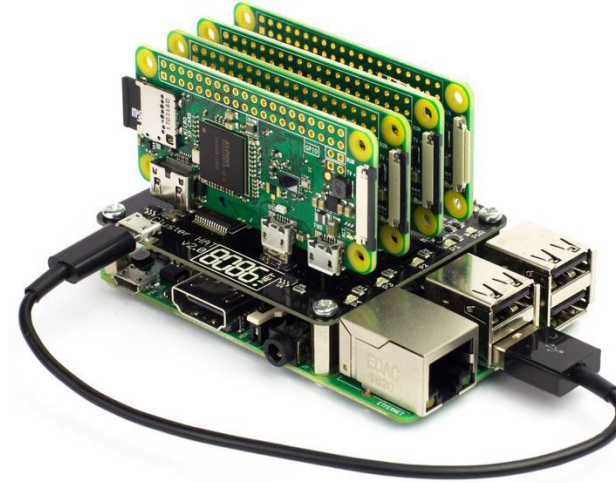
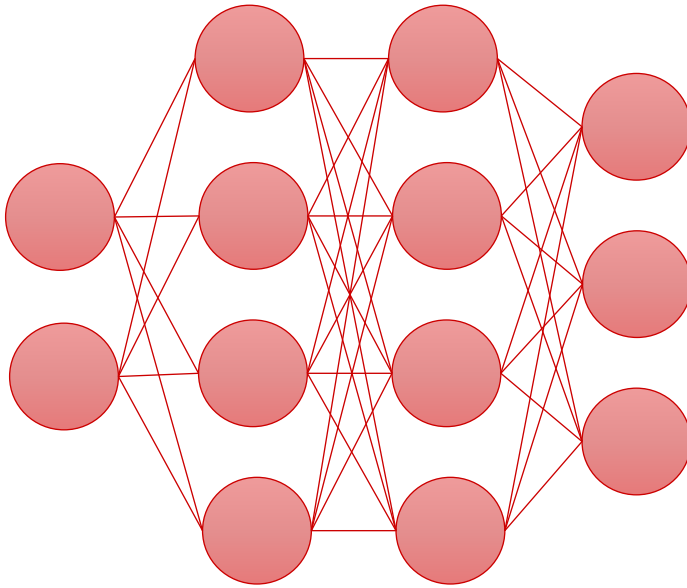


INF2009 – Edge Computing and Analytics [2024/25 T2]

Edge Analytics – Deep Learning on the Edge

Deep Learning on the Edge

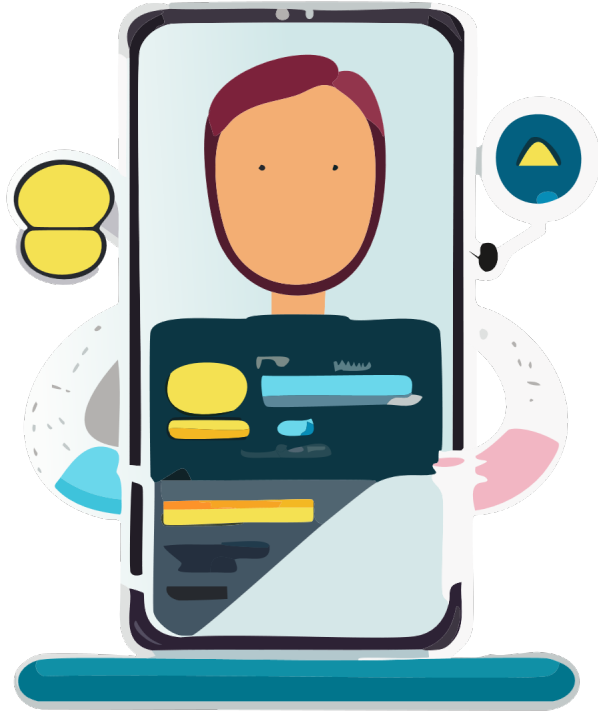


[This Photo](#) by Unknown Author is licensed under [CC BY-SA](#)

Reference:

- [TinyML and Efficient Deep Learning Computing \(MIT\)](#)
- [Machine Learning Hardware and Systems \(Cornell Tech\)](#)
- [tinyML Talks: A Practical Guide to Neural Network Quantization](#)

Contents



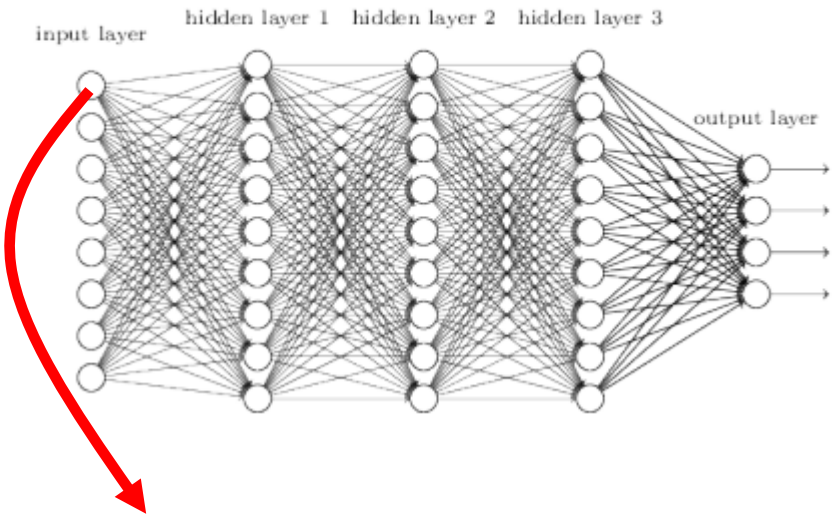
- *Strategies for DL on edge*
- *Pruning*
- *Quantization*
- *Weight Sharing*
- *Advanced Approaches*

Deep Learning on the Edge

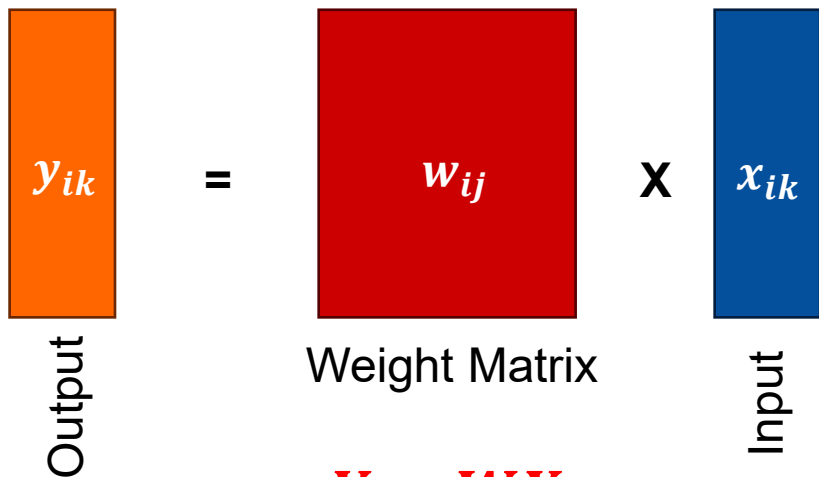


- *Model Size*
 - *Over-the-air update is important*
- *Speed*
 - *Model's inference speed is critical*
- *Energy Efficiency*
 - *Should consume lower battery*

What's Inside a DL and how to make it run on edge?



$$y_i = f\left(\sum_j w_{ij} x_i\right)$$



$$Y = WX$$

*ignoring bias

- Could reduce number of weights (edges) → Pruning
 - Can be done within or across neurons (layers)
- Could reduce the bit size of the weights → Quantization
- Could share weights post quantization → Weight sharing

Pruning (in Humans)

50 Trillion Synapses



Newborn

1000 Trillion Synapses



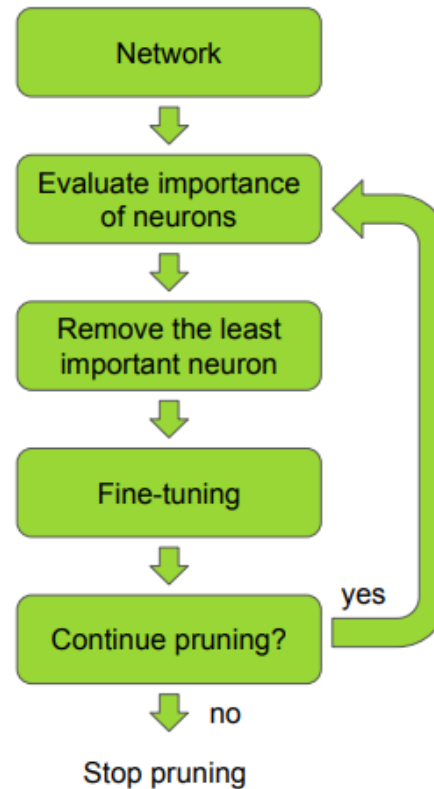
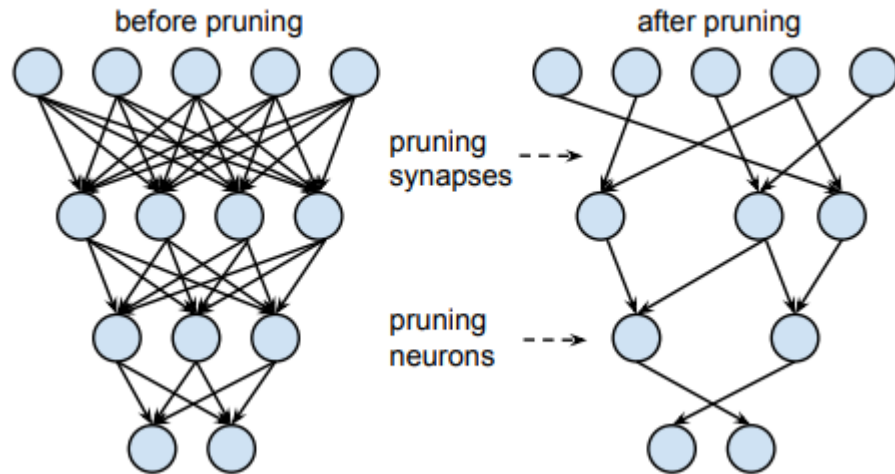
Child

500 Trillion Synapses



Adulthood

Pruning – How?



- Consider a set of training examples
 $D = \langle X = \{x_0, x_1, \dots, x_N\} \mid Y = \{y_0, y_1, \dots, y_N\} \rangle$
- The network's parameters
 $W = \{(w_1^1, b_1^1), (w_2^2, b_2^2), \dots, (w_L^{C_L}, b_L^{C_L})\}$ are optimized to minimize a cost value $C(D|W)$.
- The most common choice for a cost function $C(\cdot)$ is a **negative log-likelihood function**.

$$\min_{W'} |C(D|W') - C(D|W)| \quad \text{s.t. } \|W'\|_{l=0,1,2} \leq B$$

Pruning - Results

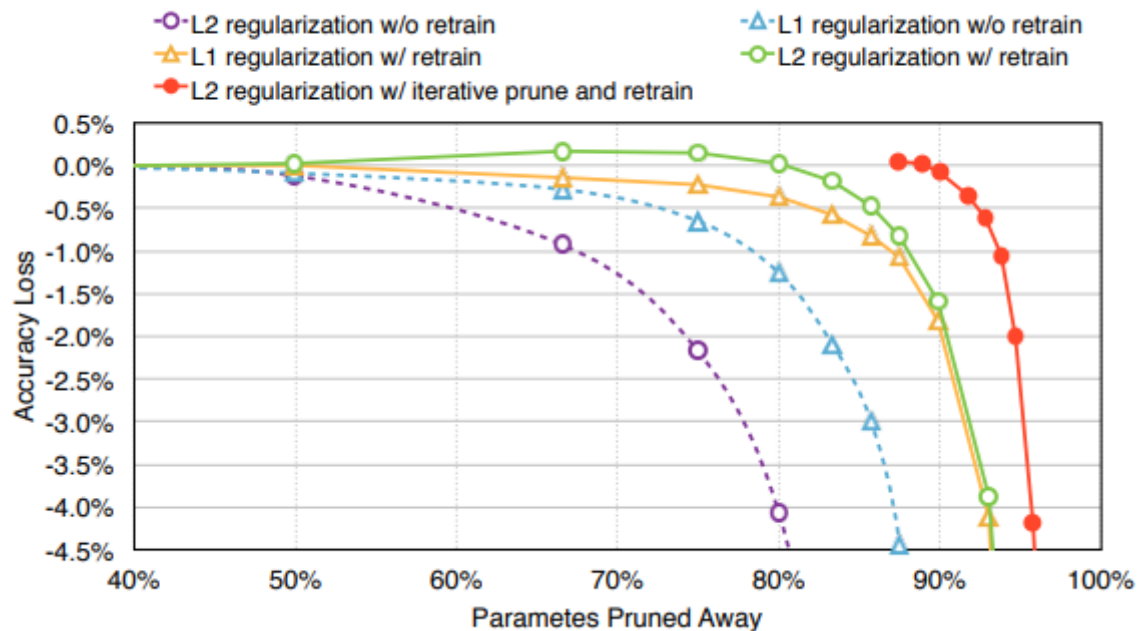
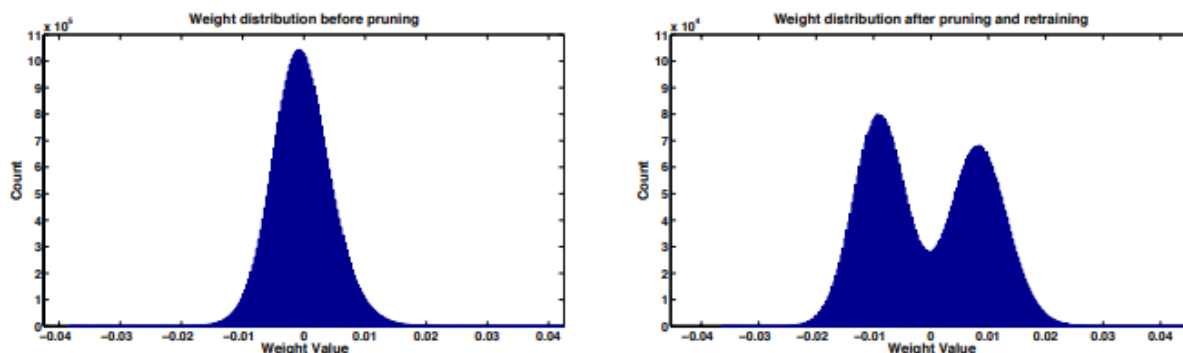


Figure 7: Generate sentence with pretrained model and 90% pruned model
 Pretrained: a brown dog is running through a grassy **field**
 Pruned: a brown dog is running through a grassy **area**



Figure 13: Generate sentence with pretrained model and 95% pruned model
 Pretrained: a soccer player in red is running in the field
 Pruned: a man in a red shirt **and black and white black shirt** is running through a field



Pruning - Results

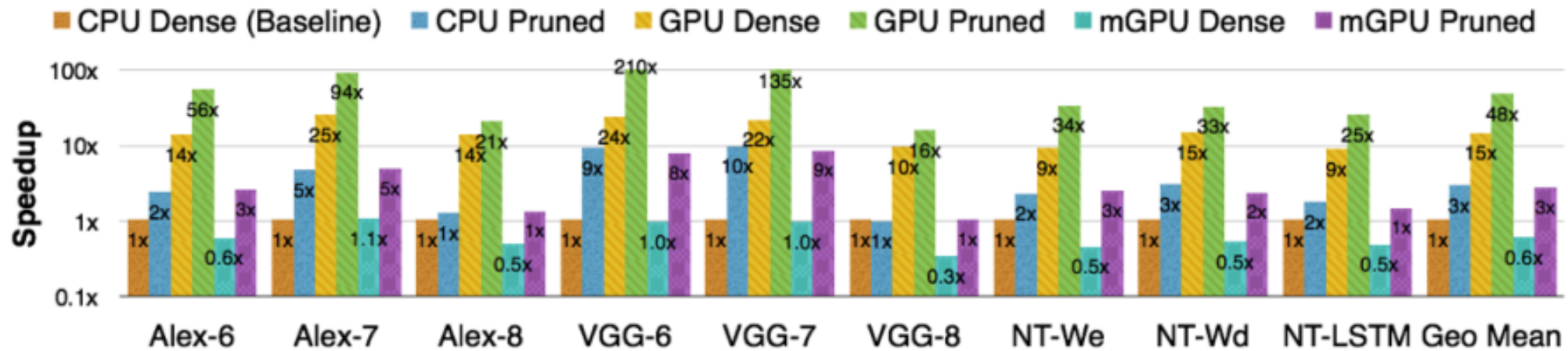


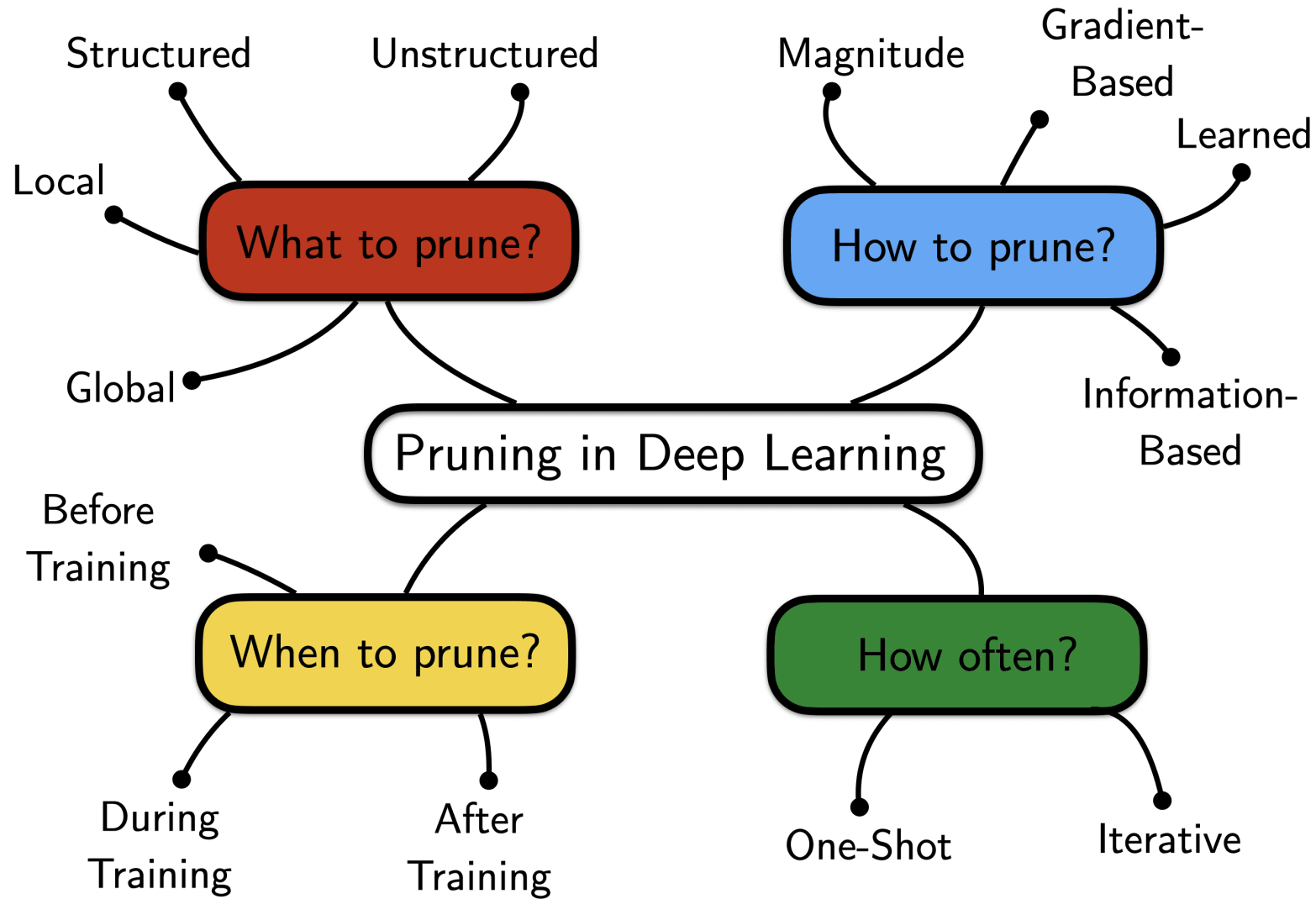
Figure 9: Compared with the original network, pruned network layer achieved $3\times$ speedup on CPU, $3.5\times$ on GPU and $4.2\times$ on mobile GPU on average. Batch size = 1 targeting real time processing. Performance number normalized to CPU.

Intel Core i7 5930K: MKL CBLAS GEMV, MKL SPBLAS CSRMMV

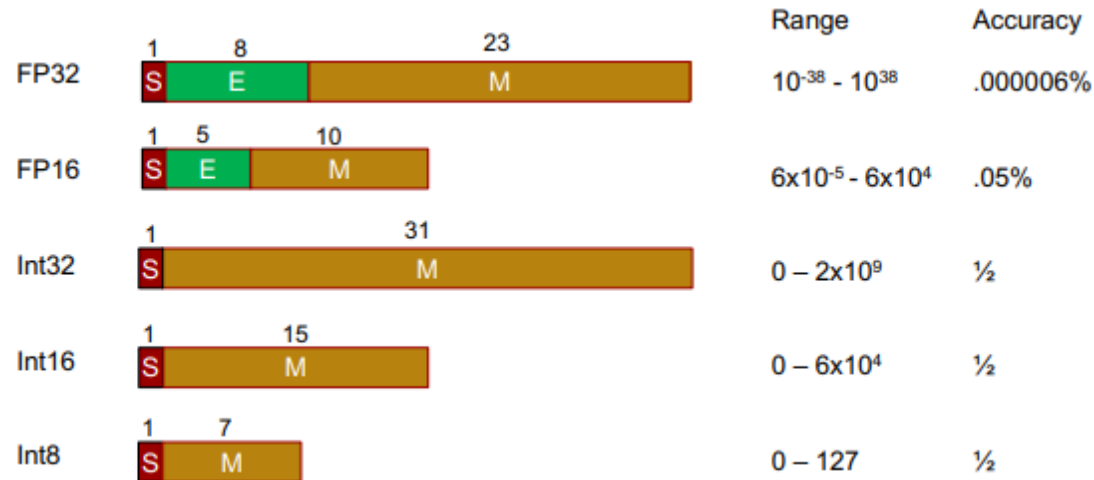
NVIDIA GeForce GTX Titan X: cuBLAS GEMV, cuSPARSE CSRMMV

NVIDIA Tegra K1: cuBLAS GEMV, cuSPARSE CSRMMV

What, When and How often to prune?



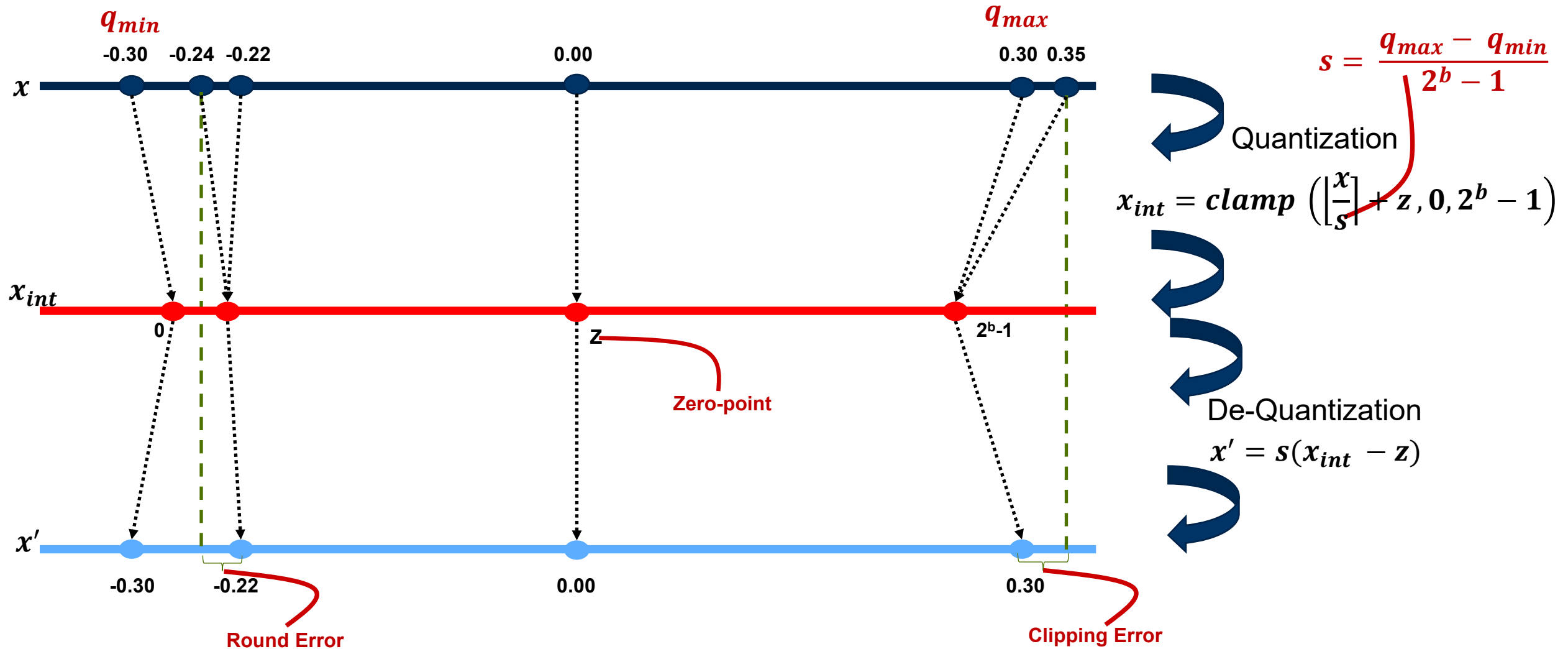
Quantization - Significance



Operation	Precision	Energy (pJ)
Addition	INT8	0.03
	INT16	0.05
	INT32	0.1
	FP16	0.4
	FP32	0.9
Multiplication	INT8	0.2
	INT32	3.1
	FP16	1.1
	FP32	3.7

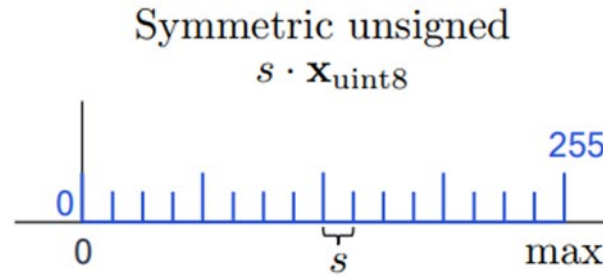
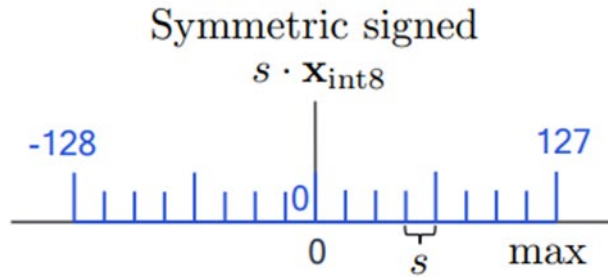
Significant savings in energy when moving from FP32 (typical NNs) to INT8

Quantization – How?

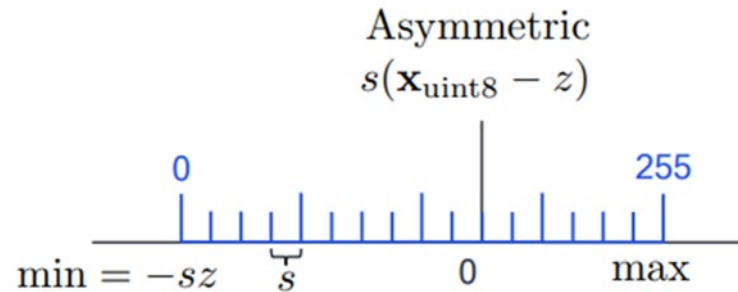


Types of Quantization

Symmetric



Asymmetric



$$Y' = (S_W S_X W_{int}) X_{int}$$

$$Y' = S_W W_{int} S_X X_{int}$$

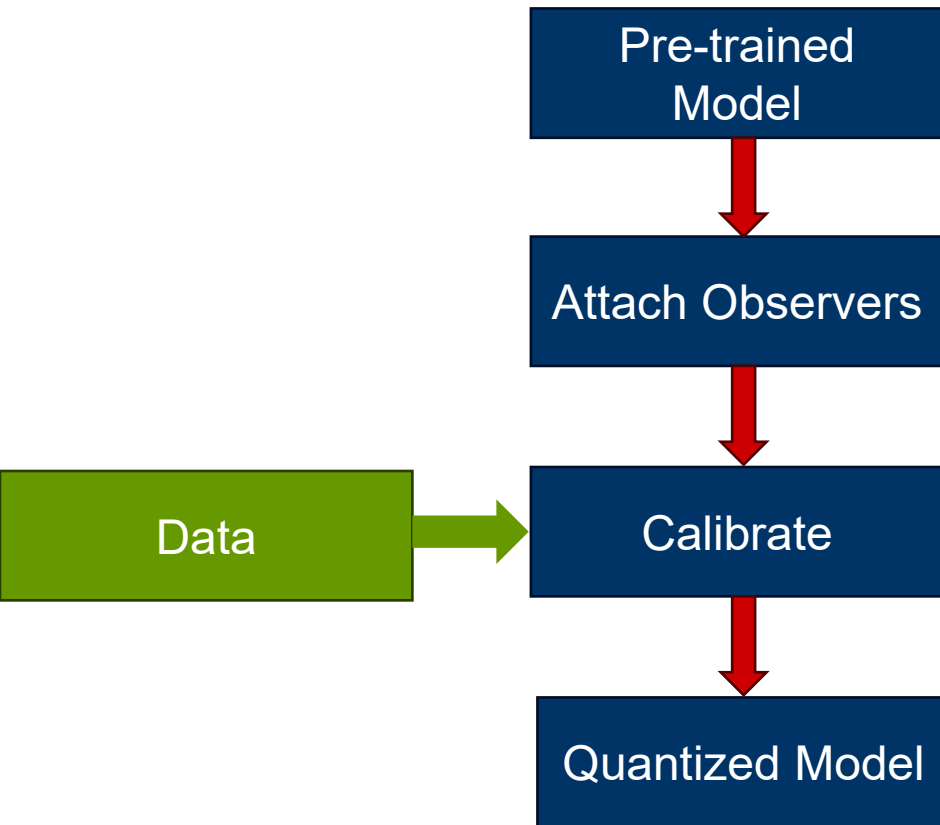
$$Y = W X$$

$$Y' = S_W (W_{int} - Z_W) S_X (X_{int} - Z_X)$$

$$Y' = (S_W S_X W_{int}) X_{int} + S_W S_X Z_W Z_X - S_W S_X Z_X W_{int} - (S_W S_X Z_W) X_{int}$$

Additional computational complexity for asymmetric quantization

[Option 1] Post Training Quantization

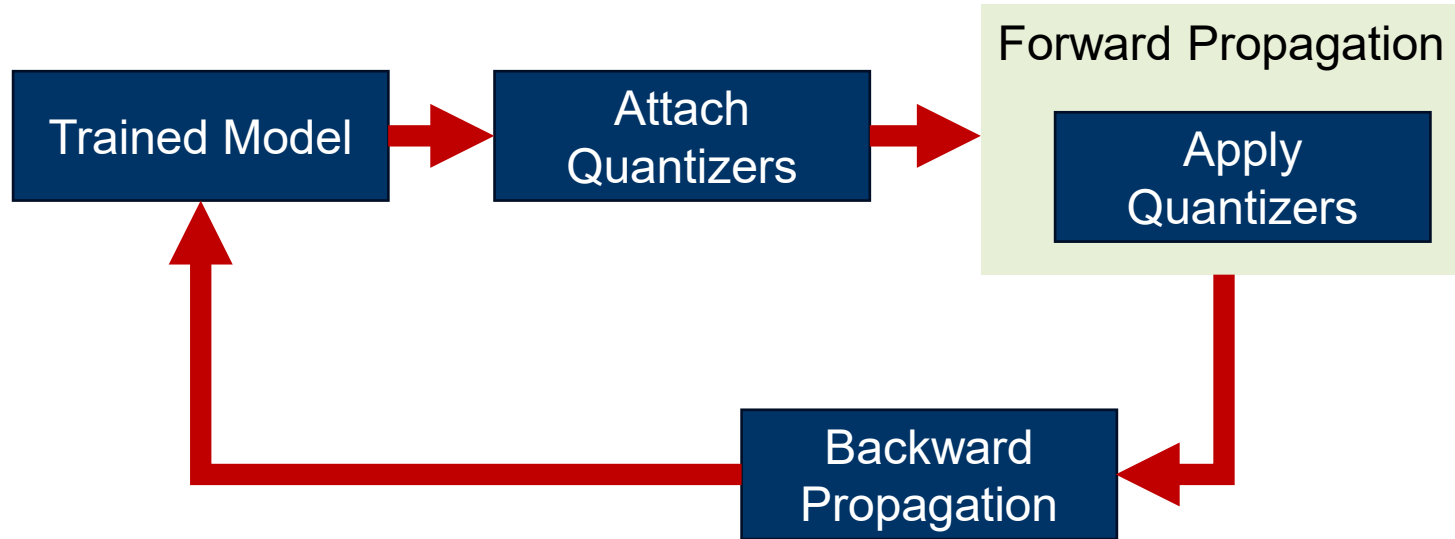


Calculate the S, Z parameters

Models	FP32	Per-tensor		Per-channel	
		W8A8	W4A8	W8A8	W4A8
ResNet18	69.68	69.60	68.62	69.56	68.91
ResNet50	76.07	75.87	75.15	75.88	75.43
MobileNetV2	71.72	70.99	69.21	71.16	69.79
InceptionV3	77.40	77.68	76.48	77.71	76.82
EfficientNet lite	75.42	75.25	71.24	75.39	74.01
DeeplabV3	72.94	72.44	70.80	72.27	71.67
EfficientDet-D1	40.08	38.29	0.31	38.67	35.08
BERT-base [†]	83.06	82.43	81.76	82.77	82.02

W→ Weights, A→ Activations

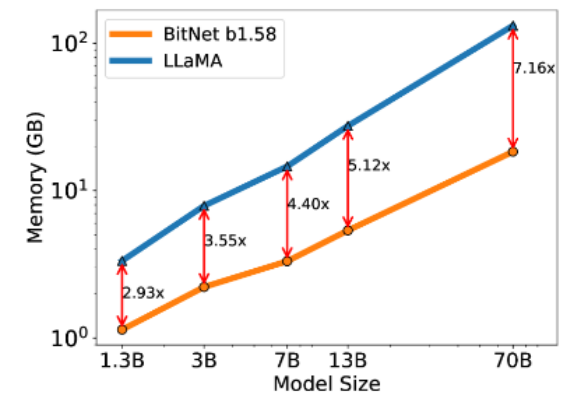
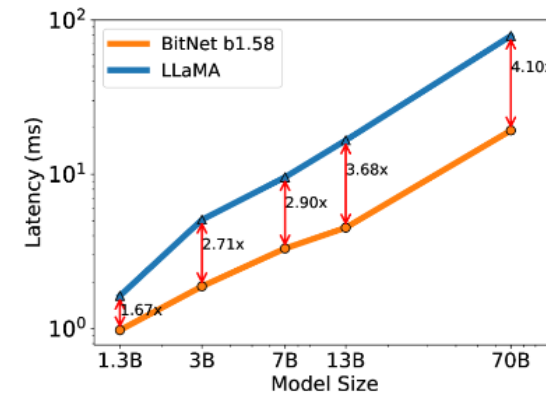
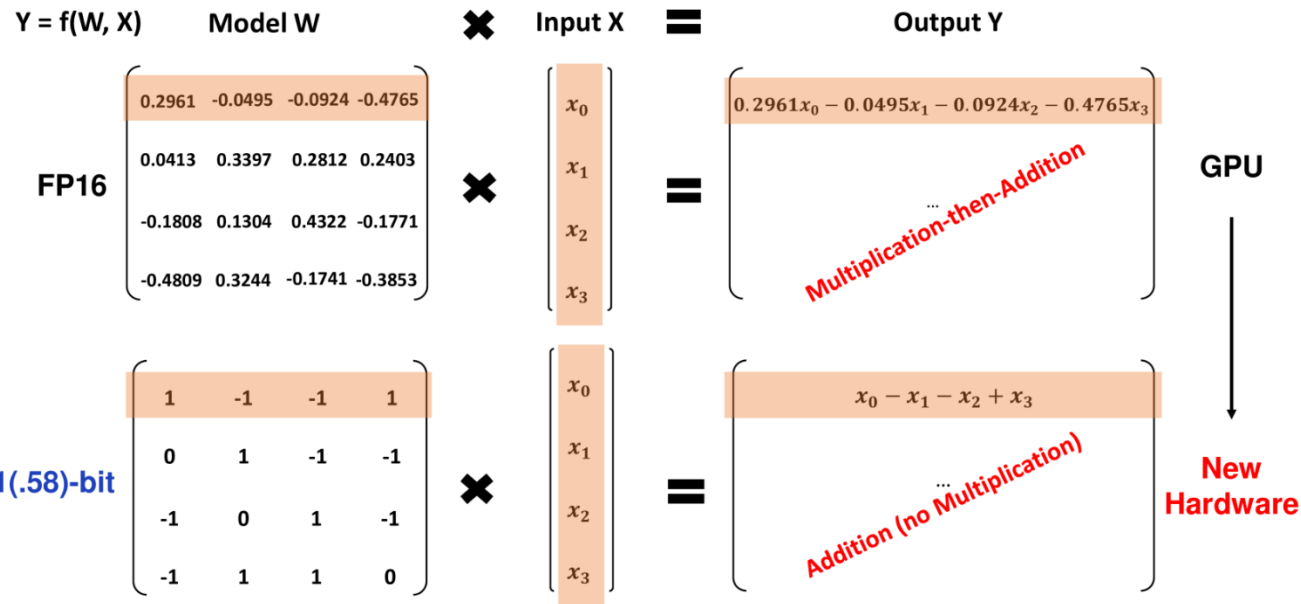
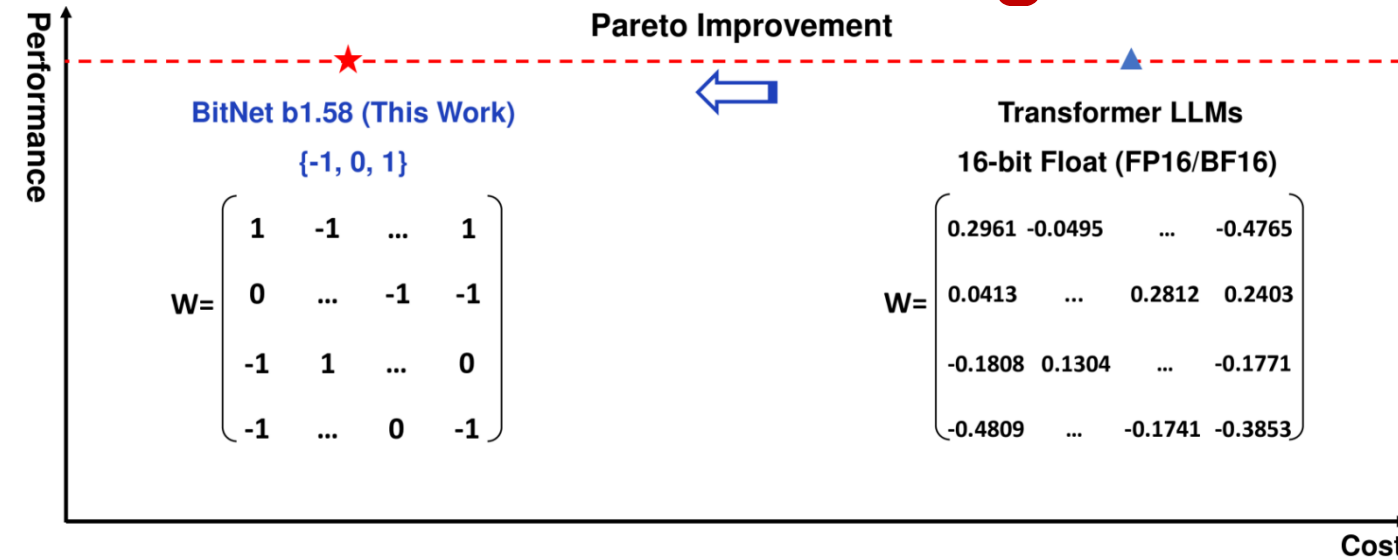
[Option 2] Quantization Aware Training



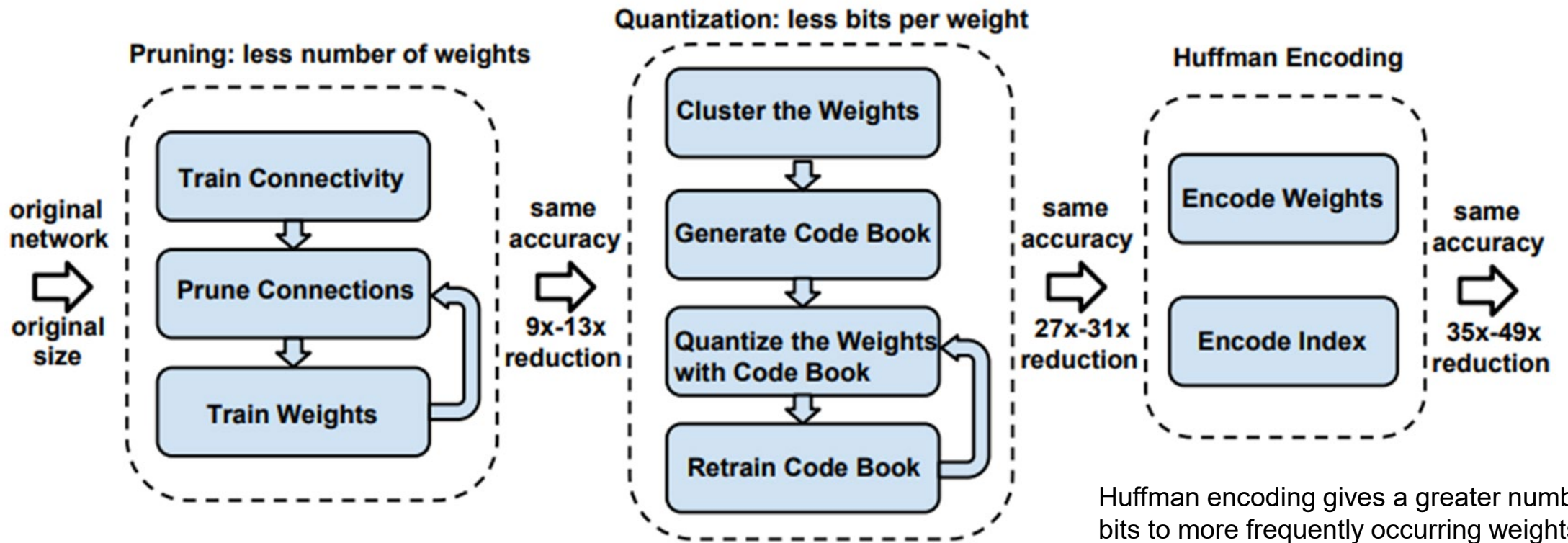
Models	FP32	Per-tensor			Per-channel		
		W8A8	W4A8	W4A4	W8A8	W4A8	W4A4
ResNet18	69.68	70.38	69.76	68.32	70.43	70.01	68.83
ResNet50	76.07	76.21	75.89	75.10	76.58	76.52	75.53
InceptionV3	77.40	78.33	77.84	77.49	78.45	78.12	77.74
MobileNetV2	71.72	71.76	70.17	66.43	71.82	70.48	66.89
EfficientNet lite	75.42	75.17	71.55	70.22	74.75	73.92	71.55
DeeplabV3	72.94	73.99	70.90	66.78	72.87	73.01	68.90
EfficientDet-D1	40.08	38.94	35.34	24.70	38.97	36.75	28.68
BERT-base	83.06	83.26	82.64	78.83	82.44	82.39	77.63

W→ Weights, A→ Activations

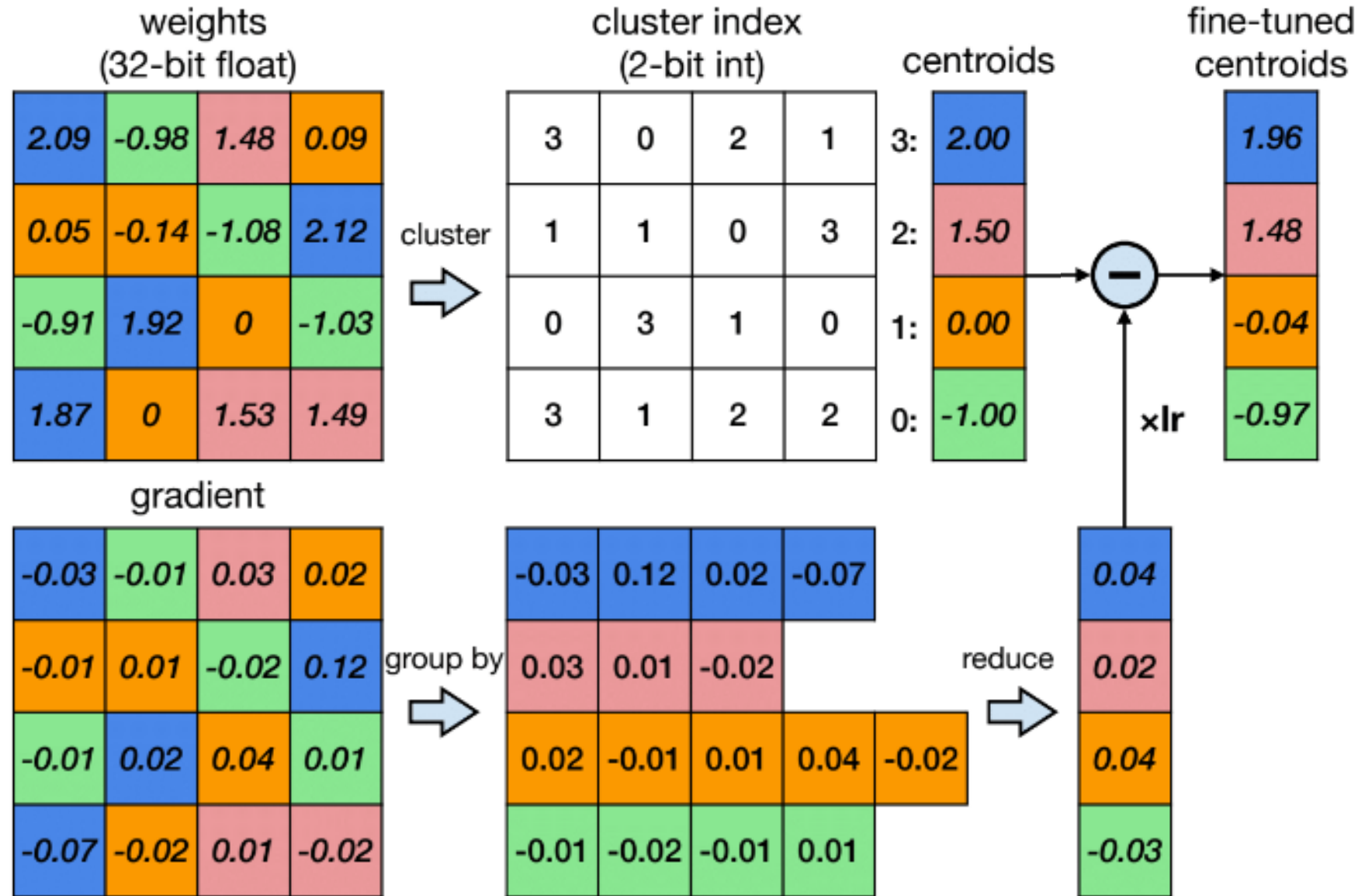
Quantization gone extreme!



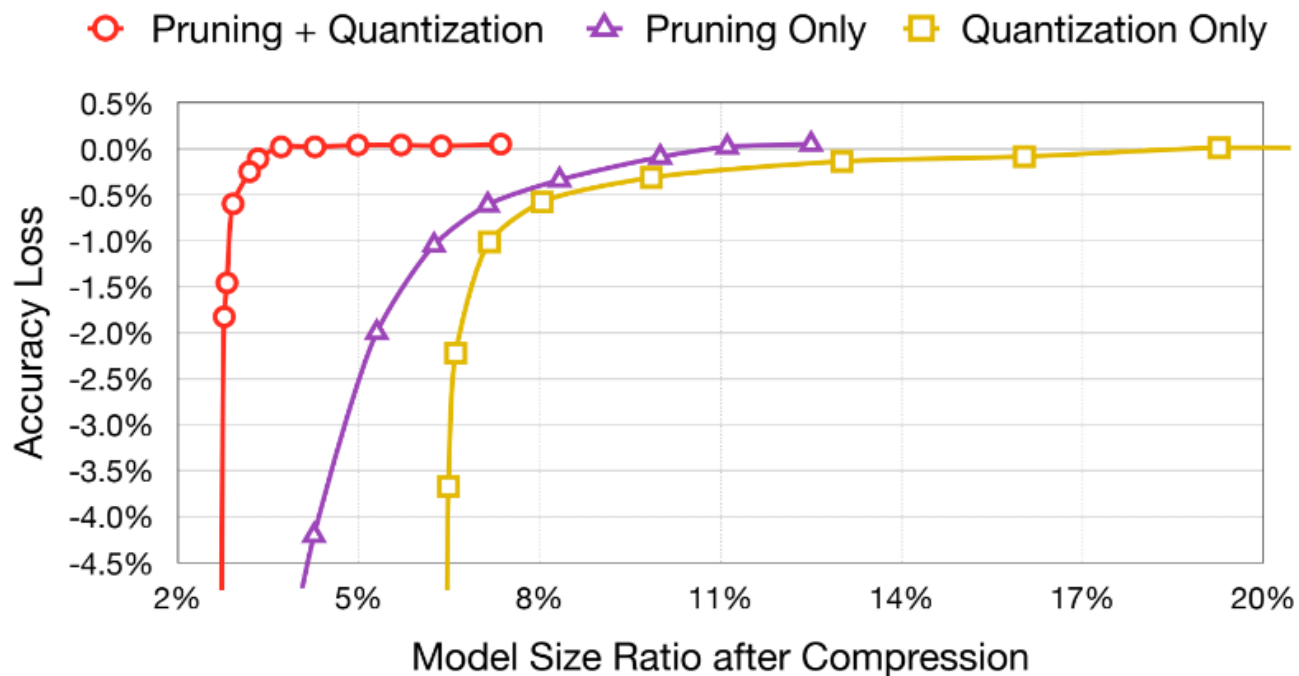
Weight Sharing



Weight Sharing – How?

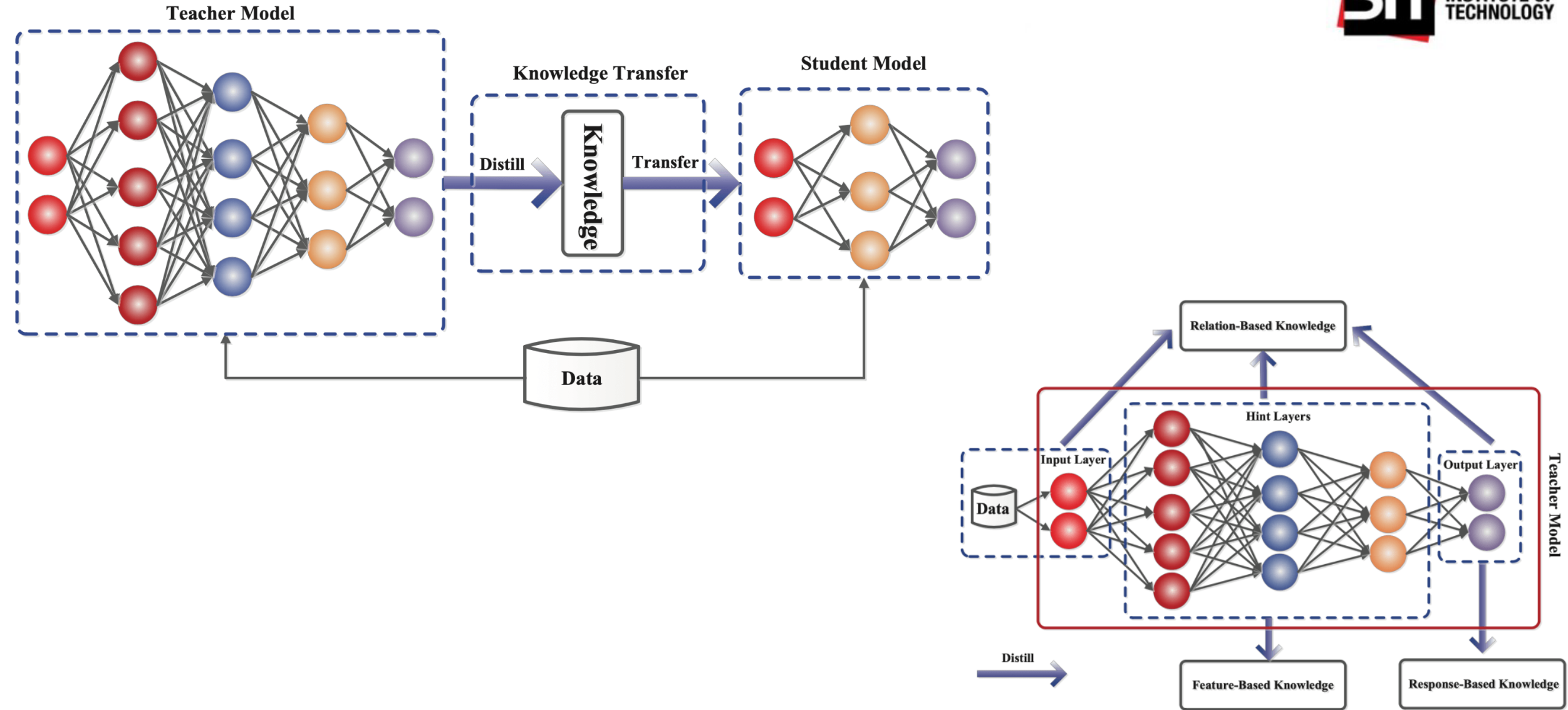


Weight Sharing - Results

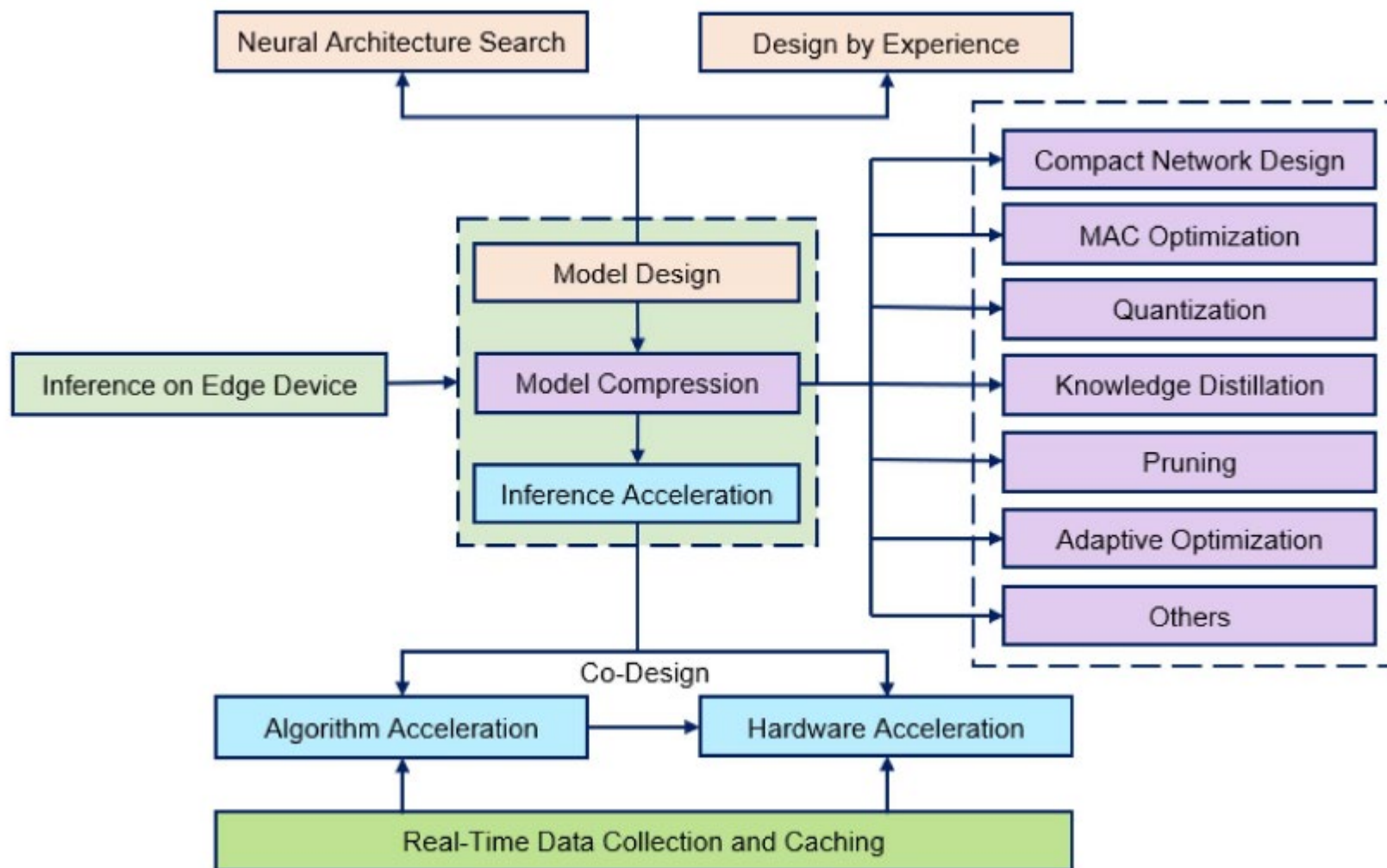


Network	Original Size	Compressed Size	Compression Ratio	Original Accuracy	Compressed Accuracy
LeNet-300	1070KB	27KB	40x	98.36%	98.42%
LeNet-5	1720KB	44KB	39x	99.20%	99.26%
AlexNet	240MB	6.9MB	35x	80.27%	80.30%
VGGNet	550MB	11.3MB	49x	88.68%	89.09%
GoogleNet	28MB	2.8MB	10x	88.90%	88.92%
ResNet-18	44.6MB	4.0MB	11x	89.24%	89.28%

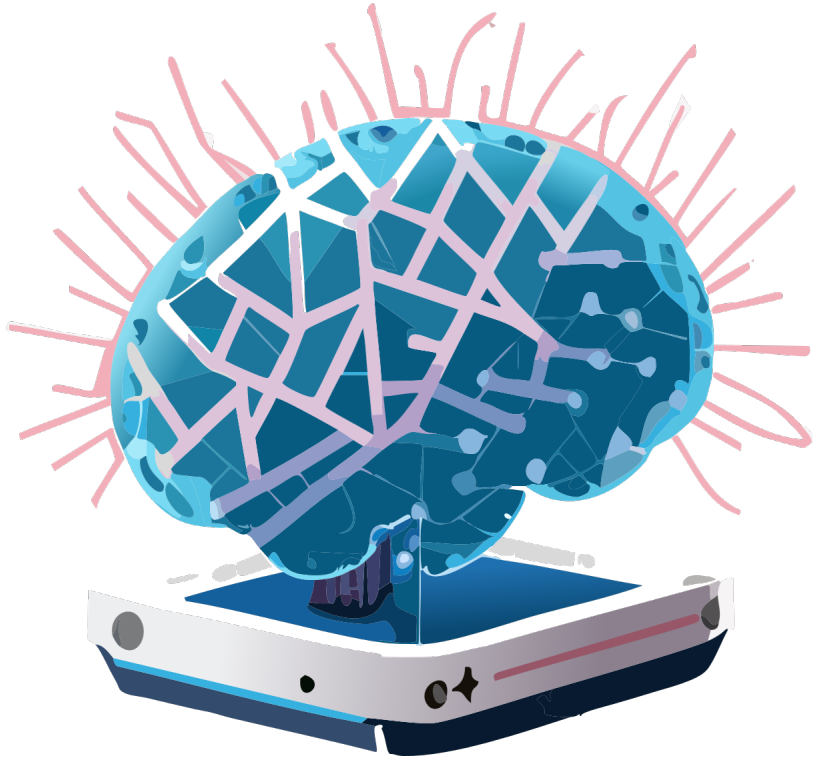
[Other] Options – Knowledge Distillation



Summary – DL on Edge



Summary



- *Adapting deep neural networks on edge devices is important*
- *Pruning, Quantization and their combinations is a popular approach*