

Edge Computing

Lecture 06: Quantization and Pruning

Recap

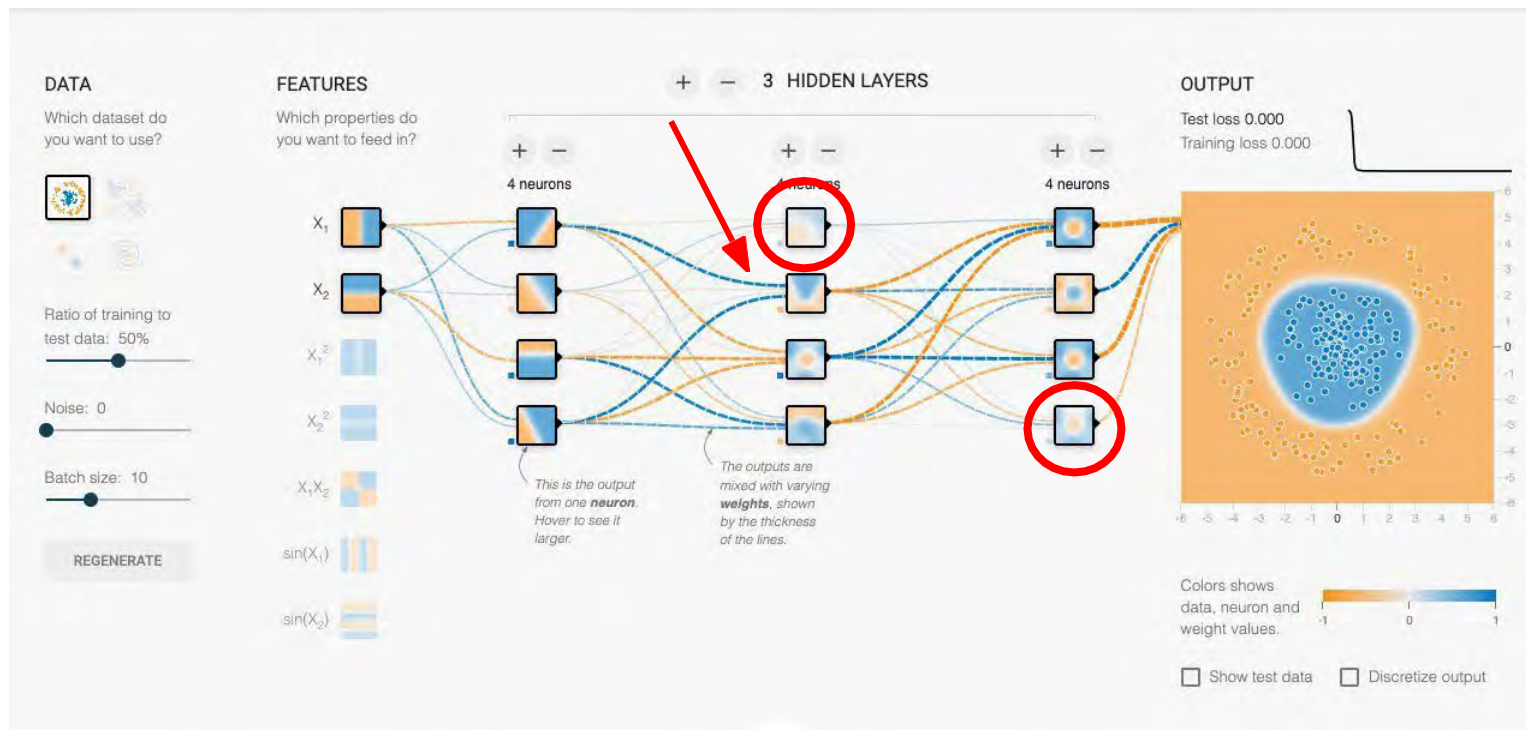
- Neural network: terminology
 - Neuron, Synapse, Activation, Weight, etc.
- Common building block: layer
 - FC, Conv, Conv2D, Depthwise Conv
 - Receptive field, padding, strides
 - Pooling, Normalization
- Convolution neural network
 - Alexnet, VGG16, MobileNetv2
 - Tensorflow, Tensorflow Lite

Agenda

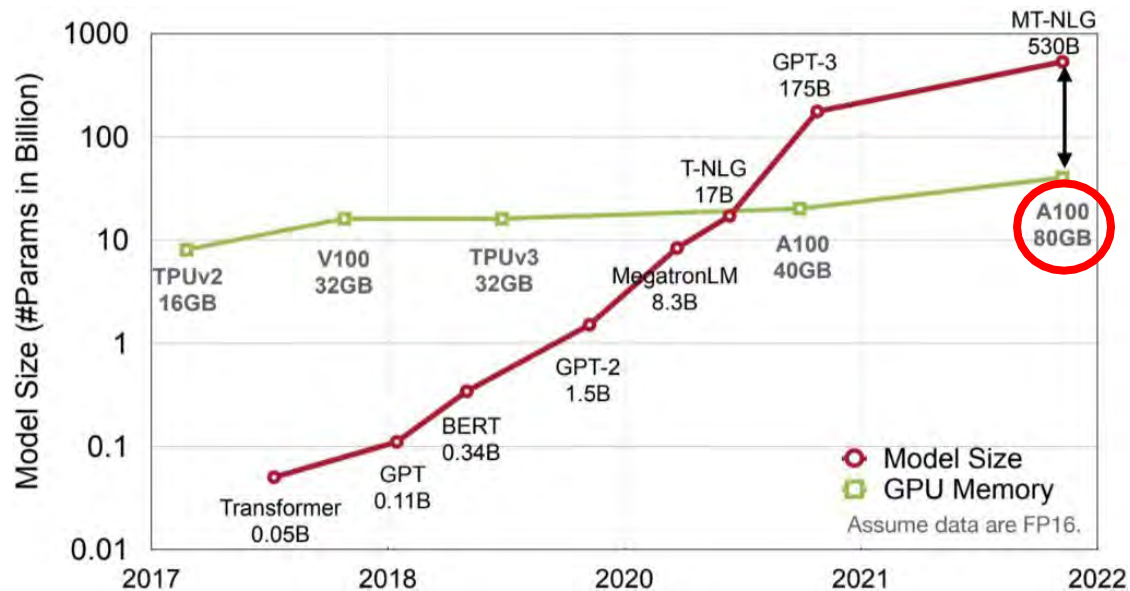
- Pruning
 - What and why?
 - How? Pruning criterion
 - Fine-tune/retrain acceleration
- Quantization
 - What and why?
 - Linear quantization
 - Quantization granularity
 - Calibration and Clipping

What is Pruning?

- Shrinking models by removing synapses and neurons.



Why pruning?



Why pruning? Model Footprint

- Nvidia Jetson
 - Custom boards vs Dev boards



Nano

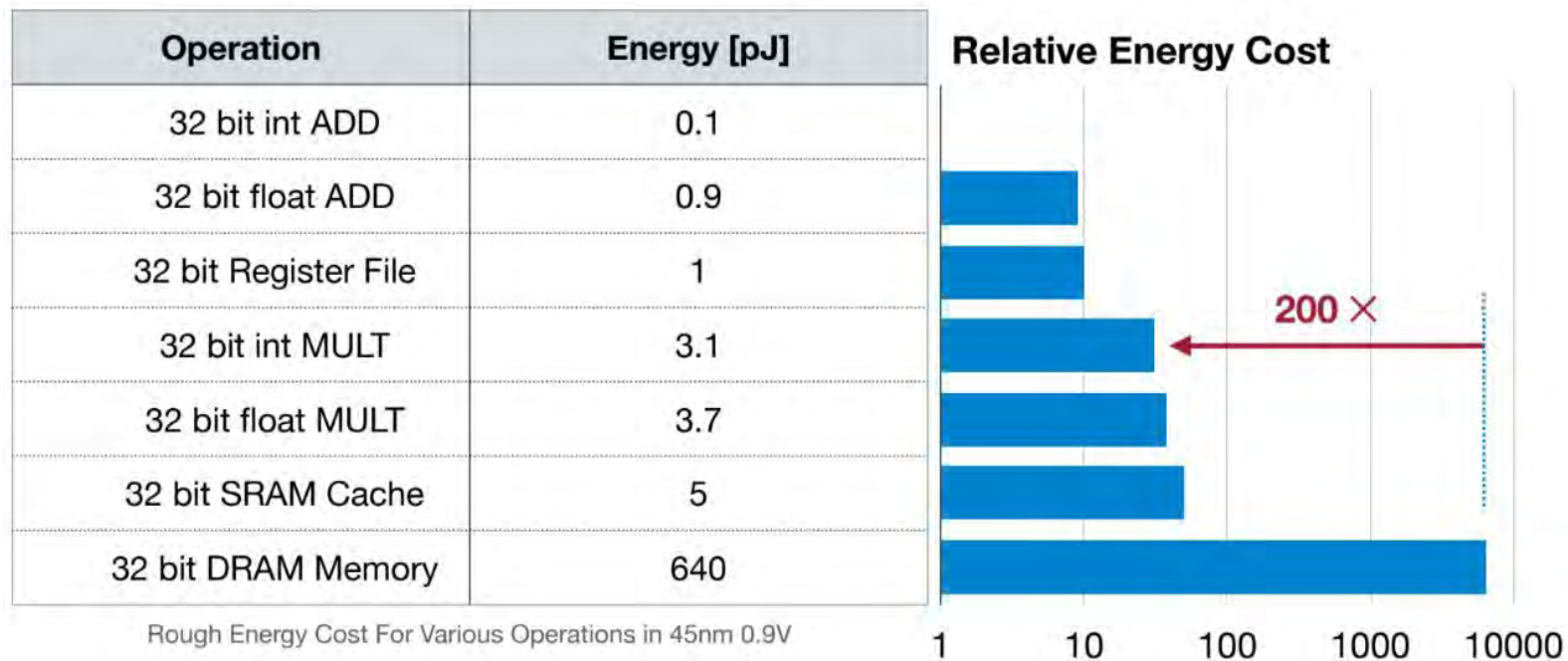
Xavier

AGX
Xavier

Orin



Why pruning? Energy

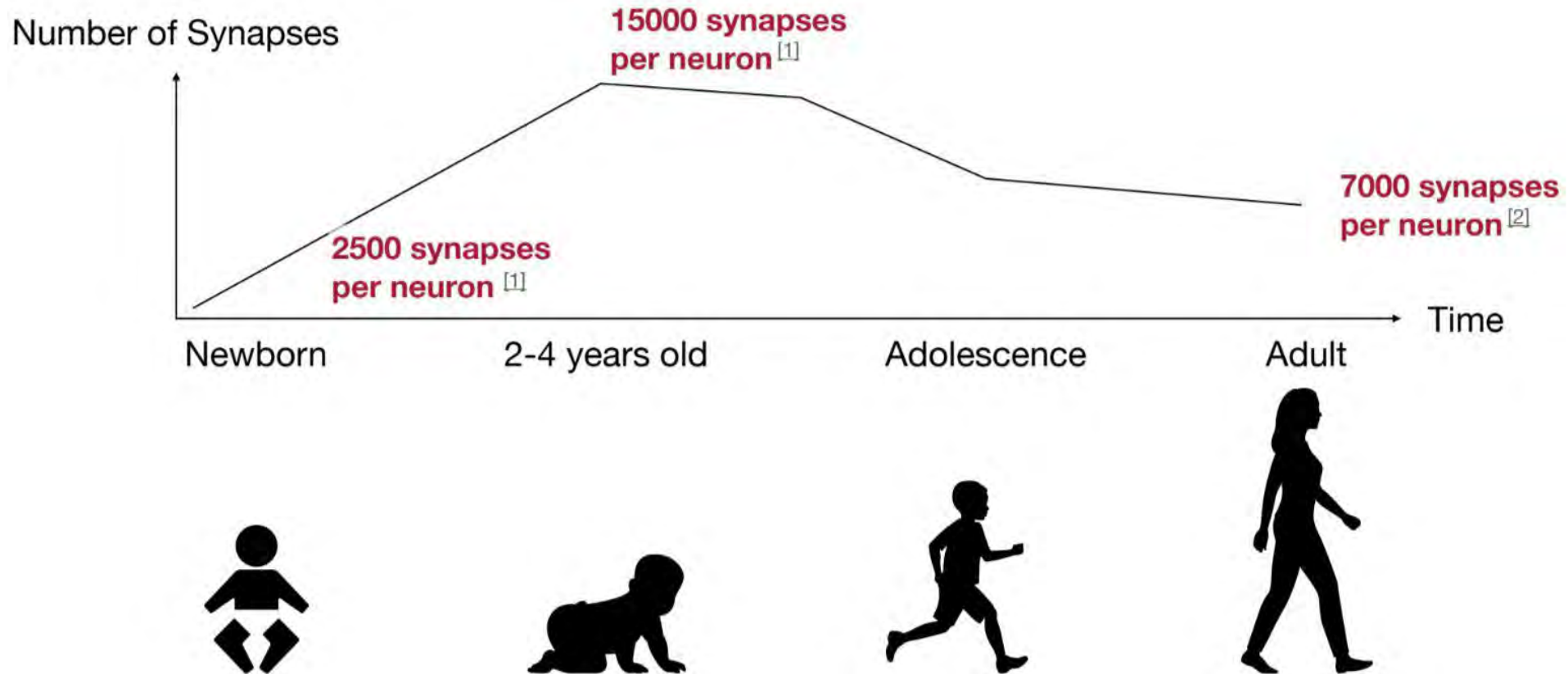


Smaller network, will it work?

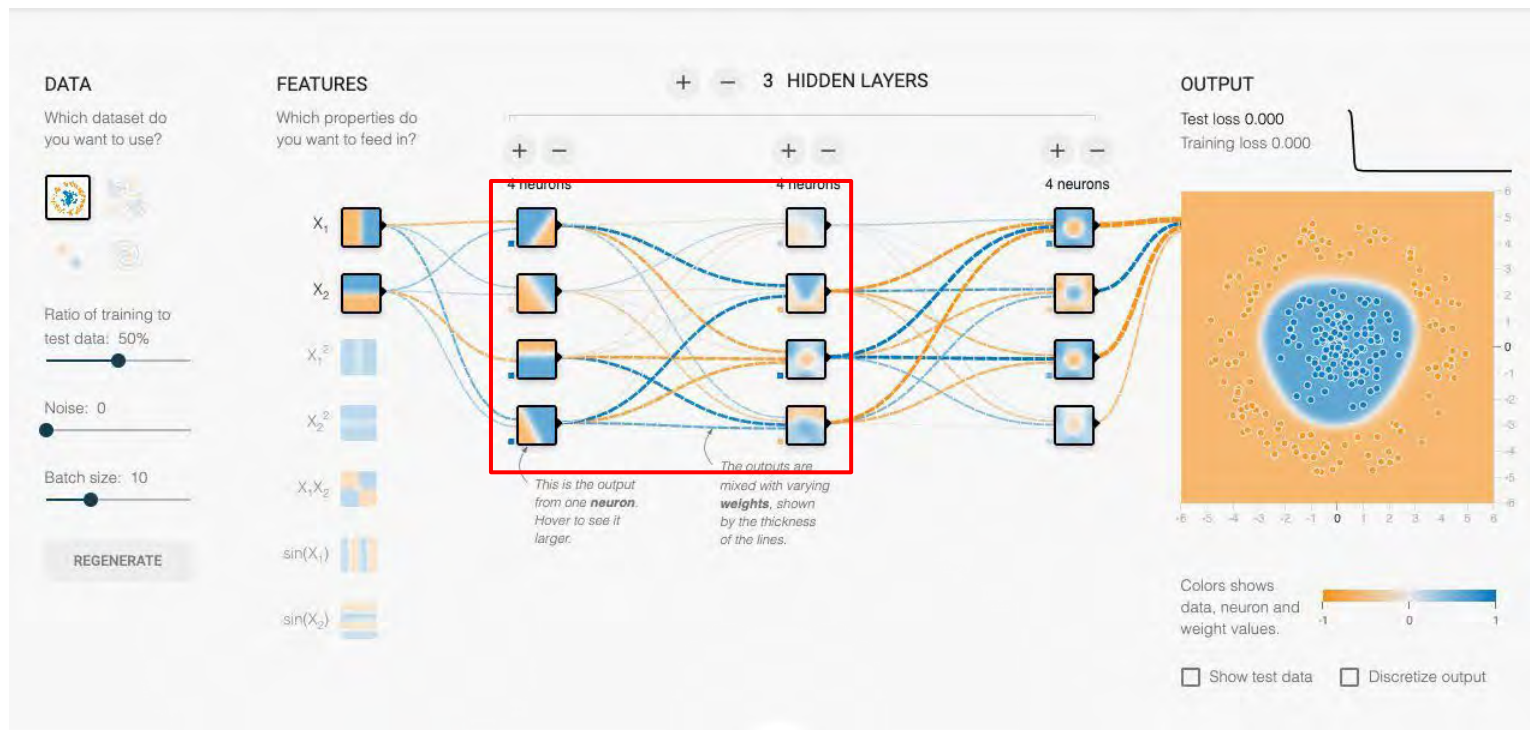
- MAC, Multiply-ACCumulation operations ~ FLOPS

Neural Network	#Parameters			MACs
	Before Pruning	After Pruning	Reduction	Reduction
AlexNet	61 M	6.7 M	9 ×	3 ×
VGG-16	138 M	10.3 M	12 ×	5 ×
GoogleNet	7 M	2.0 M	3.5 ×	5 ×
ResNet50	26 M	7.47 M	3.4 ×	6.3 ×
SqueezeNet	1 M	0.38 M	3.2 ×	3.5 ×

Human brain prunes too!

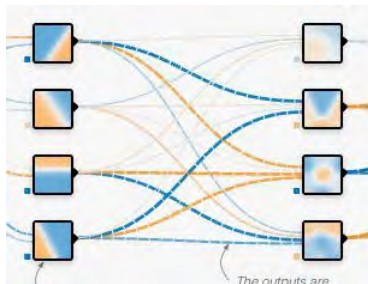


How to prune?

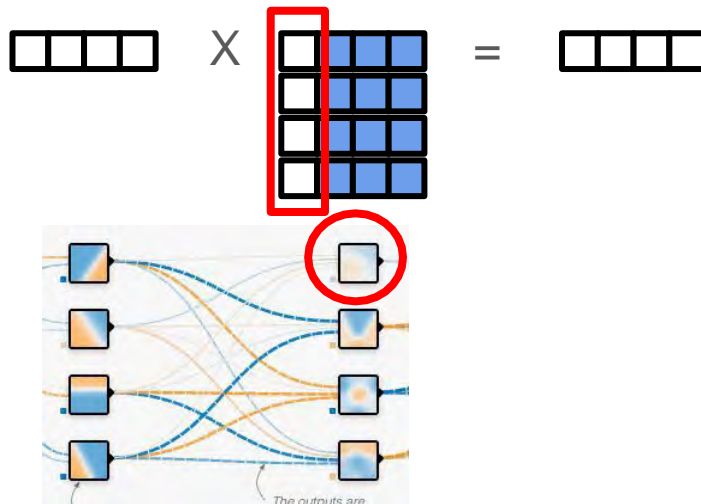


How to prune?

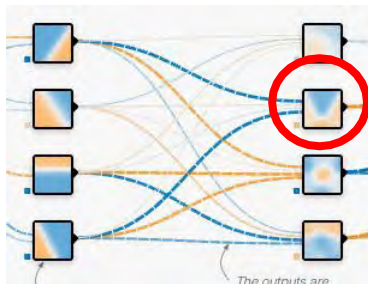
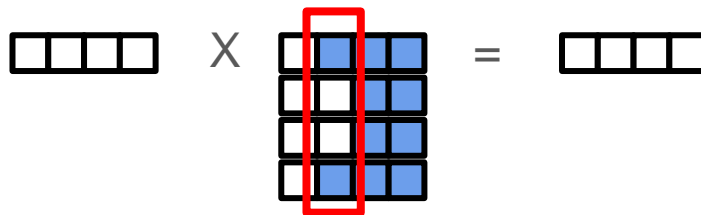
$$\begin{array}{|c|c|c|c|} \hline & & & \\ \hline \end{array} \times \begin{array}{|c|c|c|c|} \hline & & & \\ \hline & & & \\ \hline & & & \\ \hline & & & \\ \hline \end{array} = \begin{array}{|c|c|c|c|} \hline & & & \\ \hline \end{array}$$



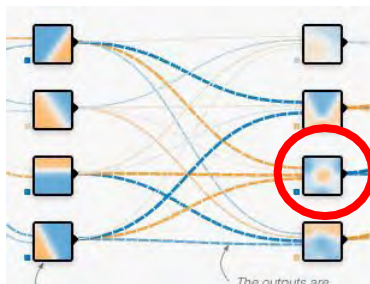
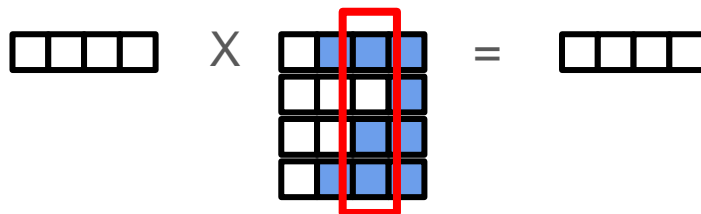
How to prune?



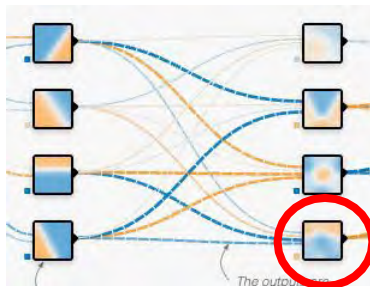
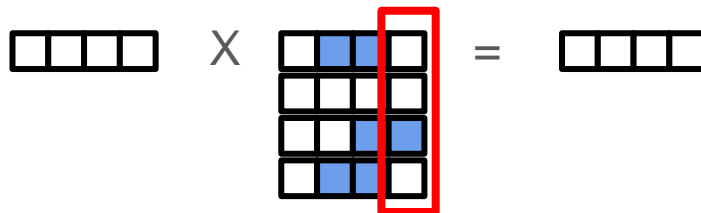
How to prune?



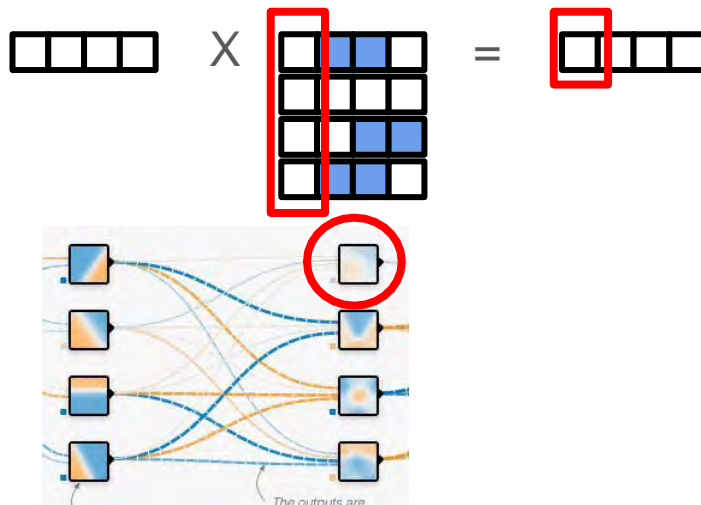
How to prune?



How to prune?

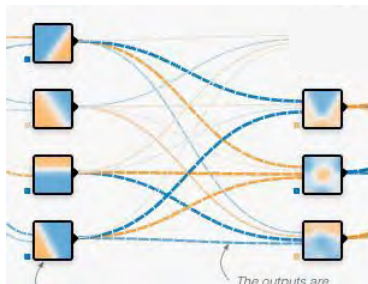


How to prune?

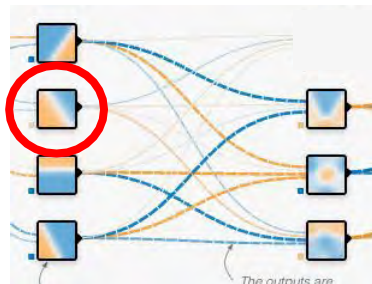
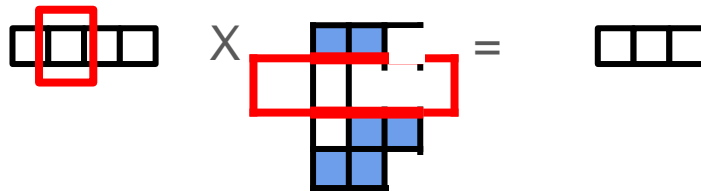


How to prune?

$$\begin{bmatrix} \square & \square & \square & \square \end{bmatrix} \times \begin{bmatrix} \blacksquare & \blacksquare & \square \\ \square & \square & \square \\ \square & \blacksquare & \blacksquare \\ \blacksquare & \blacksquare & \square \end{bmatrix} = \begin{bmatrix} \square & \square & \square \end{bmatrix}$$

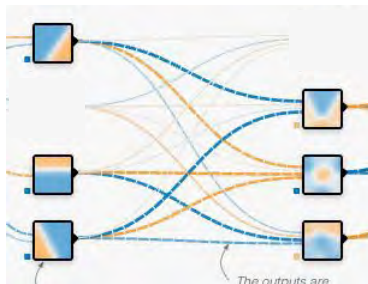


How to prune?

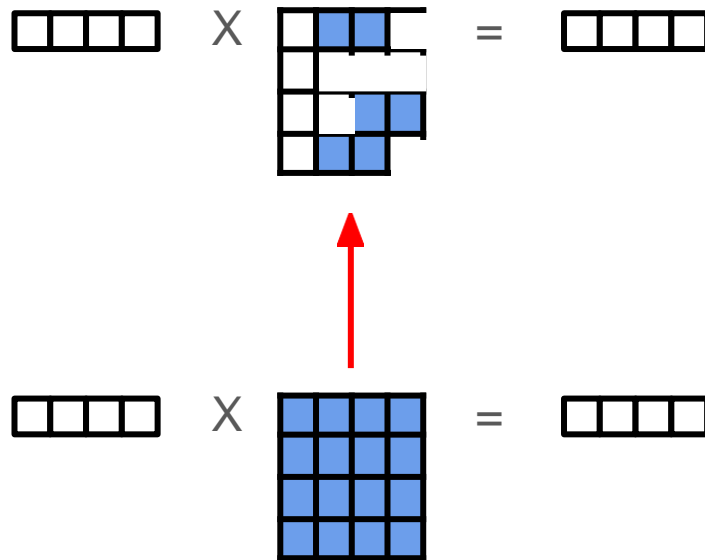


How to prune?

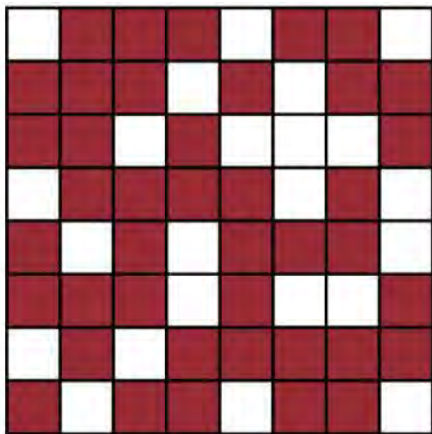
$$\begin{bmatrix} \square & \square & \square \end{bmatrix} \times \begin{bmatrix} \blacksquare & \blacksquare & \square \\ \square & \blacksquare & \blacksquare \end{bmatrix} = \begin{bmatrix} \square & \square & \square \end{bmatrix}$$



How to prune?

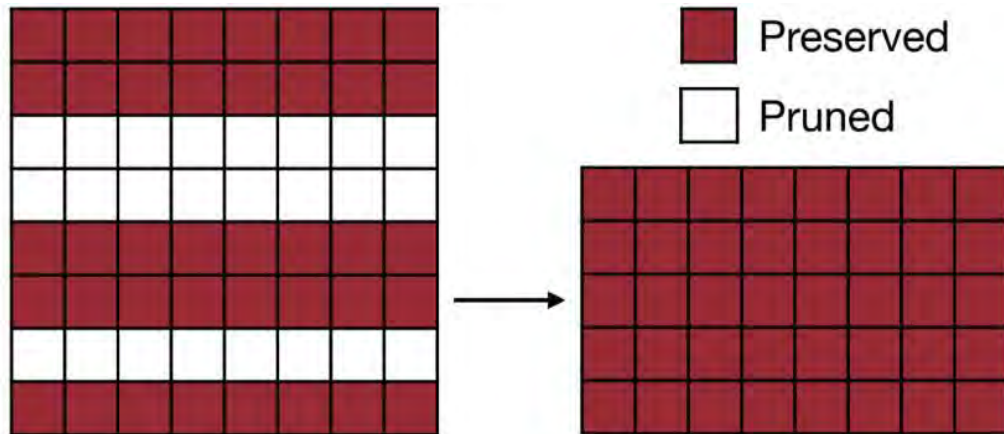


Pruning Granularity



Fine-grained/Unstructured

- More flexible pruning index choice
- Hard to accelerate (irregular)

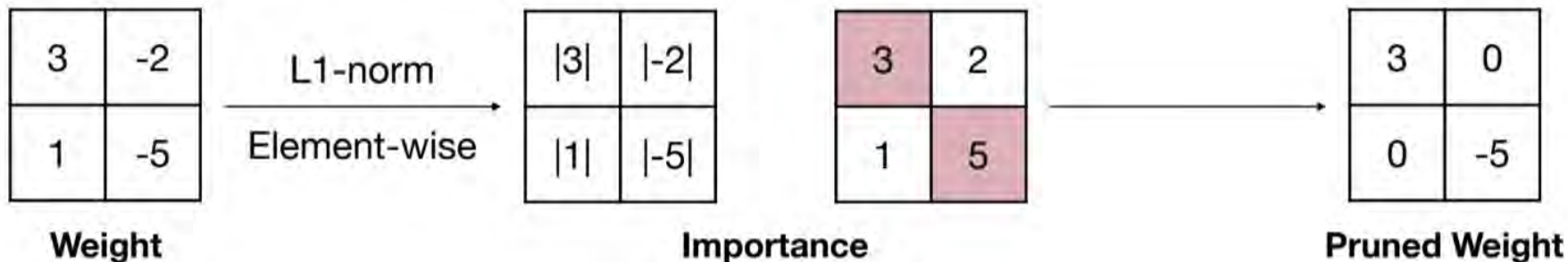


Coarse-grained/Structured

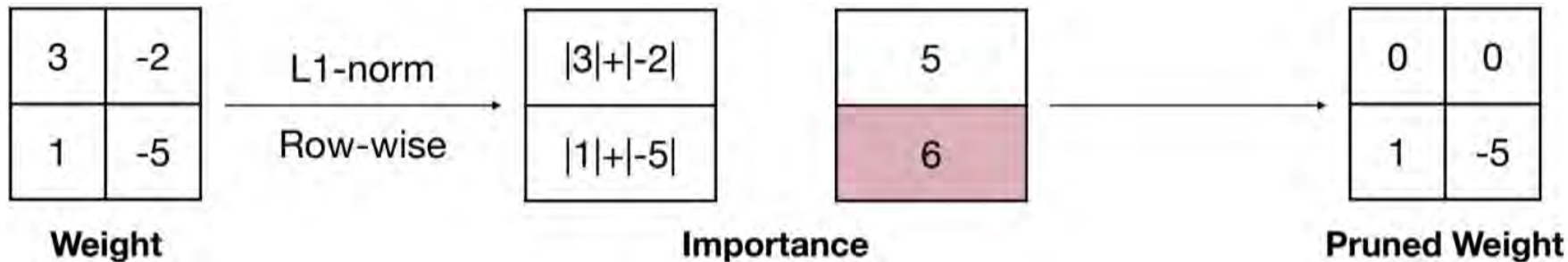
- Less flexible pruning index choice (a subset of the fine-grained case)
- Easy to accelerate (just a smaller matrix!)

Magnitude-based Pruning

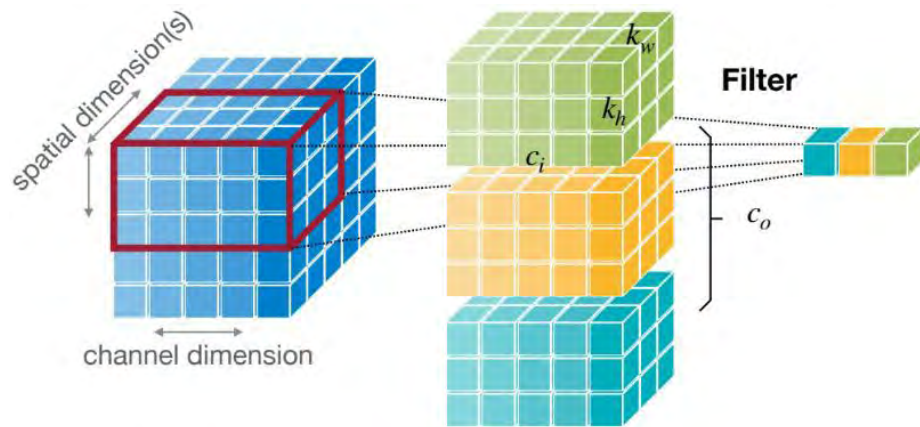
- Fine-grained



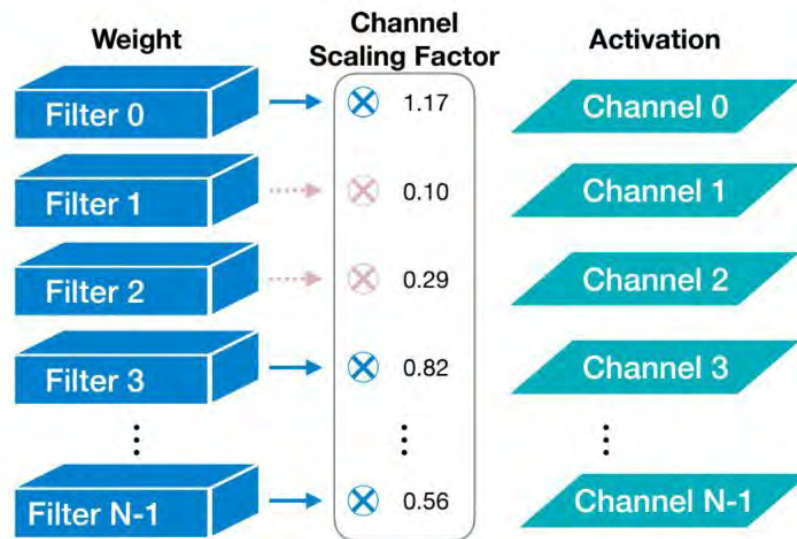
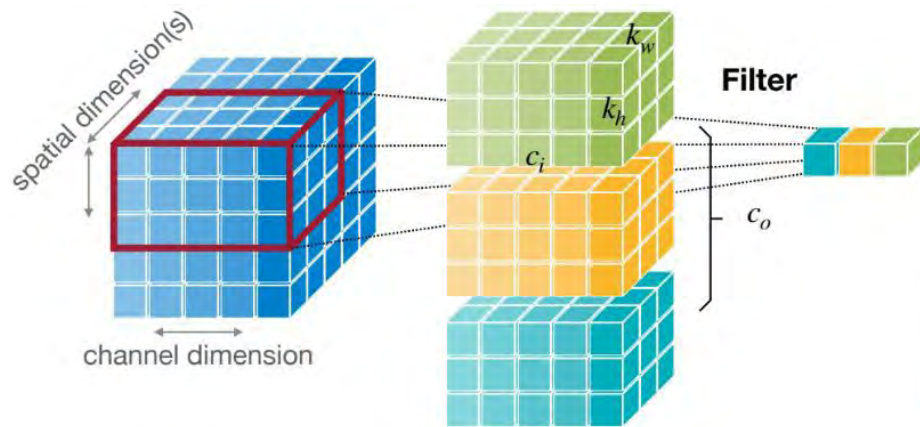
- Coarse-grained / structured



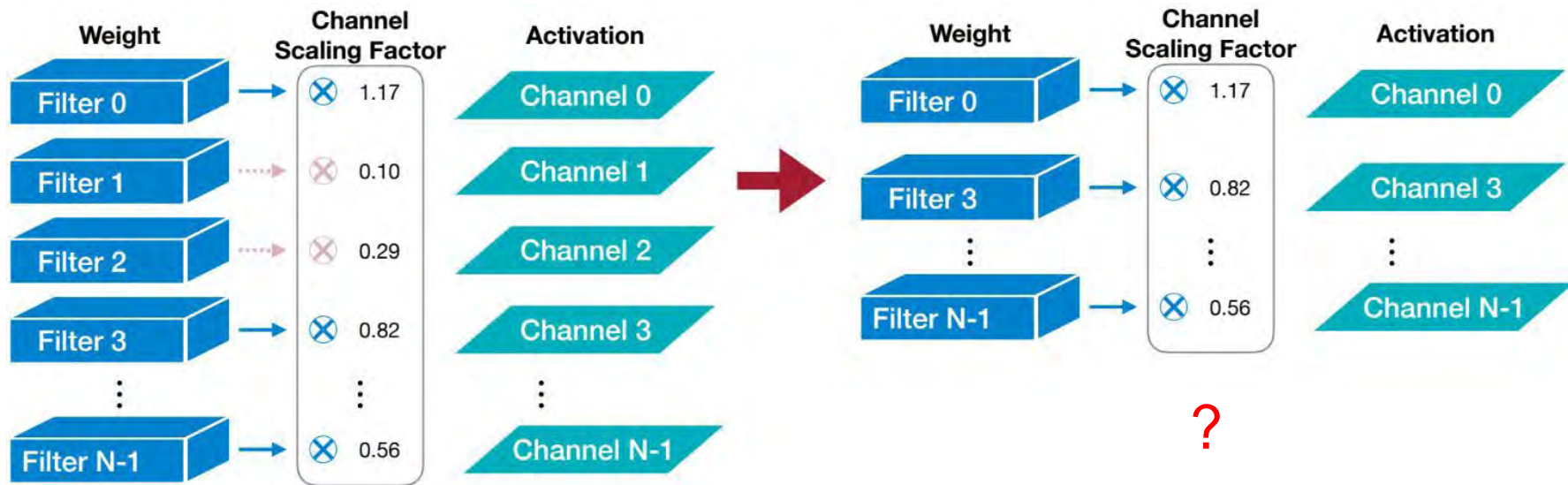
How to prune CNN?



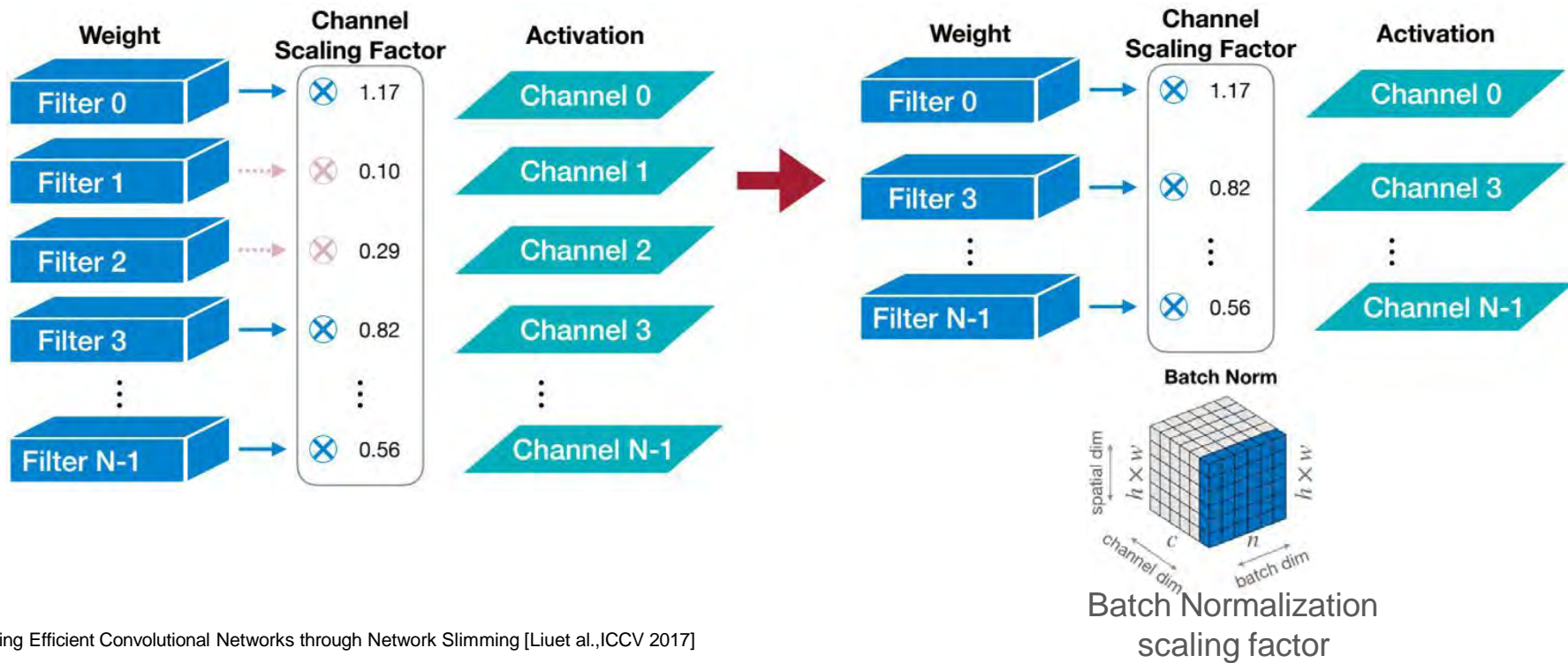
Scaling-based Pruning



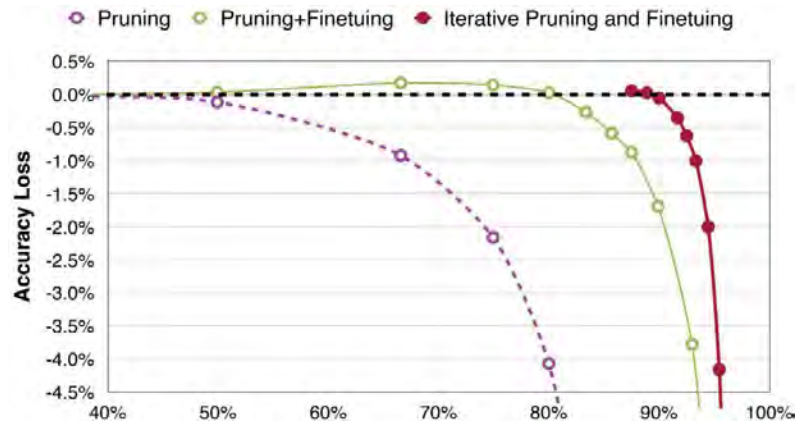
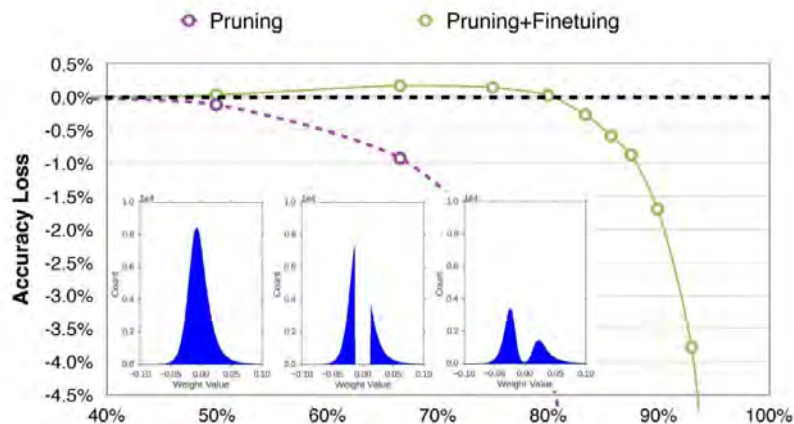
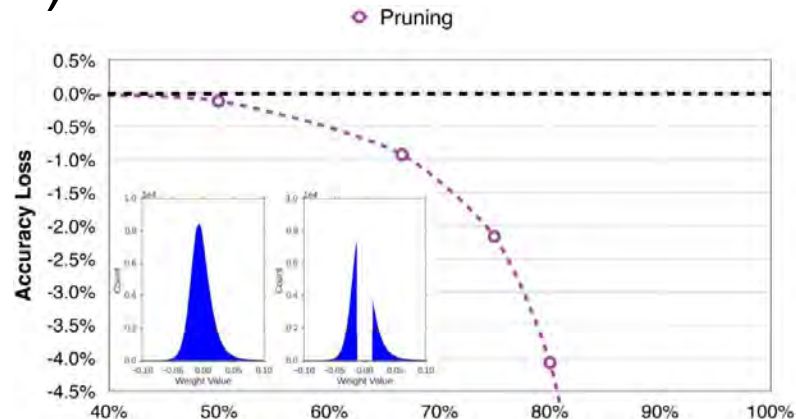
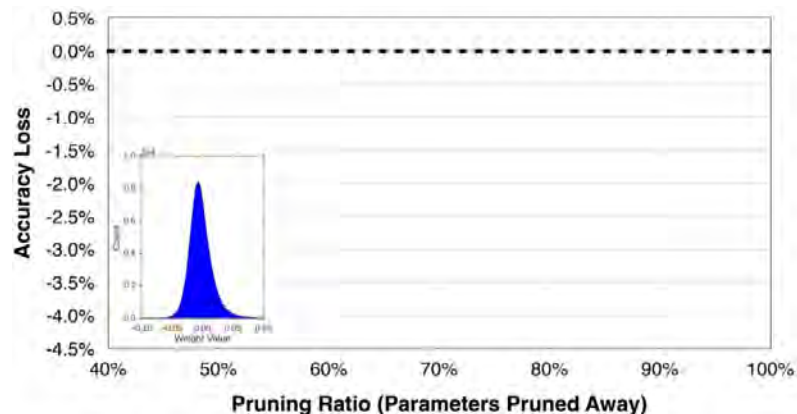
Scaling-based Pruning



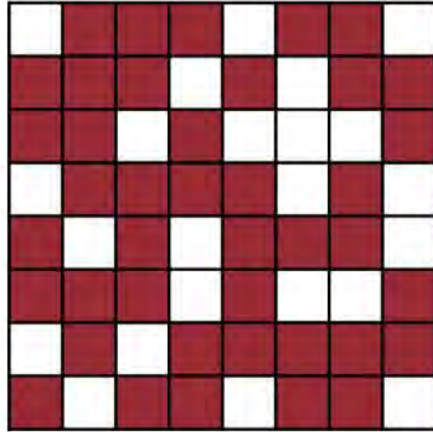
Scaling-based Pruning



Fine-tune/Retrain (e.g. AlexNet)



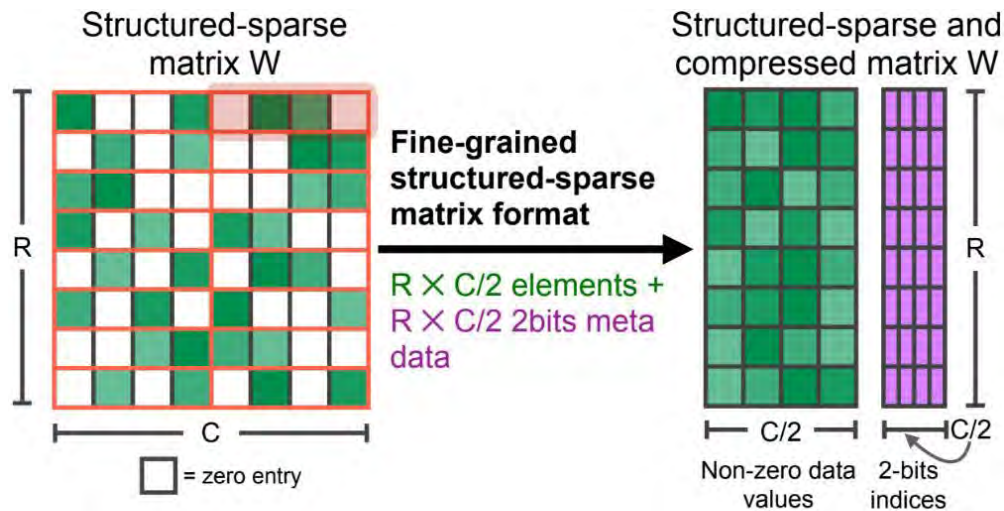
Q: How to accelerate irregular pruning?



Fine-grained/Unstructured

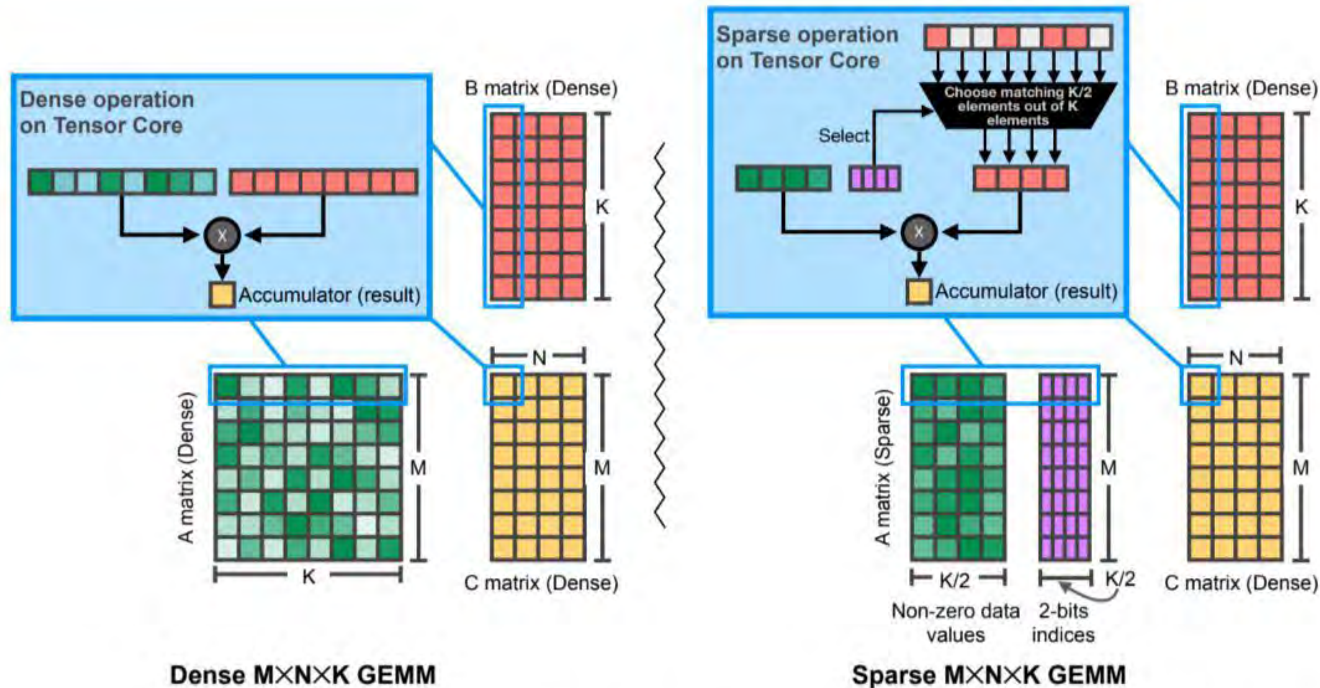
- More flexible pruning index choice
- Hard to accelerate (irregular)

M:N Sparsity



Two weights are nonzero out of four consecutive weights (2:4 sparsity).

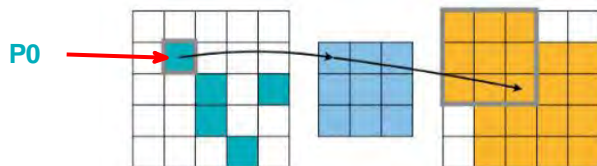
M:N Sparsity



The indices are used to mask out the inputs. Only 2 multiplications will be done out of four.

Sparse Conv

Conventional Convolution



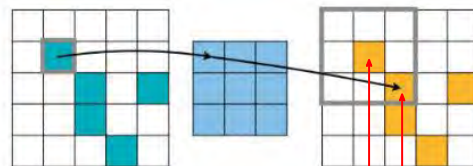
Maps
(In, Out, Wgt)

Computation
($f_{Out} = f_{Out} + f_{In} \times W_{Wgt}$) for
each entry in the maps

(P₀, Q₀, W_{1,1})
 (P₀, Q₁, W_{1,0})
 (P₀, Q₂, W_{1,-1})
 (P₀, Q₃, W_{0,1})
 (P₀, Q₄, W_{0,0})
 (P₀, Q₅, W_{0,-1})
 (P₀, Q₈, W_{-1,1})
 (P₀, Q₉, W_{-1,0})
 (P₀, Q₁₀, W_{-1,-1})

9 matrix multiplications

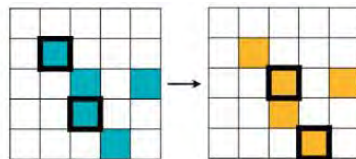
Sparse Convolution



No compute
 No compute
 No compute
 No compute
 (P₀, Q₀, W_{0,0})
 No compute
 No compute
 No compute
 (P₀, Q₁, W_{-1,-1})

2 matrix multiplications

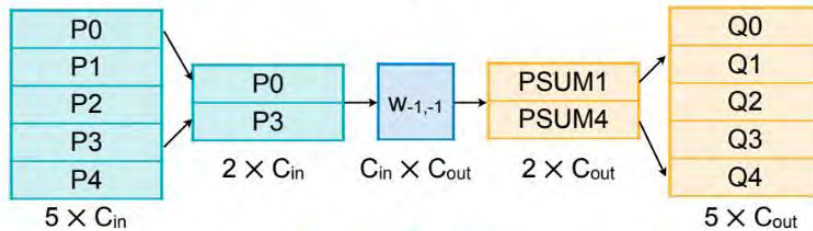
Sparse Conv



Workload

Maps (In, Out, Wgt)		
$(P_0, Q_1, W_{-1,-1})$		
$(P_3, Q_4, W_{-1,-1})$		
$(P_1, Q_3, W_{-1,0})$		
$(P_0, Q_0, W_{0,0})$		
$(P_1, Q_1, W_{0,0})$		
$(P_2, Q_2, W_{0,0})$		
$(P_3, Q_3, W_{0,0})$		
$(P_4, Q_4, W_{0,0})$		
$(P_3, Q_1, W_{1,0})$		
$(P_1, Q_0, W_{1,1})$		
$(P_4, Q_3, W_{1,1})$		

Input Features Input Buffer Weight Partial Sum Output Features



$$f_1 = f_1 + f_0 \times W_{-1,-1}$$

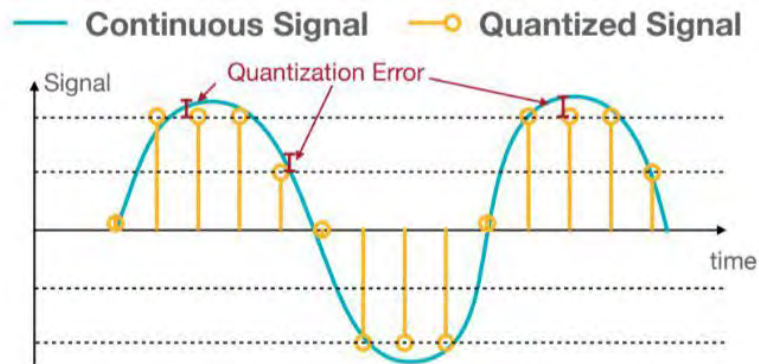
$$f_4 = f_4 + f_3 \times W_{-1,-1}$$

Agenda

- Pruning
 - What and why?
 - How? Pruning criterion
 - Fine-tune/retrain acceleration
- Quantization
 - What and why?
 - Linear quantization
 - Quantization granularity
 - Calibration and Clipping

What is quantization?

- Quantization is the process of constraining an input from a continuous or otherwise large set of values to a discrete set.



The difference between an input value and its quantized value is referred to as quantization error.

Original Image

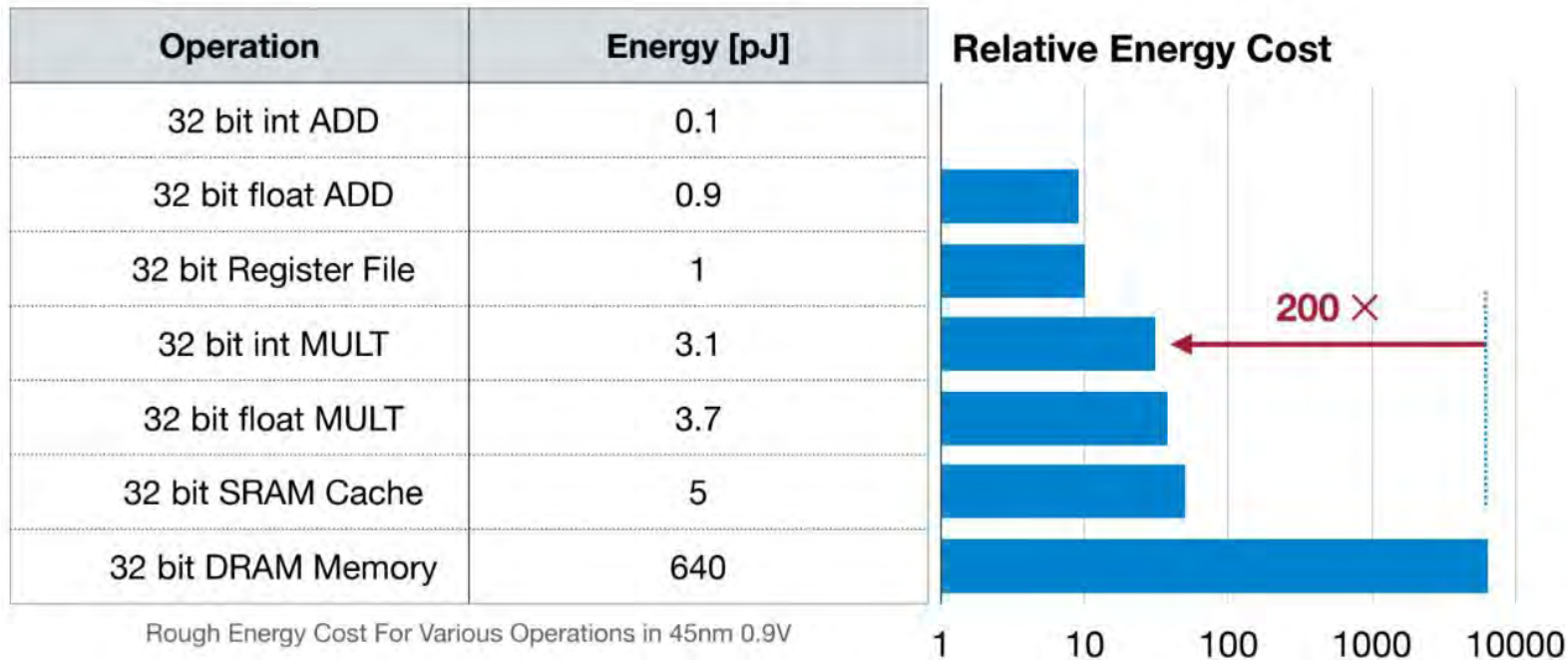


16-Color Image



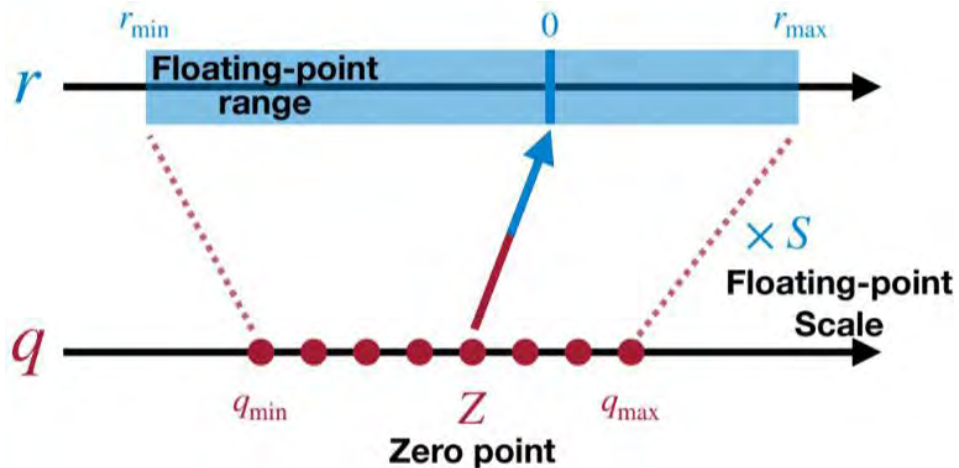
Images are in the public domain.
“Palettization”

Why quantization? Small footprint & low energy



Linear Quantization

- An affine mapping of integers to real numbers
 - $R = S(q - Z)$



$$S = \frac{r_{\max} - r_{\min}}{q_{\max} - q_{\min}}$$

$$Z = \text{round} \left(q_{\min} - \frac{r_{\min}}{S} \right)$$

$$r_{\max} = S (q_{\max} - Z)$$

$$r_{\min} = S (q_{\min} - Z)$$

Linear Quantization: Example

2.09	-0.98	1.48	0.09
0.05	-0.14	-1.08	2.12
-0.91	1.92	0	-1.03
1.87	0	1.53	1.49

Bit Width	q_{\min}	q_{\max}
2	-2	1
3	-4	3
4	-8	7
N	-2^{N-1}	$2^{N-1}-1$

$$r_{\max} = S (q_{\max} - Z)$$

$$r_{\min} = S (q_{\min} - Z)$$

$$S = \frac{r_{\max} - r_{\min}}{q_{\max} - q_{\min}}$$

$$Z = \text{round} \left(q_{\min} - \frac{r_{\min}}{S} \right)$$

Linear Quantization: Example

2.09	-0.98	1.48	0.09
0.05	-0.14	-1.08	2.12
-0.91	1.92	0	-1.03
1.87	0	1.53	1.49

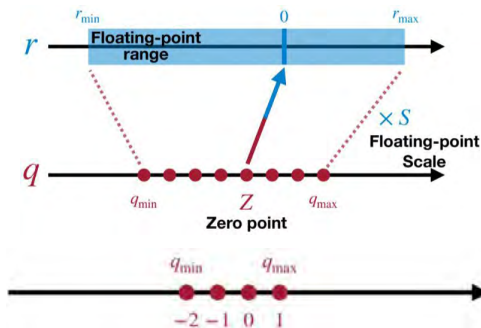
$$S = (2.12 - (-1.08)) / (1 - (-2)) \\ = 1.07$$

$$Z = \text{round}(-2 - -1.08 / 1.07) \\ = -1$$

$$r_{\max} = S (q_{\max} - Z)$$

$$r_{\min} = S (q_{\min} - Z)$$

Bit Width	q_{\min}	q_{\max}
2	-2	1
3	-4	3
4	-8	7
N	-2^{N-1}	$2^{N-1}-1$



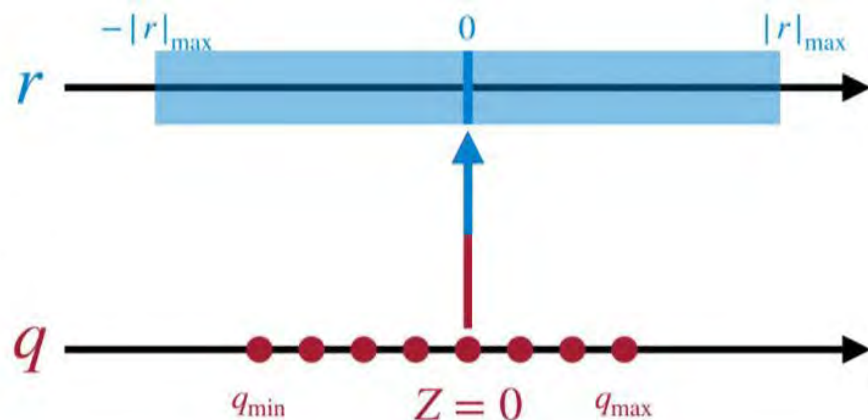
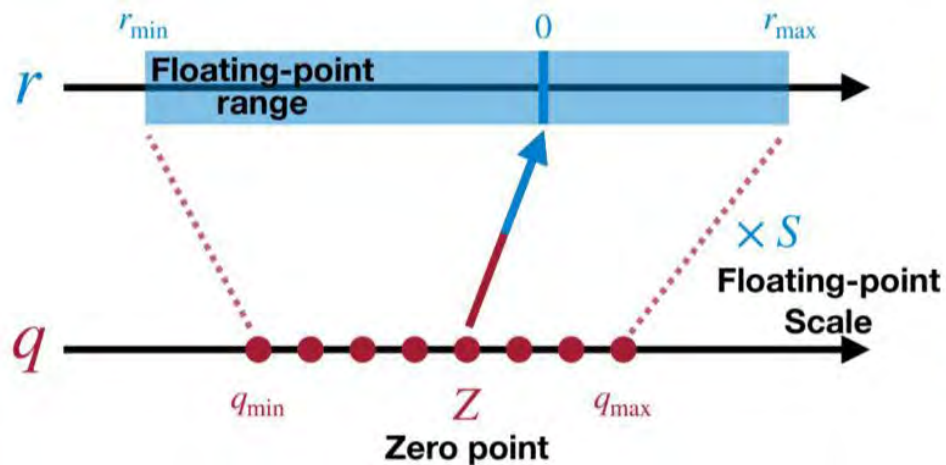
$$S = \frac{r_{\max} - r_{\min}}{q_{\max} - q_{\min}}$$

$$Z = \text{round} \left(q_{\min} - \frac{r_{\min}}{S} \right)$$

Linear Quantization

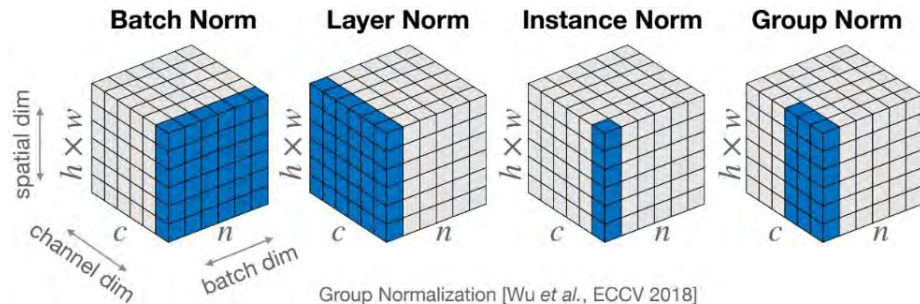


Symmetric Linear Quantization

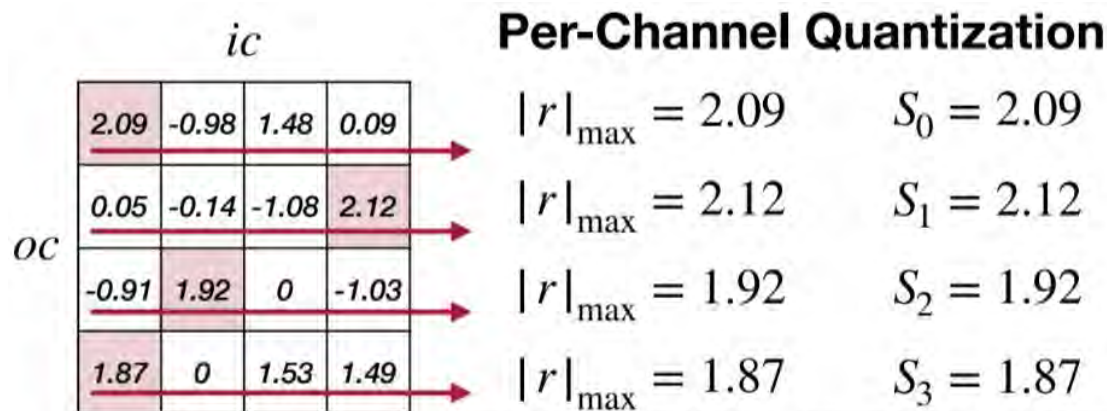


Quantization Granularity

- Per-tensor quantization
- Per-channel quantization
- Group quantization



Per-channel vs Per-tensor



1	0	1	0
0	0	-1	1
0	1	0	-1
1	0	1	1

Quantized

2.09	0	2.09	0
0	0	-2.12	2.12
0	1.92	0	-1.92
1.87	0	1.87	1.87

Reconstructed

$$\|W - S \odot q_W\|_F = 2.08$$

Per-Tensor Quantization

$$|r|_{\max} = 2.12$$

$$S = \frac{|r|_{\max}}{q_{\max}} = \frac{2.12}{2^{2-1} - 1} = 2.12$$

1	0	1	0
0	0	-1	1
0	1	0	0
1	0	1	1

Quantized

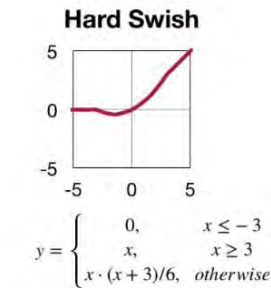
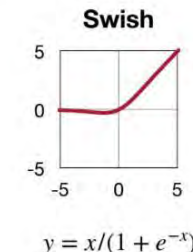
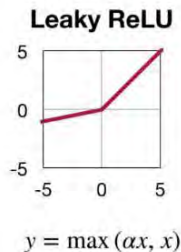
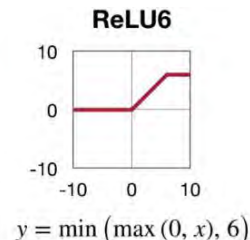
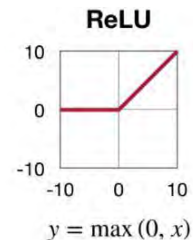
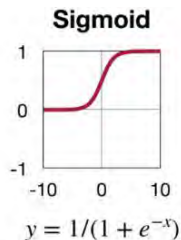
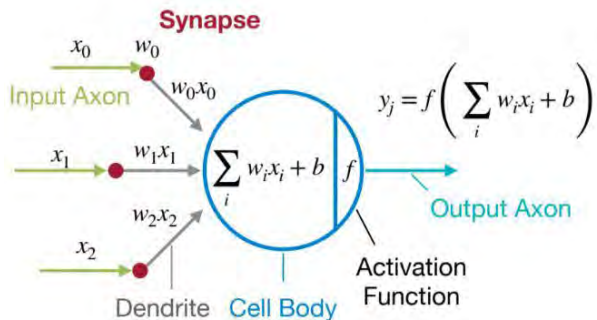
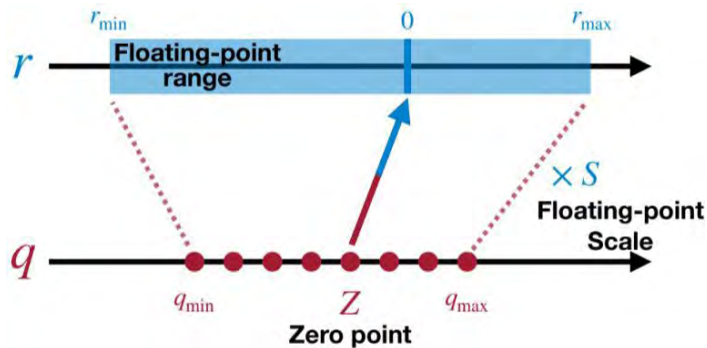
2.12	0	2.12	0
0	0	-2.12	2.12
0	2.12	0	0
2.12	0	2.12	2.12

Reconstructed

$$\|W - Sq_W\|_F = 2.28$$

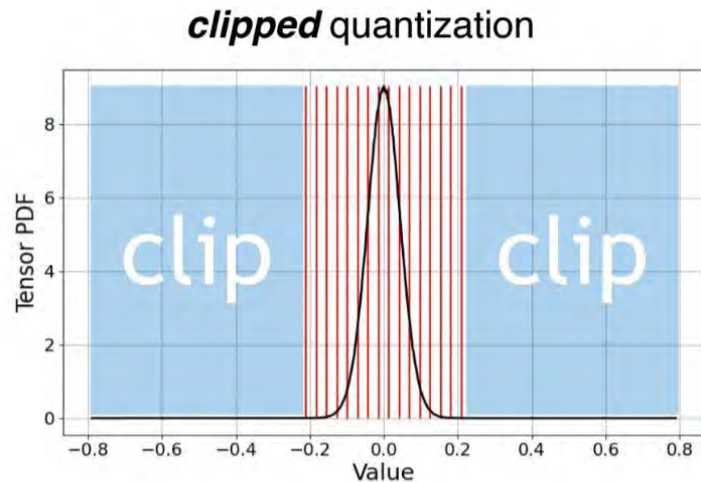
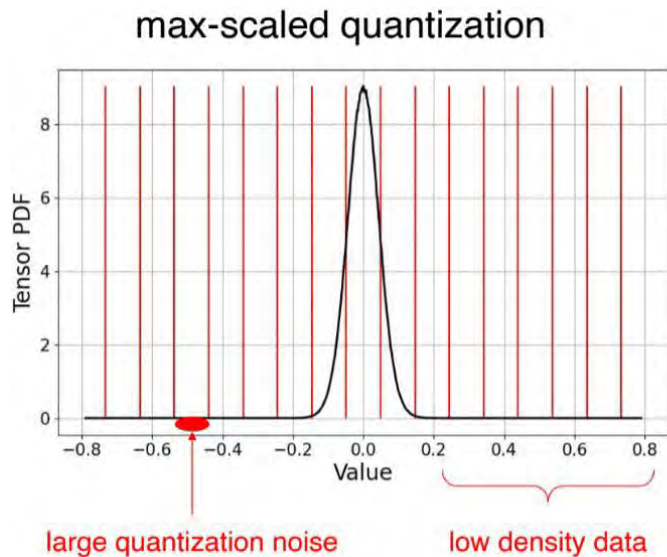
Non-static Dynamic Range

- Weights has $[r_min, r_max]$ fixed range
- Activation value has unknown range
 - depends on input
- How to determine r_min and r_max ?



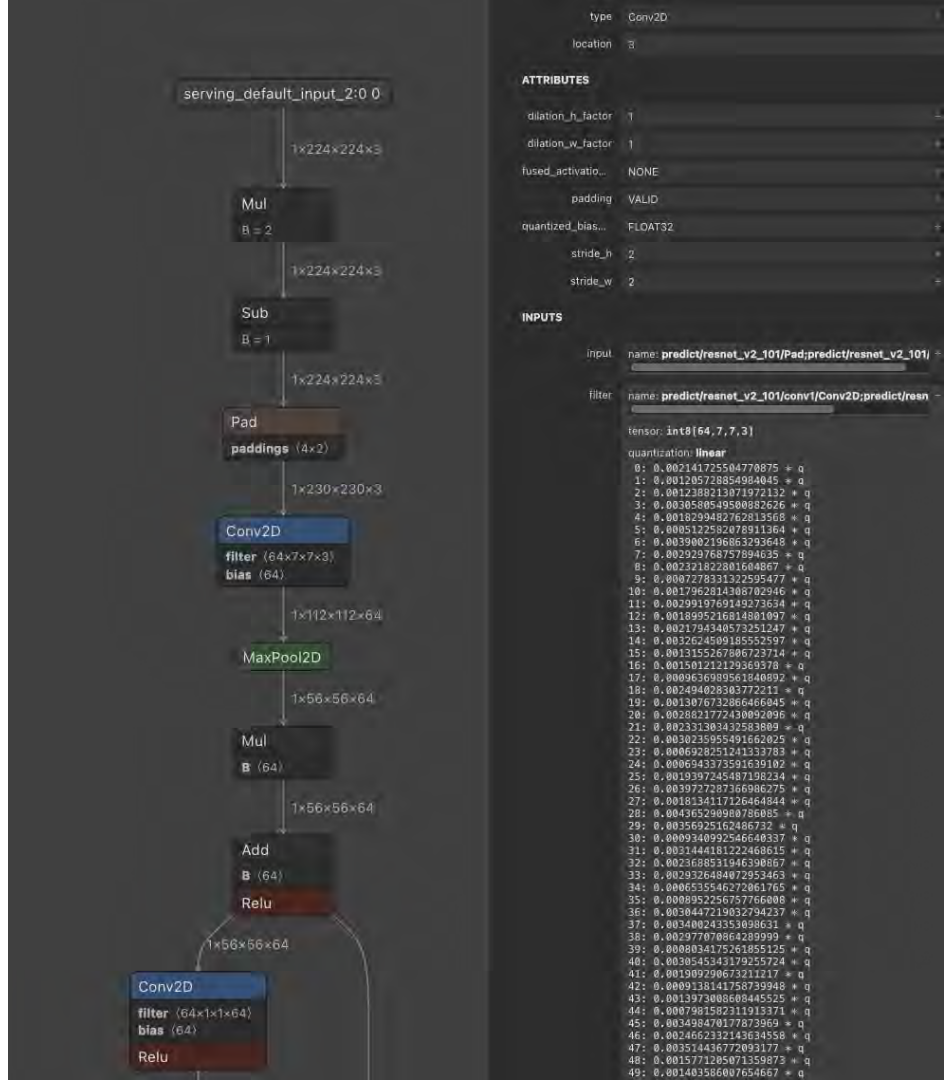
Dynamic Range: Calibration & Clipping

- Run a batch of samples
- Get the statistical distribution of activations
- Get rid of outliers



Tutorial: TF Lite quantization

- [Post-training dynamic range quantization | TensorFlow Lite](#)
- Visualize quantized vs unquantized model in Netron.app



Summary

- Pruning

- What and why?
- How? Pruning criterion
 - Magnitude-based, Scaling-based
- Fine-tune/retrain acceleration
 - M:N sparsity, SparseConv

- Quantization

- What and why?
- Linear quantization
 - Scale, zero point
- Quantization granularity
 - Per-tensor, per-channel
- Calibration and Clipping

Next Lectures

- Example paper presentation: ML-Exray
- Hardware architecture
- Special accelerators
- Quiz on Lecture 05-06