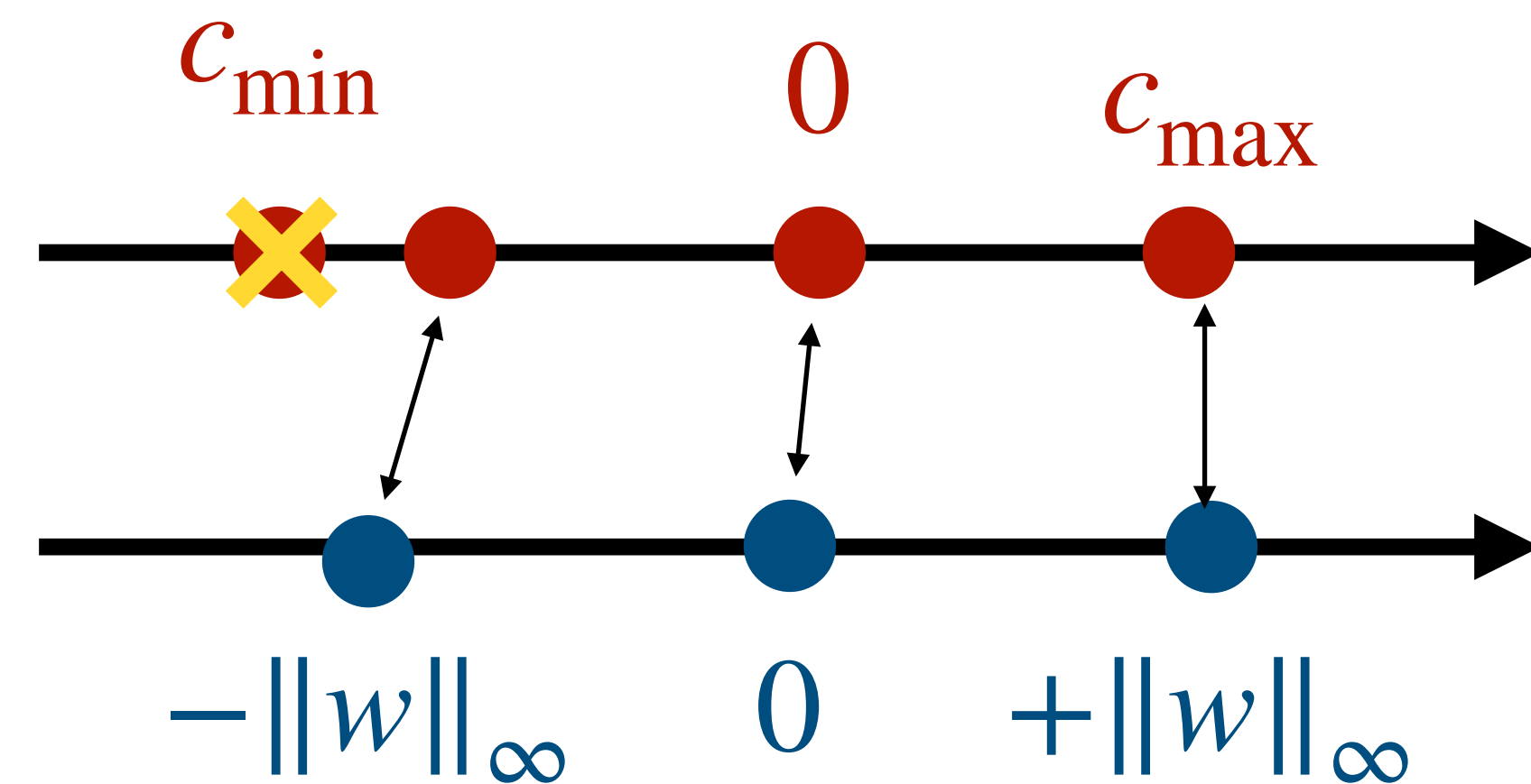
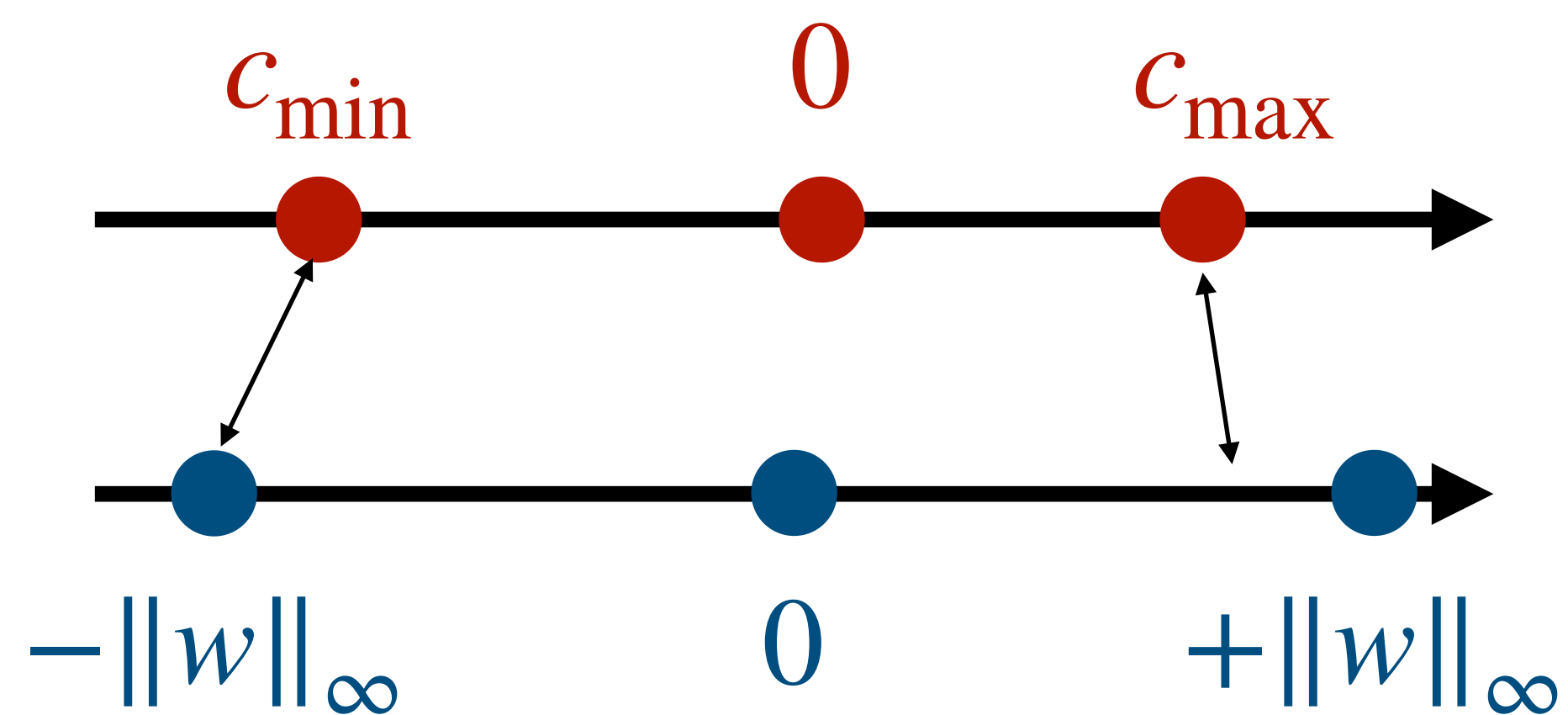
Symmetric Quantization

- **Problem.** INT is not symmetric! (e.g., $\{-128, \dots, 127\}$)

- Two different ways to do symmetric quantization:

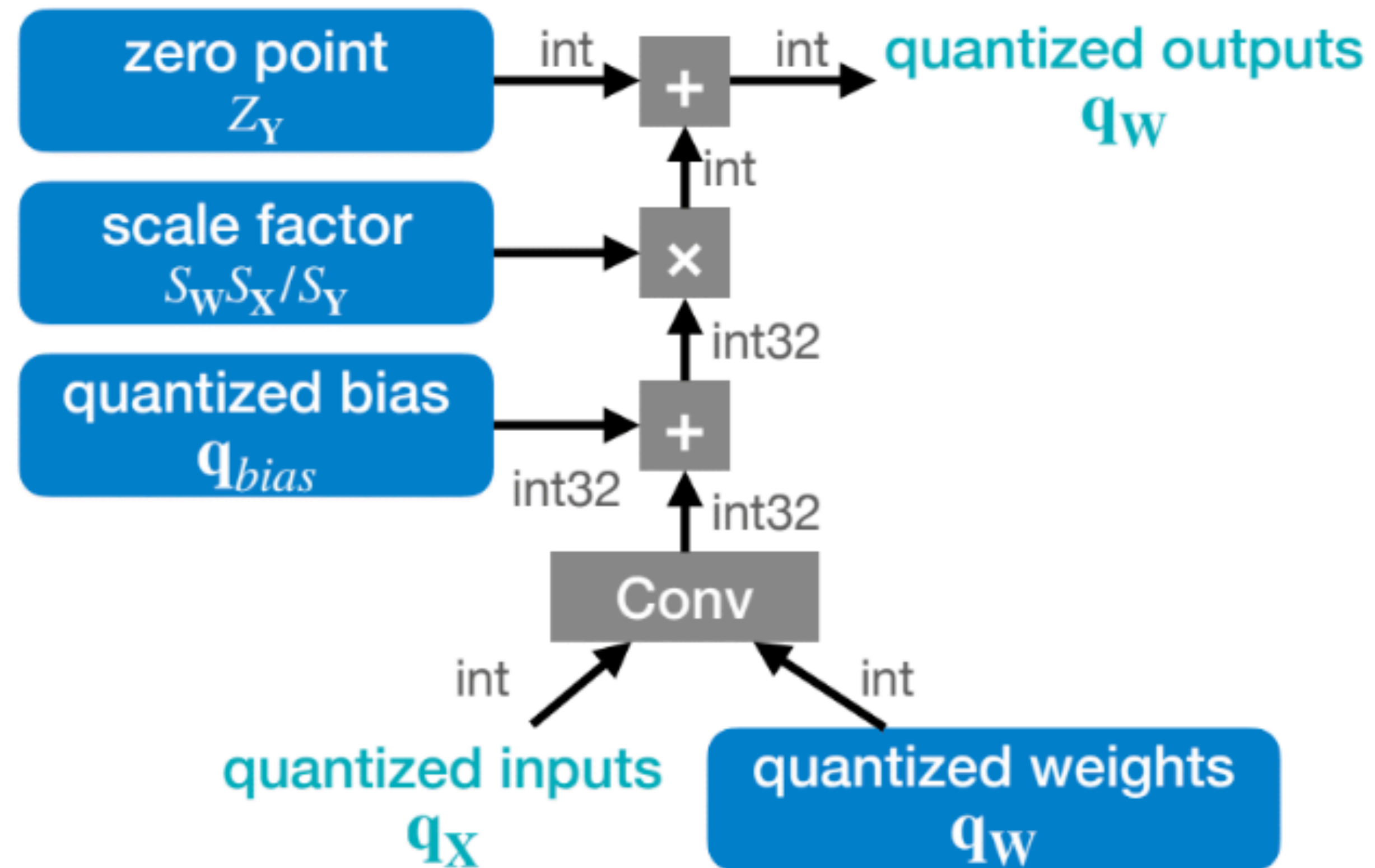
- **Full-range** (\Leftarrow). ONNX, PyTorch, ...

- **Restricted-Range** (\Rightarrow). TensorFlow, TensorRT, ...



Symmetric Quantization

- **Note.** Accumulation can take place in high-bits (e.g., INT32)



Symmetric Quantization

- Integer-only ops can dramatically reduce the latency

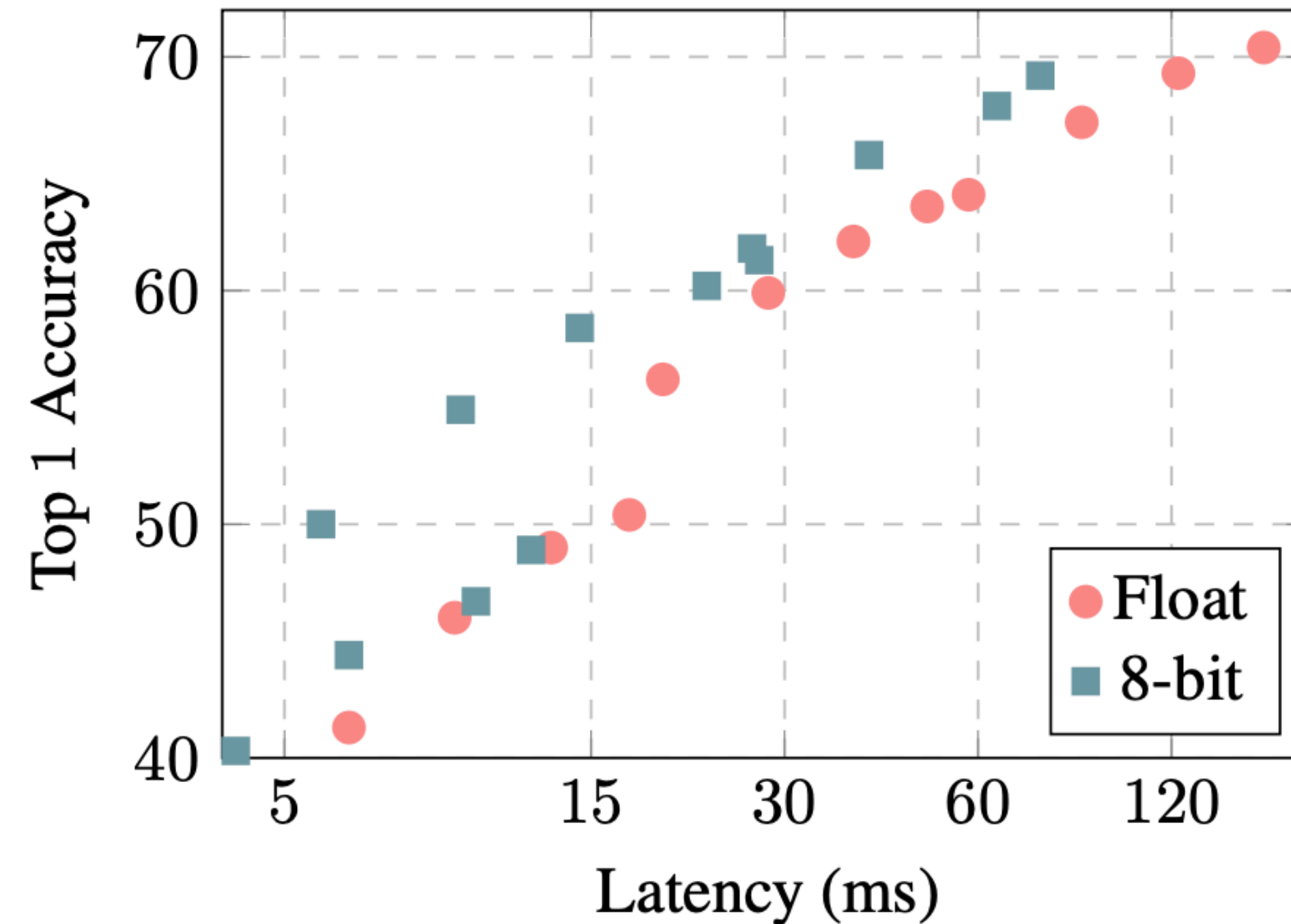


Figure 4.1: ImageNet classifier on Qualcomm Snapdragon 835 big cores: Latency-vs-accuracy tradeoff of floating-point and integer-only MobileNets.

Advanced ideas for PTQ

Agenda

- In most PTQ algorithms, we adopt more ideas:
 - Finer granularity
 - Weight rescaling
 - Activation clipping
 - Adaptive rounding

Granularity

- **Motivation.** Weight ranges are quite dissimilar in different dimensions

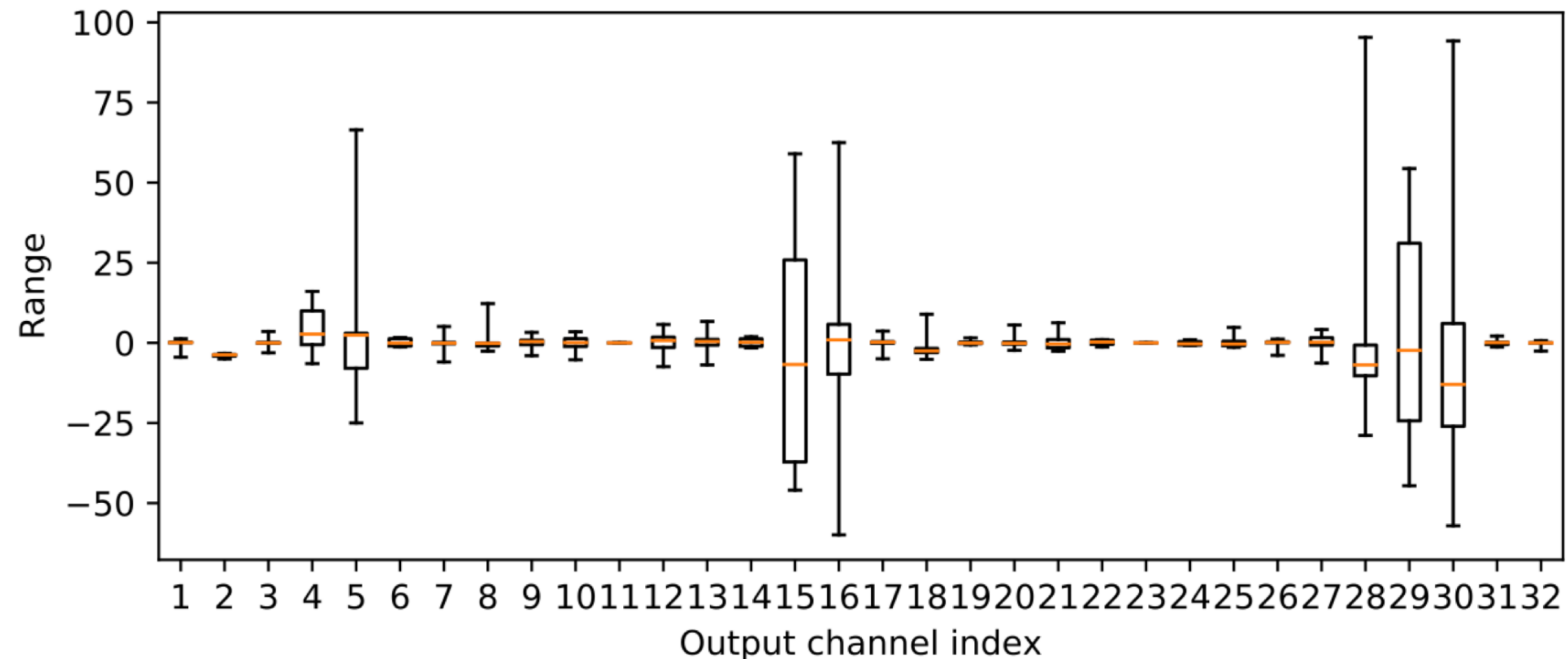
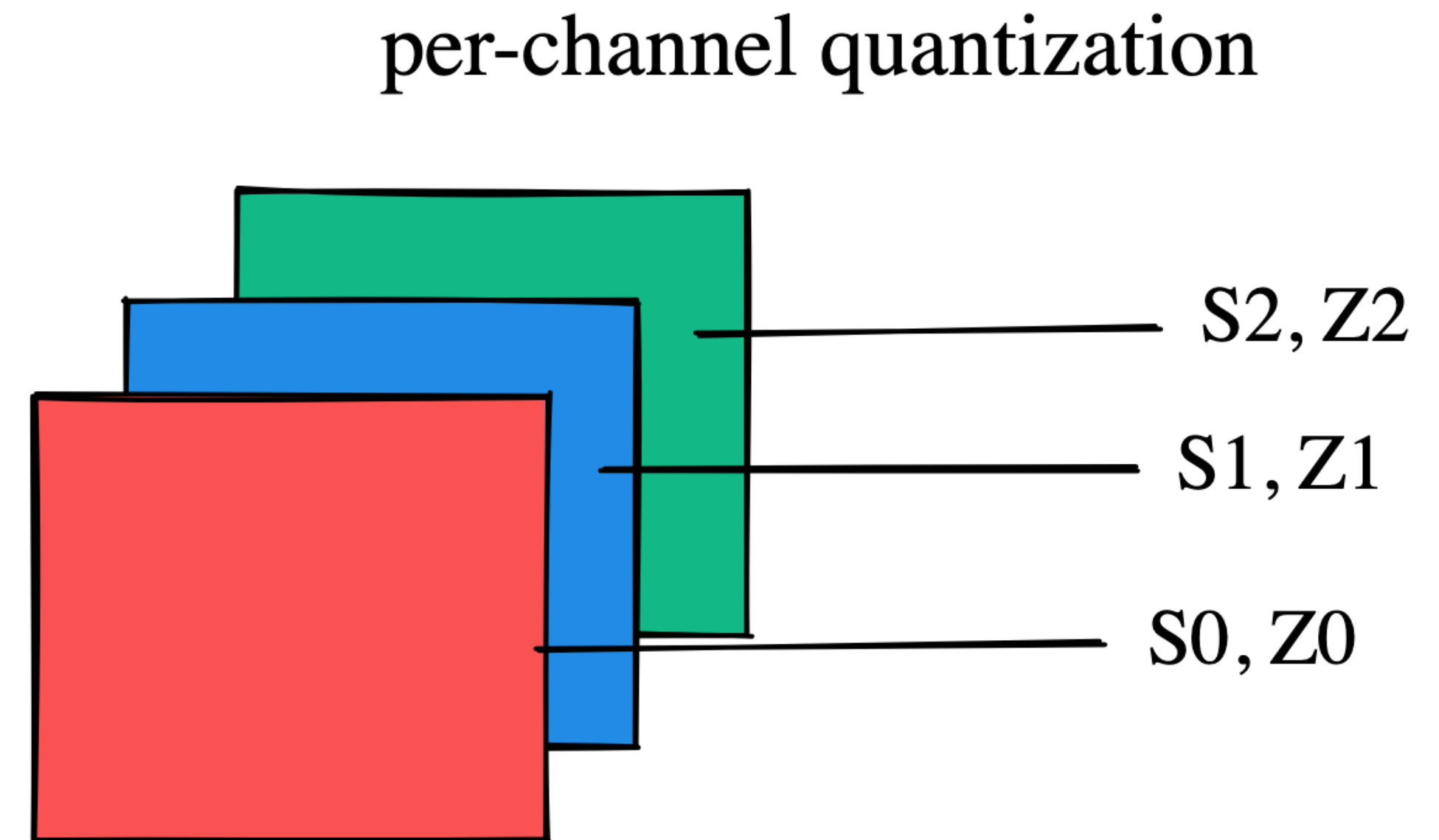
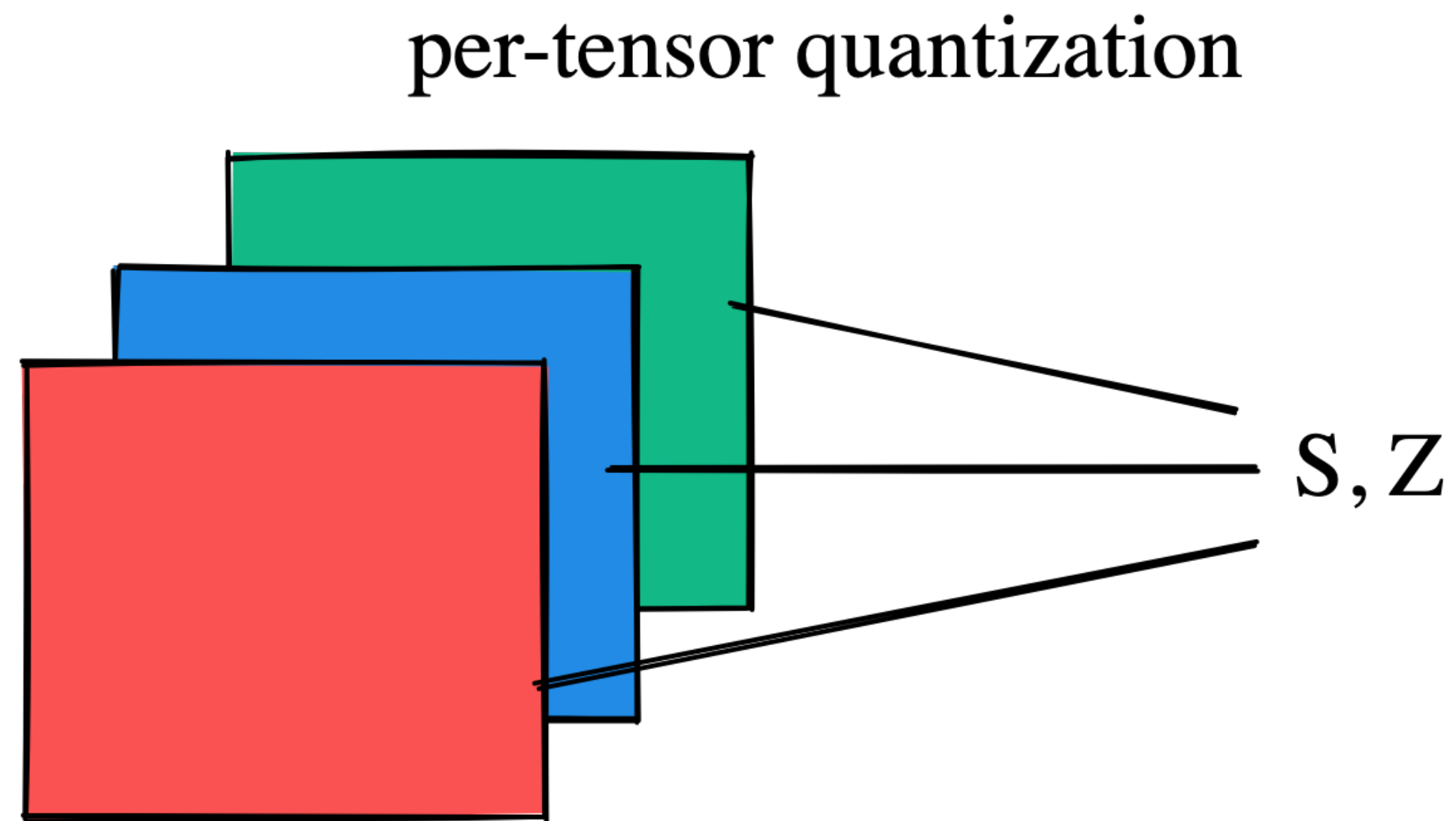


Figure 2. Per (output) channel weight ranges of the first depthwise-separable layer in MobileNetV2. In the boxplot the min and max value, the 2nd and 3rd quartile and the median are plotted for each channel. This layer exhibits strong differences between channel weight ranges.

Granularity

- **Idea.** Apply different (s, z) for different group of weights (and/or activations)
 - Example. Per-Channel Quantization



Granularity

- Example. Per-Vector Quantization
 - Not much degradation in speed, if computation is done in the units of fixed-length vectors

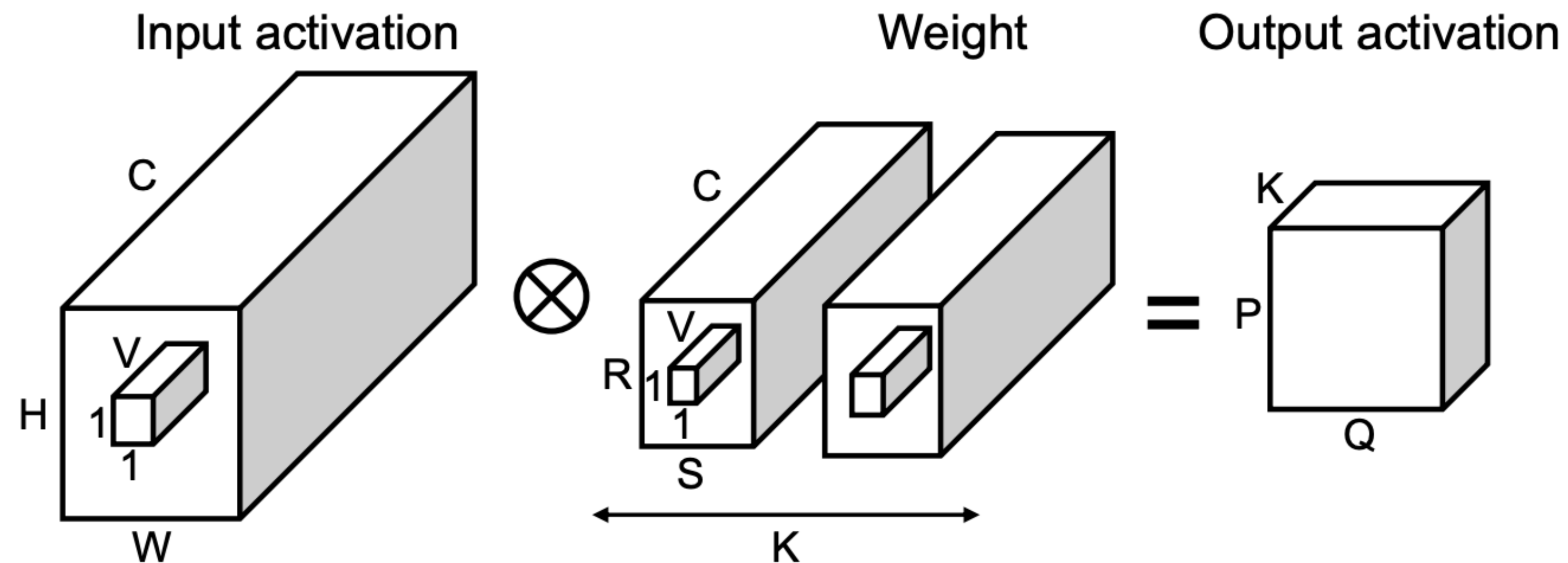


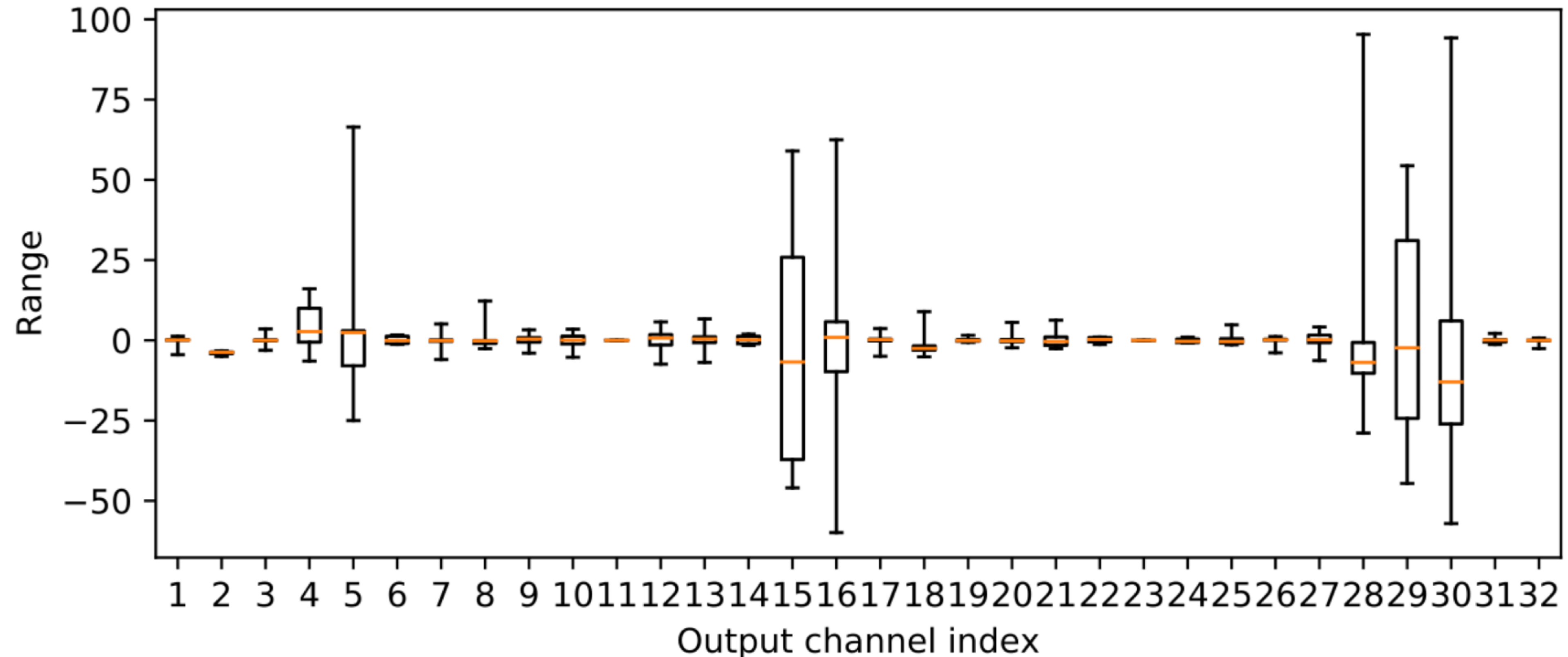
Figure 1. Convolution — Comparison between per-layer/per-output-channel scaling and per-vector scaling.

Granularity

- **Further readings**
 - Use two-stage scaling factors ([link](#))
 - Sharing micro-exponents ([link](#))

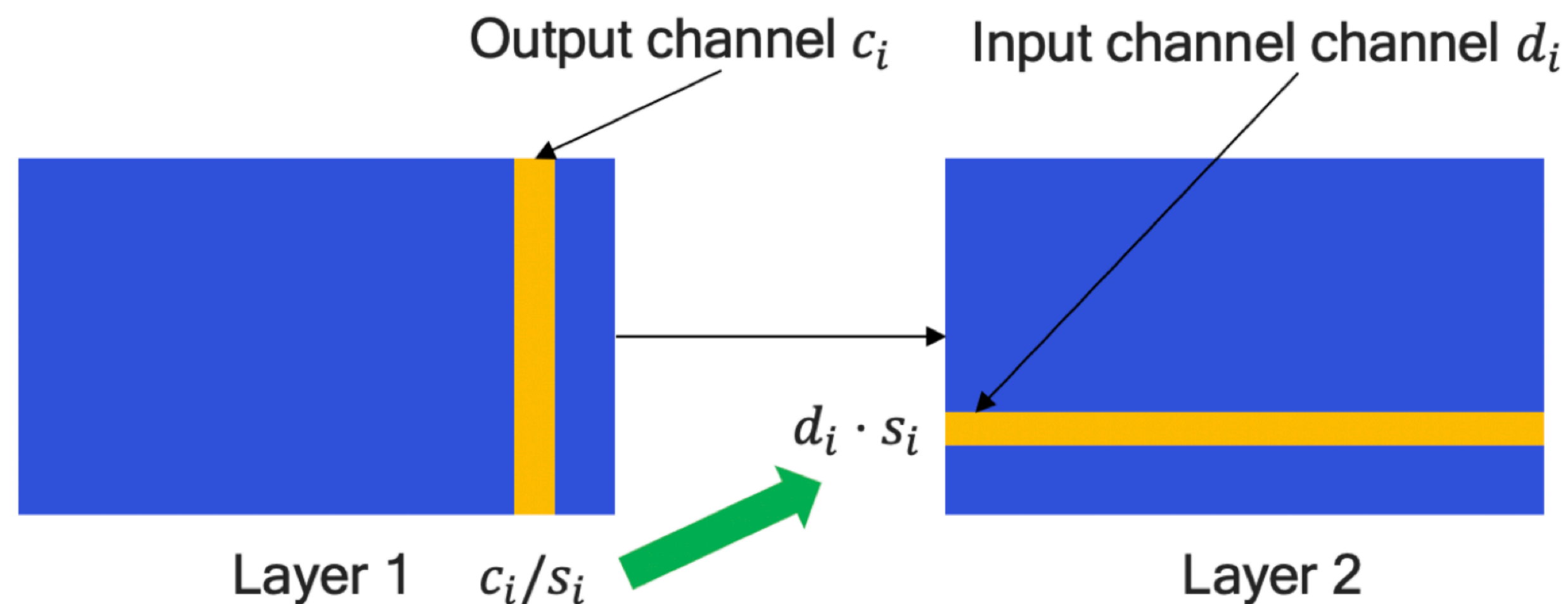
Rescaling

- **Idea.** Tackle the same problem, but use the **positive homogeneity** of ReLU.
 - i.e., $\sigma(cx) = c\sigma(x), \quad \forall c > 0$



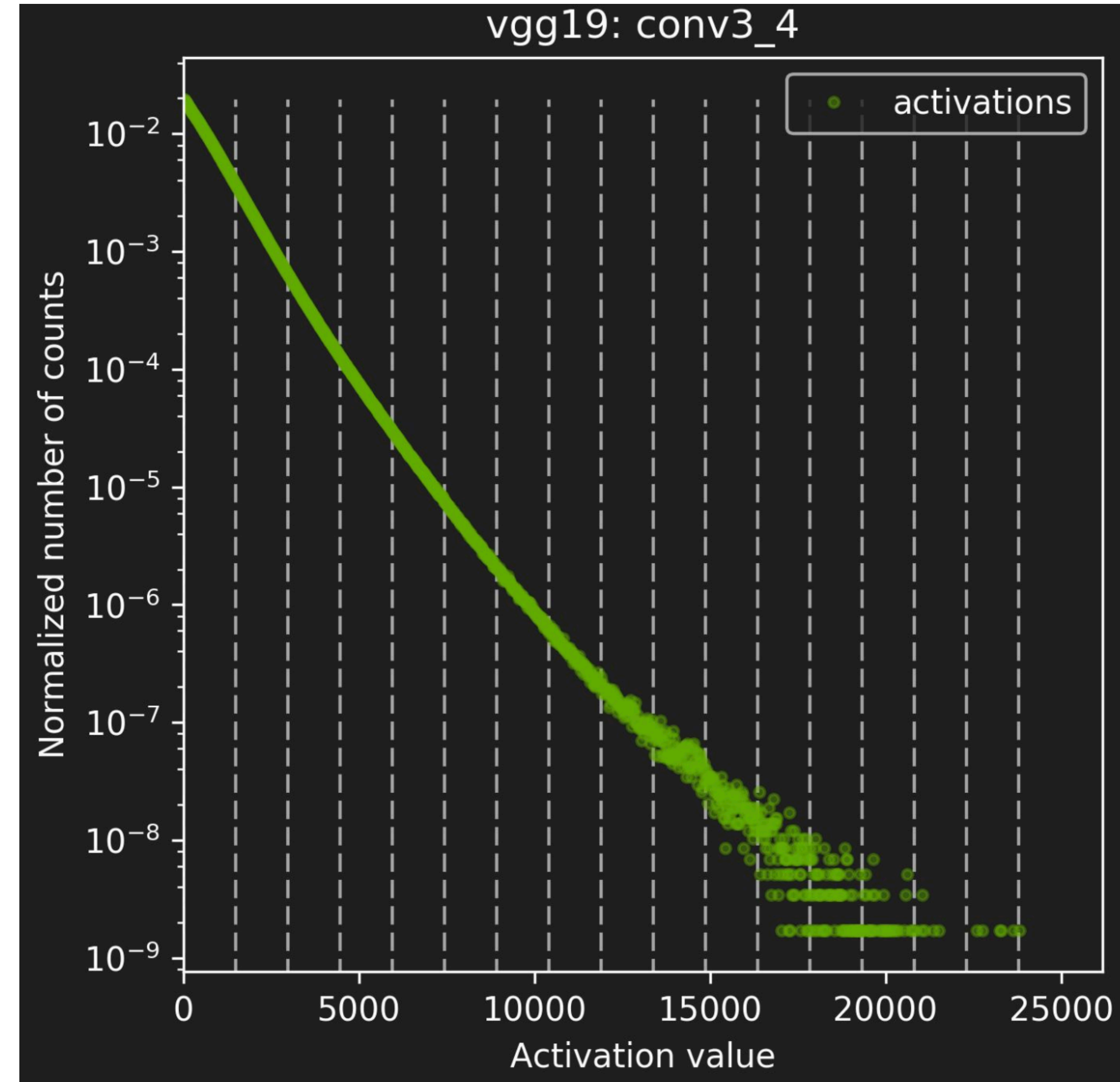
Rescaling

- Take a ReLU neural net
 - Multiply $\times 100$ to all weights in j th output channel of layer i
 - Multiply $\times 0.01$ to all weights in j th input channel of layer $i+1$
 - Identical function, with different weights
- Do this many times to match the weight range



Clipping

- **Motivation.**
- For activation, $s_{\mathbf{x}}$ is determined via:
 - During training. Take an exponential moving average
 - After training. Use calibration batches
- However, many **outliers** appear



Clipping

- **Idea.** Clip the activations
- **Question.** Where do we clip?
 - Explicit optimization.
 - Approximates w/ Gaussian/Laplace and minimize ℓ^2 (e.g., ACIQ)
 - Minimize ℓ^2 via Newton-Raphson (e.g., OCTAV)
 - Minimize the KL-divergence b/w quantized & reference dist (link)

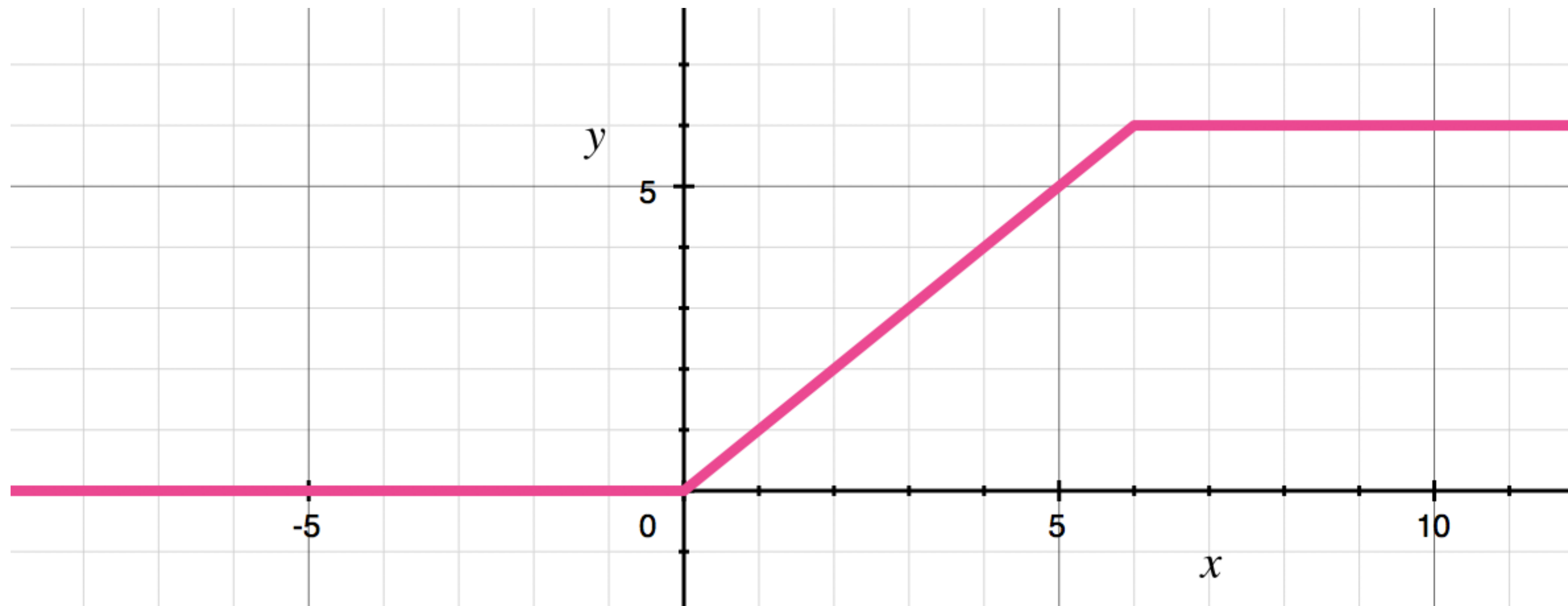
Clipping

- Use clipping activation.

- ReLU6

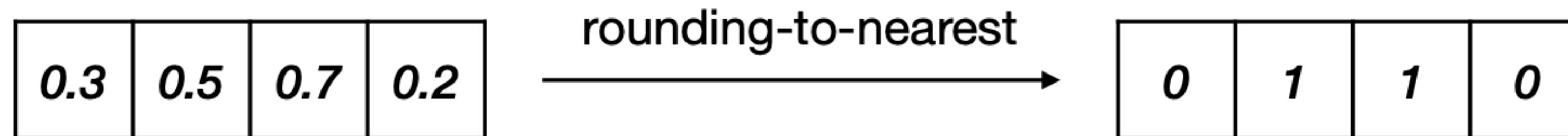
(brainteaser: why 6?)

- ReLU with learnable clipping range (e.g., PACT)



Rounding

- **Motivation.** Turns out that round-to-nearest (RTN) is suboptimal
 - In fact, **stochastic rounding** gives better options than RTN



Rounding scheme	Acc(%)
Nearest	52.29
Ceil	0.10
Floor	0.10
Stochastic	52.06±5.52
Stochastic (best)	63.06

Rounding

- **Idea.** Round up-or-down in a way that **minimizes the activation distortion**
 - More concretely, AdaRound solves

$$\min_{\mathbf{V}: v_i \in [0,1]} \|\mathbf{W}\mathbf{X} - \tilde{\mathbf{W}}\mathbf{X}\|_F^2 + \lambda \cdot f_{\text{reg}}(\mathbf{V})$$

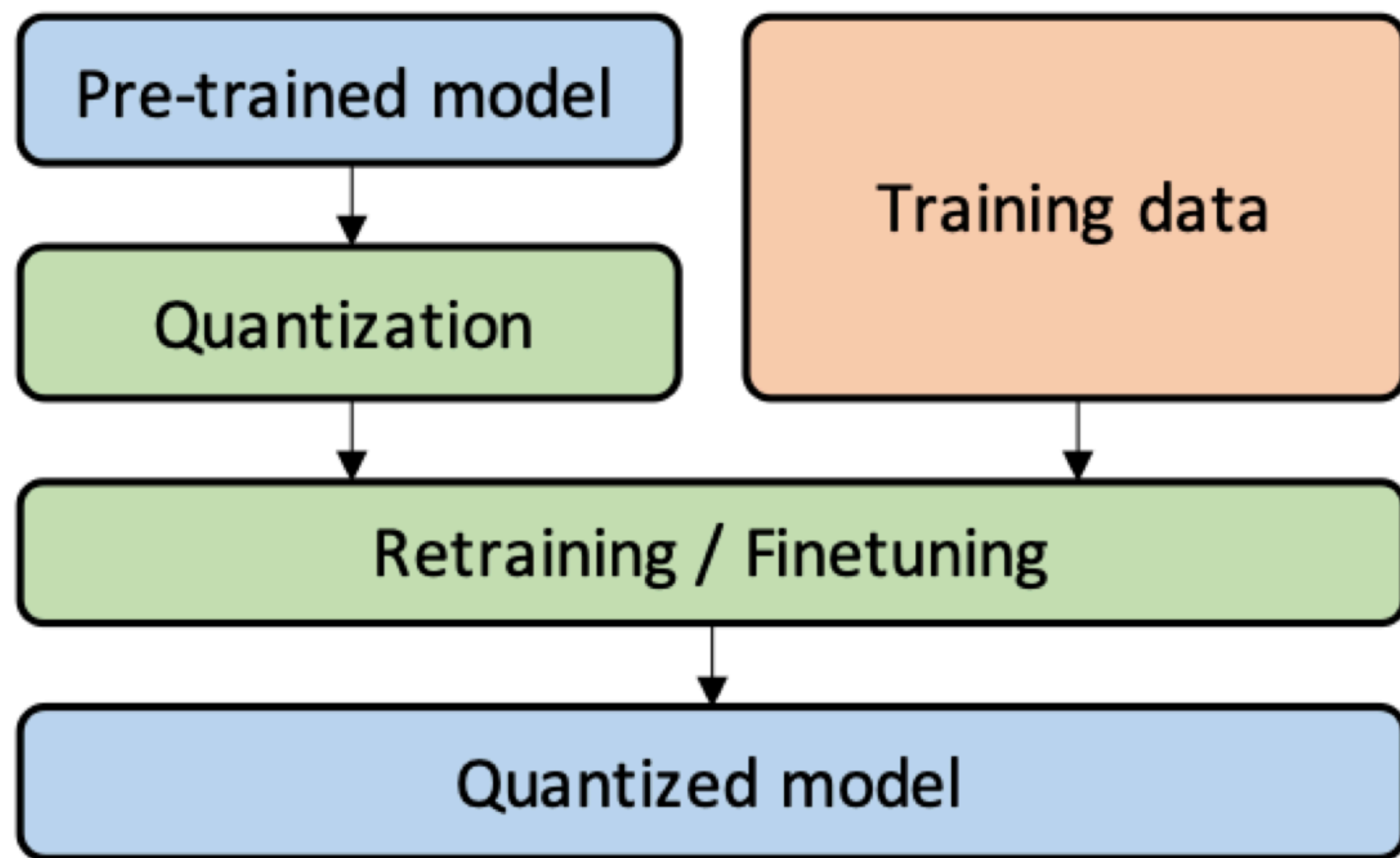
- \mathbf{V} : weight shift for round up/down
- $\tilde{\mathbf{W}} = s \cdot (\lfloor \mathbf{W}/s \rfloor + \mathbf{V})$: quantized weight
- $f_{\text{reg}}(\mathbf{V}) = \sum_i 1 - |2\mathbf{V}_i - 1|^\beta$: binary-forcing regularizer

- See also: [AdaQuant](#), [FlexRound](#)

QAT

Quantization-aware training

- After quantization, fine-tune the weight
 - Recovers much of the lost accuracy



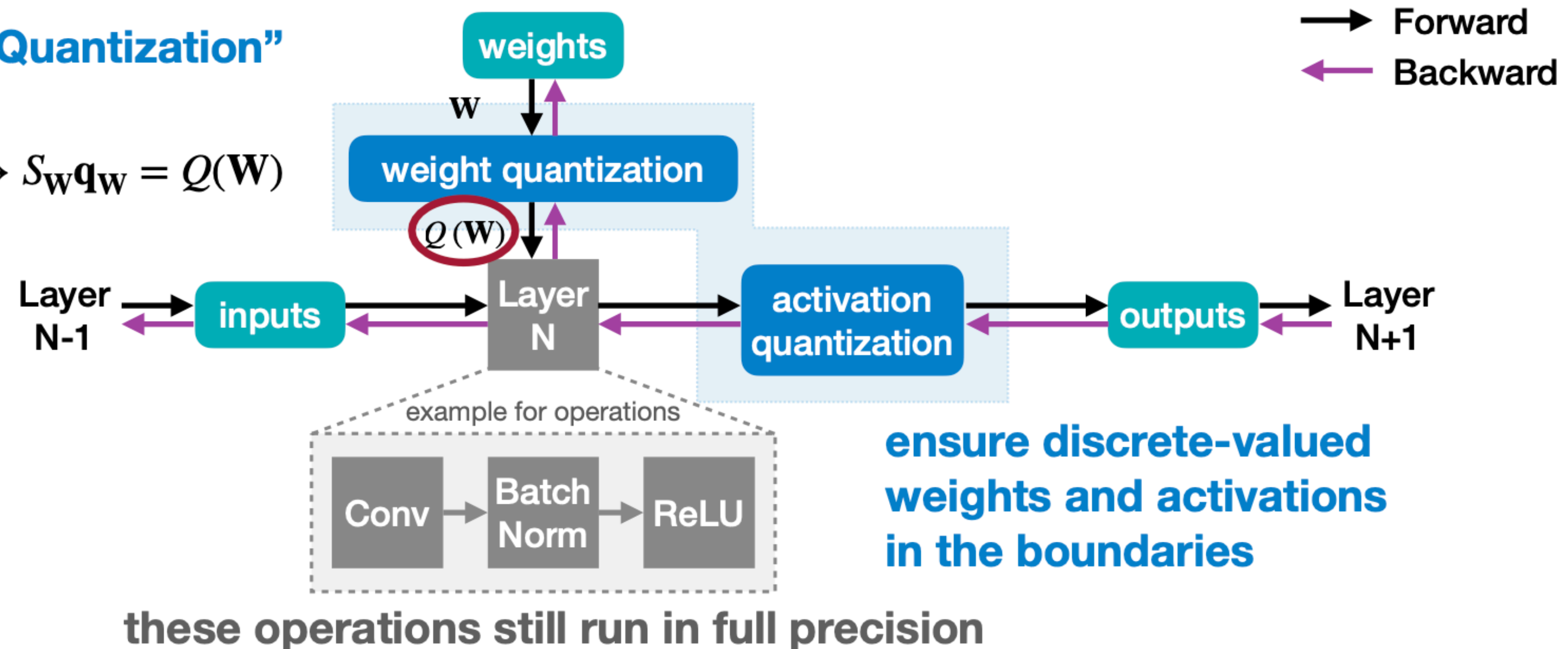
Neural Network	Floating-Point	Post-Training Quantization		Quantization-Aware Training	
		Asymmetric	Symmetric	Asymmetric	Symmetric
		Per-Tensor	Per-Channel	Per-Tensor	Per-Channel
MobileNetV1	70.9%	0.1%	59.1%	70.0%	70.7%
MobileNetV2	71.9%	0.1%	69.8%	70.9%	71.1%
NASNet-Mobile	74.9%	72.2%	72.1%	73.0%	73.0%

Quantization-aware training

- (1) Keep full-precision weight, and simulate quantization at forward
- (2) Compute gradients in full-precision
- (3) Update full-precision weights

“Simulated/Fake Quantization”

$$W \rightarrow S_W q_W = Q(W)$$



Quantization-aware training

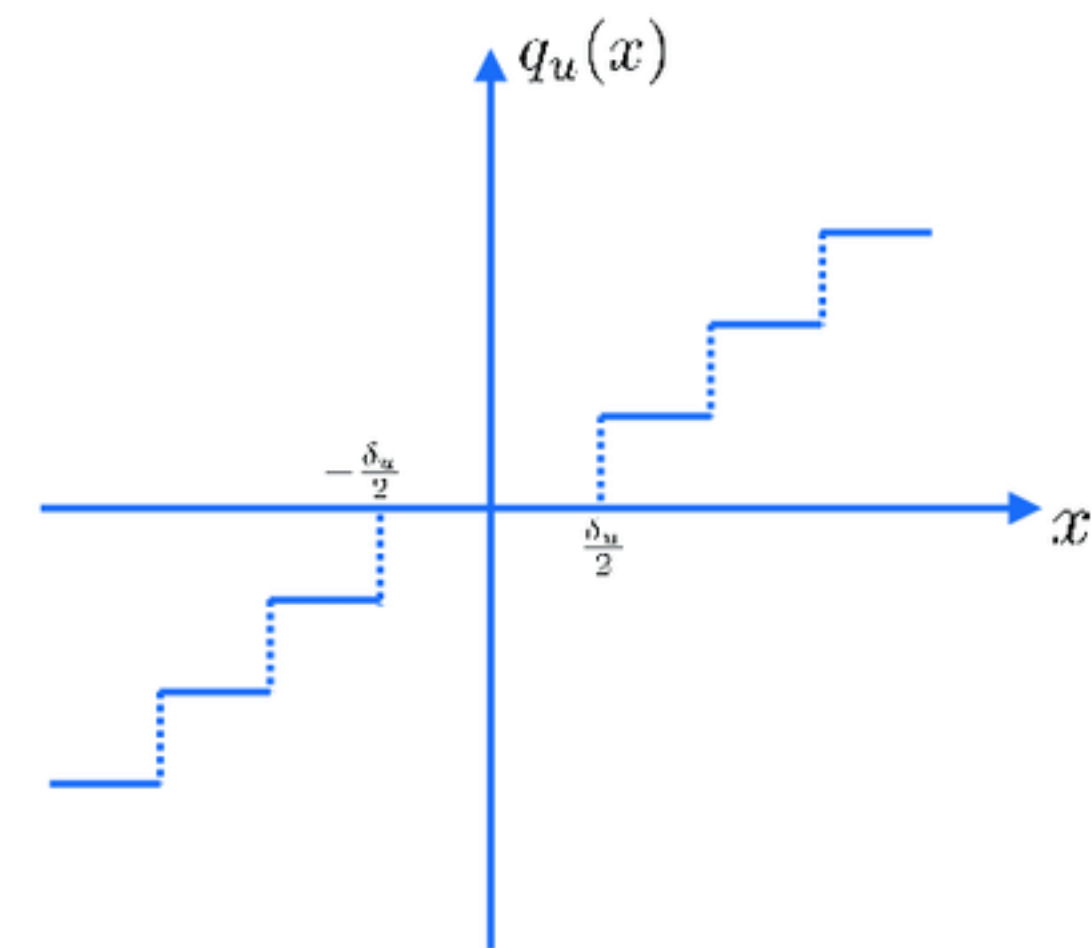
- Again, we are solving $\min_{\mathbf{w}} \hat{L}(q(\mathbf{w}))$
- **Challenge.** Computing gradients through discretizing operation $q(\cdot)$
 - Example. Consider a linear regression

$$\min_{\mathbf{w}} (y - q(\mathbf{w})^T \mathbf{x})^2$$

- The gradient is

$$2(q(\mathbf{w})^T \mathbf{x} - y) \cdot \mathbf{x}^T \nabla_{\mathbf{w}} q(\mathbf{w})$$

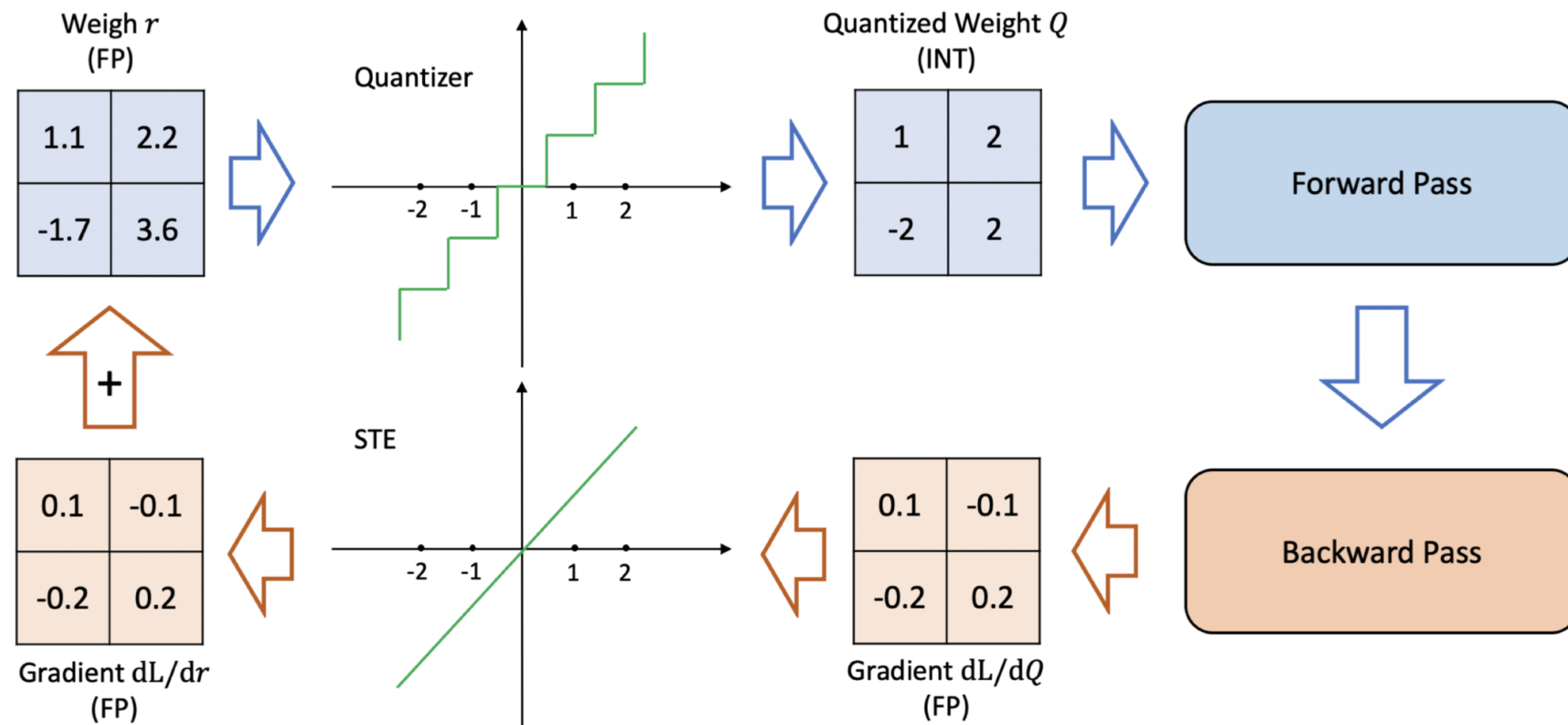
- The **red term** is always zero-or-infinity!



Straight-through estimator

- Again, we use STE
 - i.e., ignore quantization during backward

$$2(q(\mathbf{w})^\top \mathbf{x} - y) \cdot \mathbf{x}^\top \nabla_{\mathbf{w}} q(\mathbf{w}) \Rightarrow 2(q(\mathbf{w})^\top \mathbf{x} - y) \cdot \mathbf{x}^\top \nabla_{\mathbf{w}} (\mathbf{w})$$



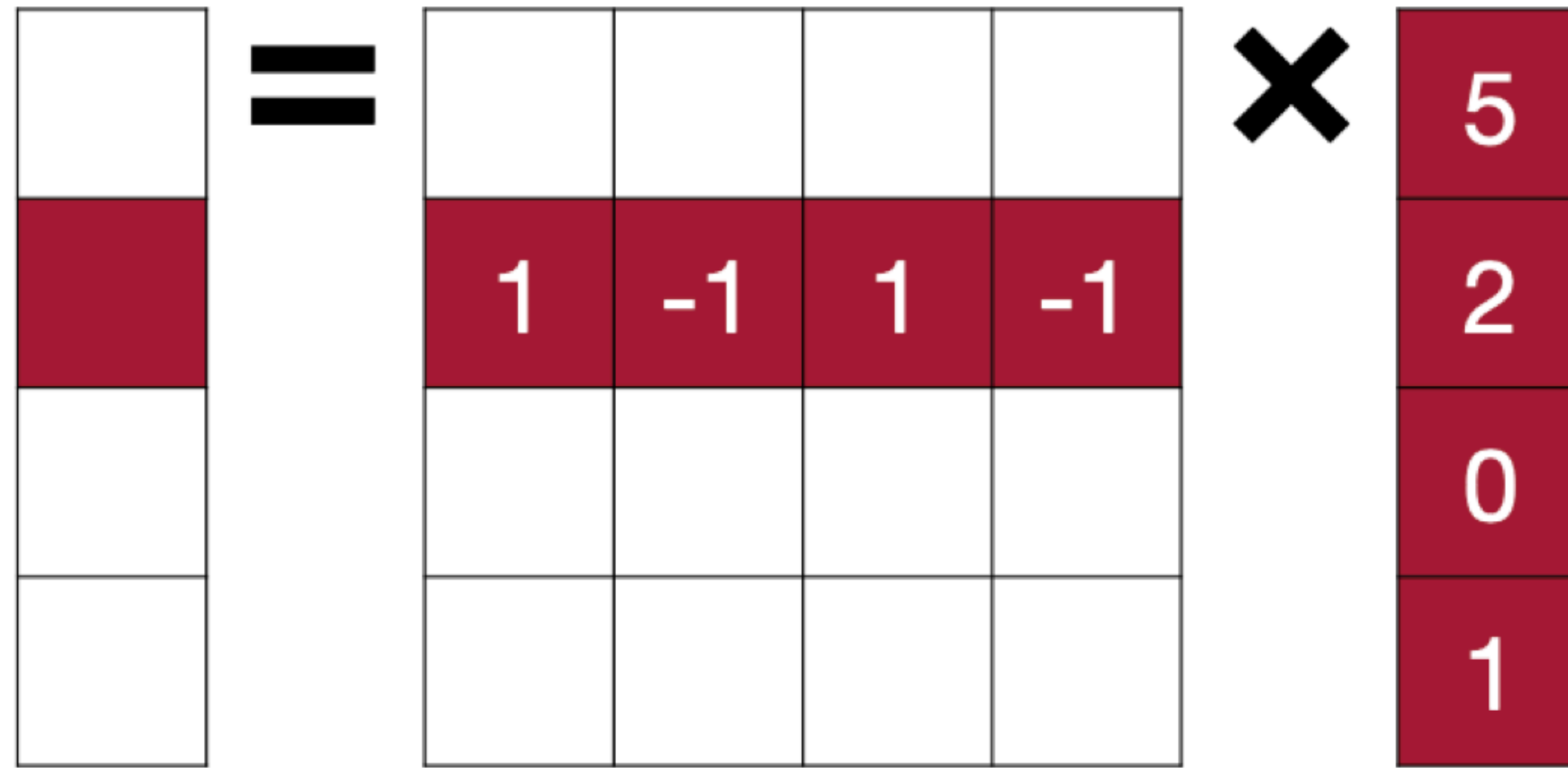
Further reading

- Some theoretical arguments suggest that using (clipped) ReLU, instead of identity function, is a better choice for STE
 - <https://arxiv.org/abs/1903.05662>

Other topics

Binary nets

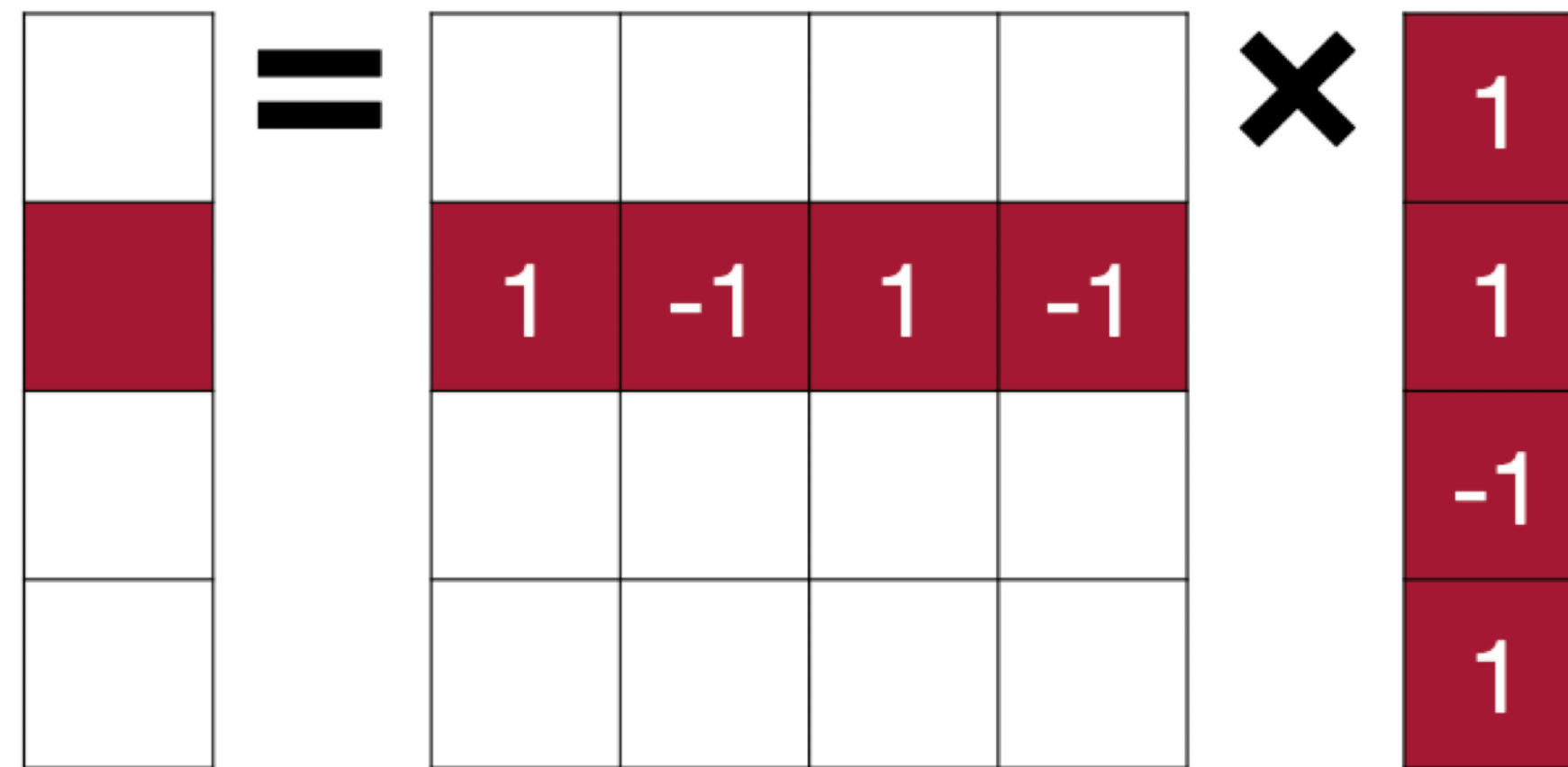
- If all weights are binary, then we need no multiplications



$$y = 5 + (-2) + 0 + (-1)$$

Binary nets

- If all activations are binary as well, then we only need XNOR + Counting

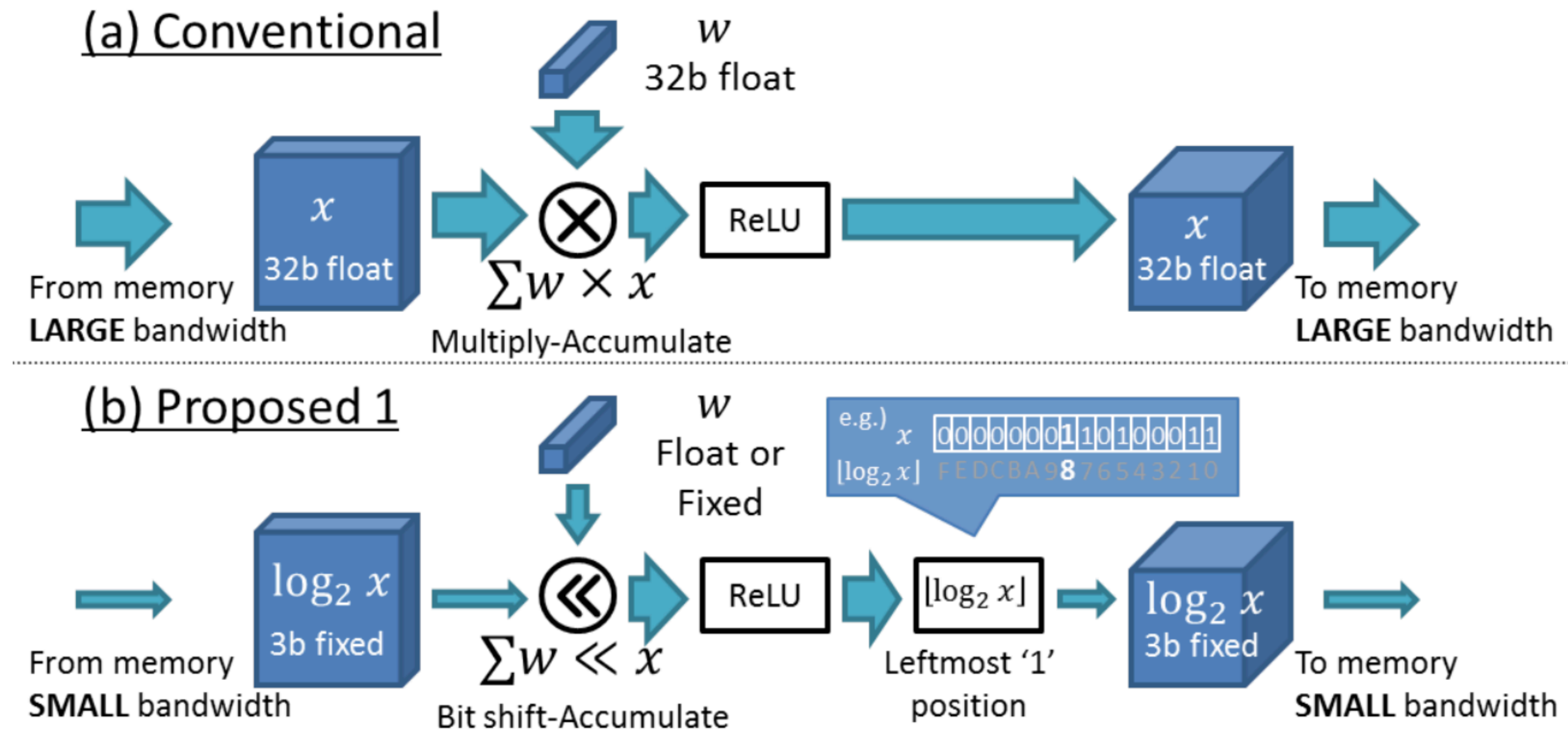


W	X	Y=W \times X
1	1	1
1	-1	-1
-1	-1	1
-1	1	-1

b _w	b _x	XNOR(b _w , b _x)
1	1	1
1	0	0
0	0	1
0	1	0

LogQuant

- Use logarithmically quantized activations
 - Multiplications are simply shifting bits



Next Class

- Knowledge distillation

That's it for today 🙌