



Accelerator Architectures for Machine Learning (AAML)

Lecture 3: Quantization

Tsung Tai Yeh

Department of Computer Science
National Yang-Ming Chiao Tung University



Acknowledgements and Disclaimer

- Slides was developed in the reference with
Joel Emer, Vivienne Sze, Yu-Hsin Chen, Tien-Ju Yang, ISCA 2019 tutorial
Efficient Processing of Deep Neural Network, Vivienne Sze, Yu-Hsin Chen,
Tien-Ju Yang, Joel Emer, Morgan and Claypool Publisher, 2020
Yakun Sophia Shao, EE290-2: Hardware for Machine Learning, UC
Berkeley, 2020
CS231n Convolutional Neural Networks for Visual Recognition, Stanford
University, 2020
- 6.5940, TinyML and Efficient Deep Learning Computing, MIT
- NVIDIA, Precision and performance: Floating point and IEEE 754
Compliance for NVIDIA GPUs, TB-06711-001_v8.0, 2017



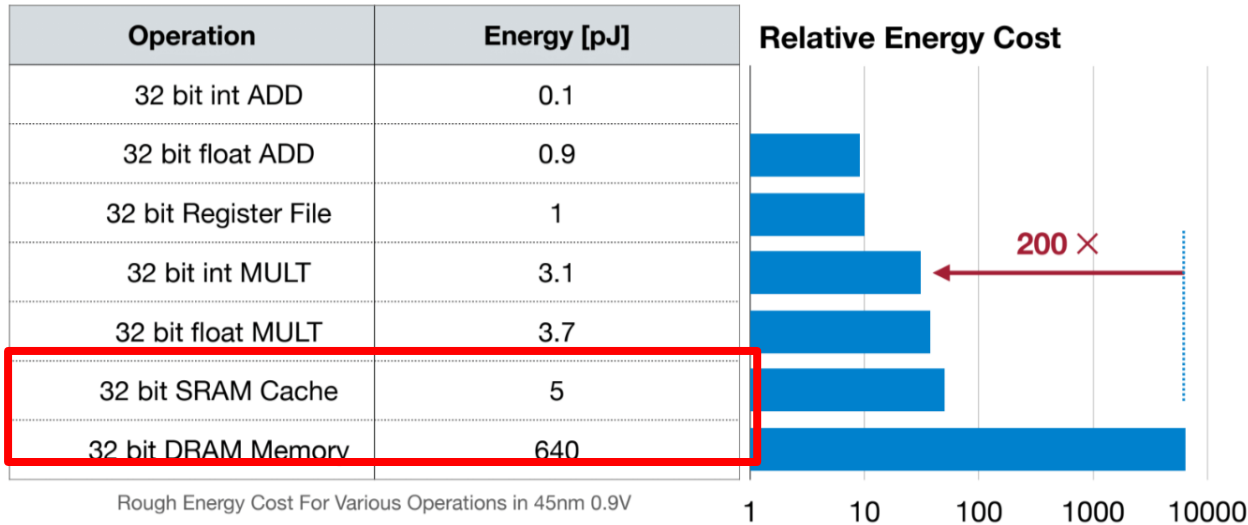
Outline

- K-Means-based Quantization
- Linear Quantization
- Binary and Ternary Quantization



Memory is Expensive !!

- **Data movement -> Move memory reference -> More energy**



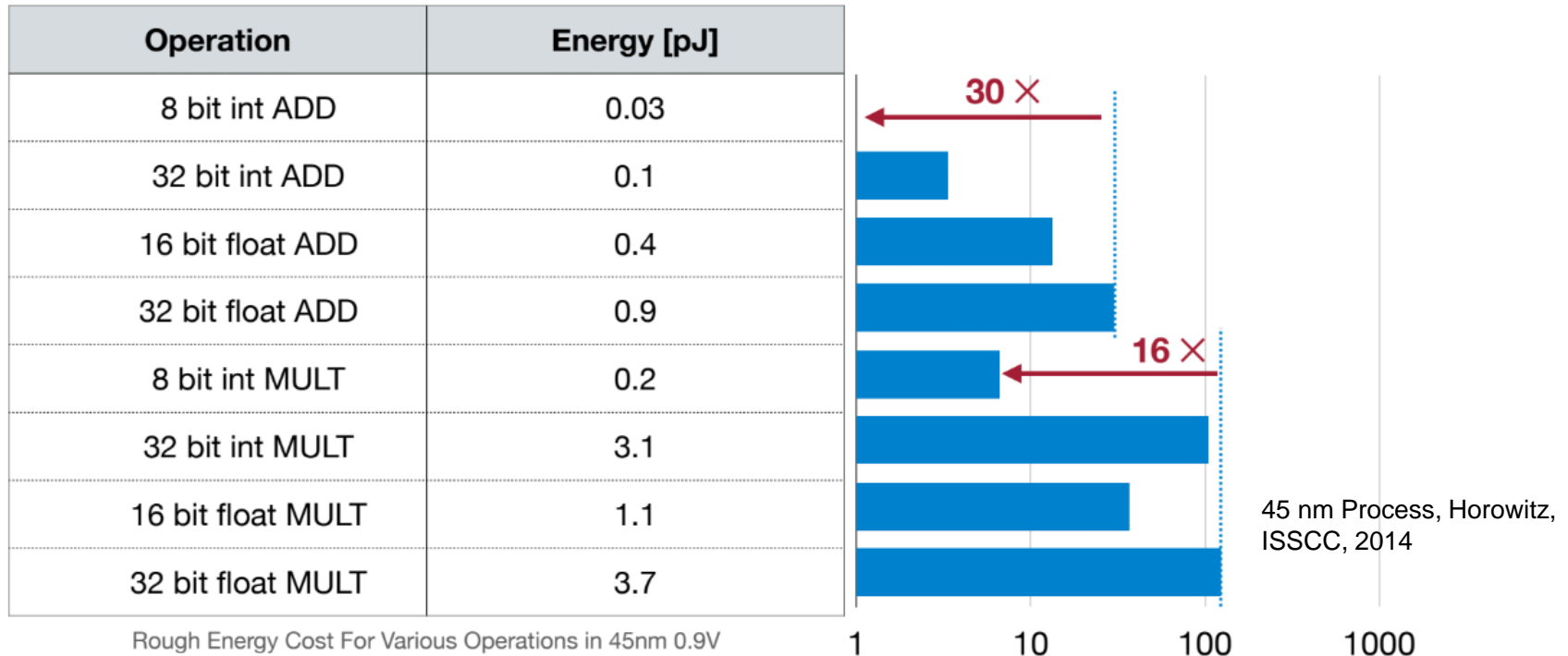
1  = 200 

This image is in the public domain



Low Bit-Width Operations are Cheap

- **Less Bit-Width -> Less energy**





Energy and Area Cost

Could we make the deep learning efficient by lowering the precision of data ?

Operation	Energy (pJ)	Area(um ²)
8b Add	0.03	36
16b Add	0.05	67
32b Add	0.1	137
16b FP Add	0.4	1360
32b FP Add	0.9	4184
16b FP Mult	1.1	1640
32b FP Mult	3.7	7700
32b SRAM Read (8KB)	5	
32b DRAM Read	640	

173X

4.7X



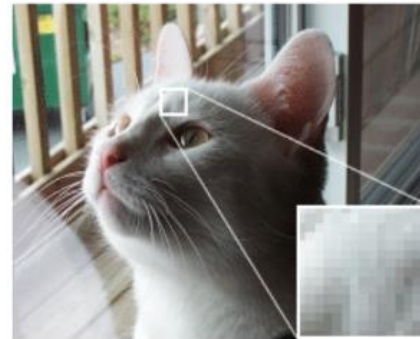
What is Quantization ?

- **Quantization**

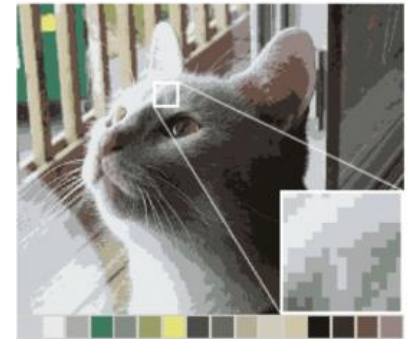
- The process of **constraining an input from a continuous or large set of values to a discrete set**



Original Image



16-Color Image



Images are in the public domain.



Numeric Data Types

- Fixed-point number



Integer . Fraction

“Decimal” Point



$$\begin{array}{cccccccc} \times & \times & \times & \times & \times & \times & \times & \times \\ -2^3 & + 2^2 & + 2^1 & + 2^0 & + 2^{-1} & + 2^{-2} & + 2^{-3} & + 2^{-4} \end{array} = 3.0625$$



$$\begin{array}{cccccccc} \times & \times & \times & \times & \times & \times & \times & \times \\ (-2^7 & + 2^6 & + 2^5 & + 2^4 & + 2^3 & + 2^2 & + 2^1 & + 2^0) \end{array} \times 2^{-4} = 49 \times 0.0625 = 3.0625$$

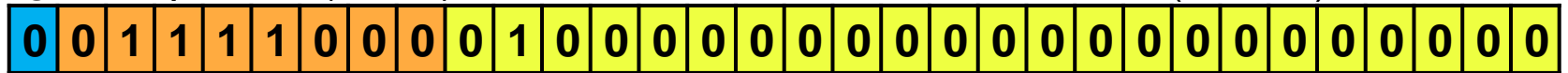


IEEE 765 Single Precision Float Point

- **Sign** determines the sign of the number
- **Exponent** (8 bit) represent -127 (all 0s) and +128 (all 1s)
- **Significand** (23 fraction bits), total precision is 24 bits (23 + 1 implicit leading bit) $\log_{10}(2^{24}) \approx 7.225$ digital bit

Sign Exponent (8 bits)

Mantissa/Fraction (23 bits)

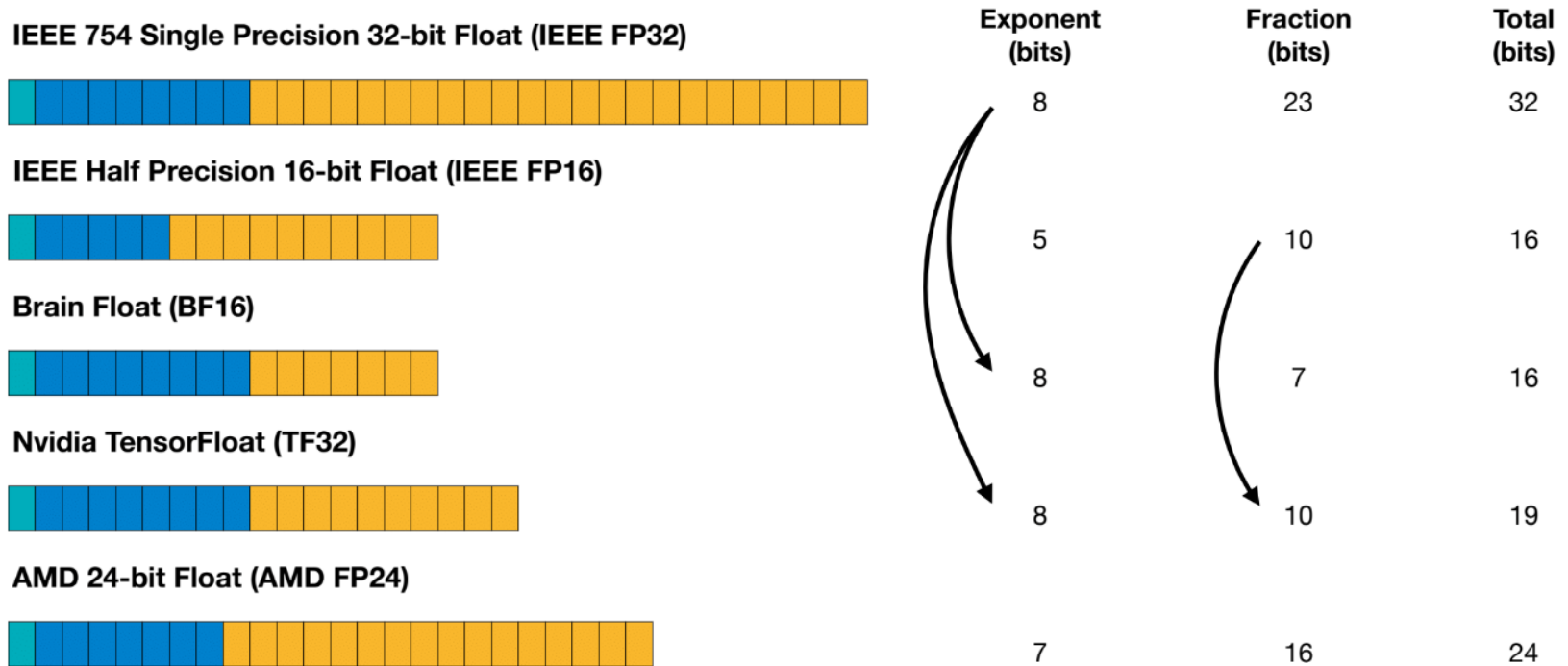


$$value = (-1)^{sign} \times 2^{(e-127)} \times \left(1 + \sum_{i=1}^{23} b_{(23-i)} 2^{-i}\right)$$



Floating-Point Number

- Exponent Width -> Range; Fraction Width-> Precision





Number Representation

	1	8	23	
FP32	S	E	M	1.2E-38 to 3.4E+38
	1	5	10	
FP16	S	E	M	6.1E-5 to 6.6E+4
	1	31		
INT32	S	M		2147483648 to 2147483647
	1	15		
INT16	S	M		-32,768 to 32,767
	1	8		
INT8	S	M		-128 ~ 127



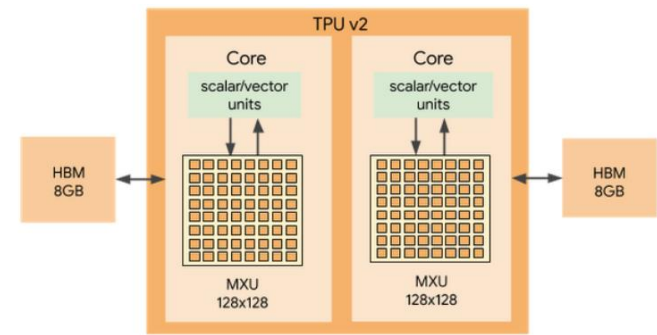
Choosing bFloat16

- **Motivation**

- The physical size of a hardware multiplier scales with the square of the mantissa width
- Mantissa bit length – FP32: 23, FP16: 10, BF16: 7

- **BF16**

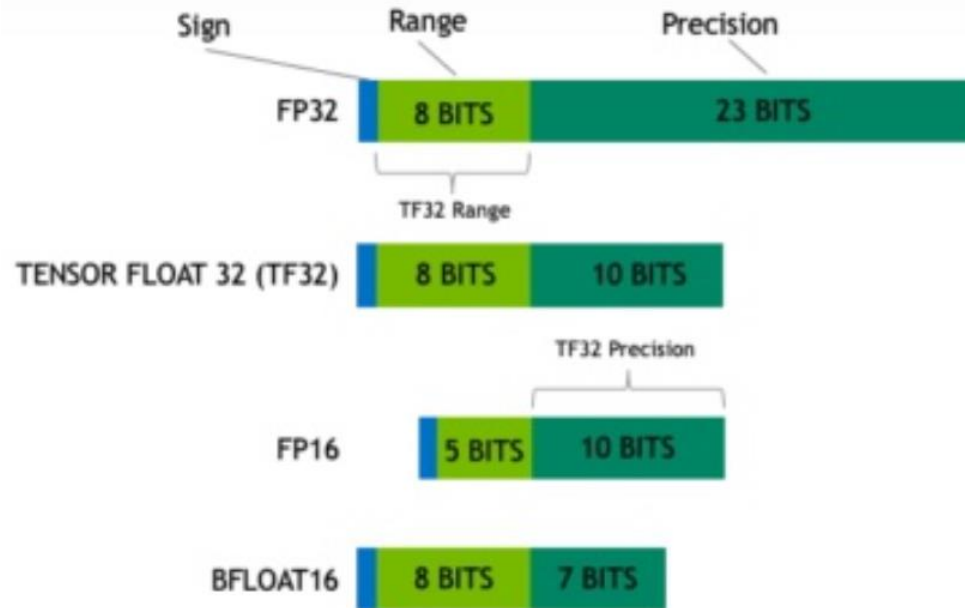
- 8 X smaller than an FP32 multiplier
- Has the same exponent size as FP32
- No require special handling (loss scaling) in the FP16 conversion
- XLA compiler's automatic format conversion
- In side the MXU, multiplications are performed in BF16 format
- Accumulations are performed in full FP32 precision





Nvidia's TF32

- **Nvidia's TF32**
 - 19-bit (BF19)
 - 1-bit sign, 8-bit exponent
10-bit fraction
 - Fuse BF16 and FP16
 - BF16: 8-bit exponent +
 - FP16: 10-bit fraction
 - Nvidia A100 Tensor Core
 - TF32: 156 TFLOPS
 - FP16/BF16: 312 TFLOPS



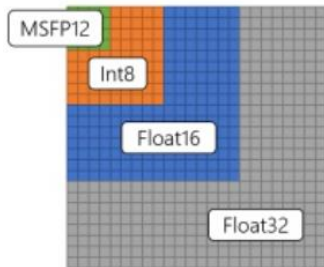
<https://zhuanlan.zhihu.com/p/449857213>



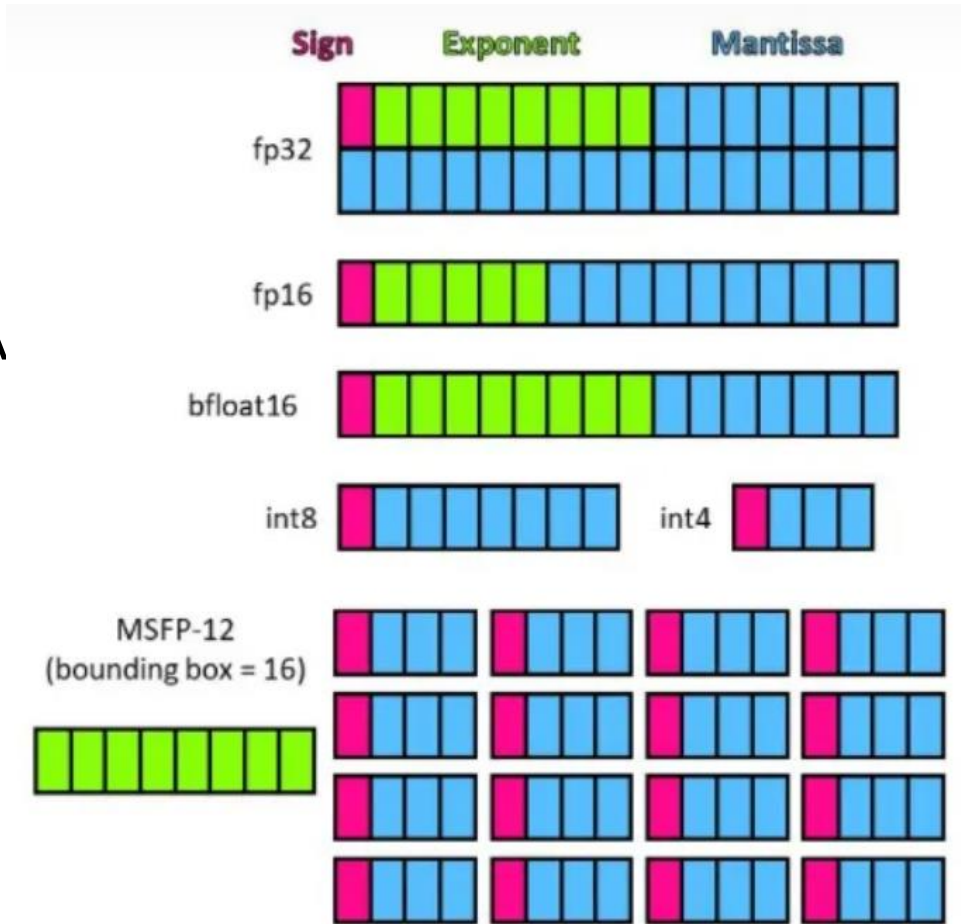
Microsoft MSFP

● Microsoft MSFP

- Used in Brainwave FPGA
- 8-bit shared exponent
- 1-bit sign, 3-bit fraction
- A group of INT4 vector shares 8-bit exponent



MAC Area & Energy





FP8 and Tesla CFloat

- **FP8 (1-5-2)**
 - Large loss in MobileNet v2
 - Hybrid FP8 (HFP8)
 - Use FP(1-4-3) in forward
 - Use FP(1-5-2) in backward
- **Tesla Dojo Cfloat (configurable float)**
 - Configurable exponent and mantissa
 - Use software to choose appropriate Cfloat format
 - CF16
 - CF8 (1-4-3), CF8 (1-5-2)

c. Trans-precision Inference Accuracy of FP32 models in FP8 1-5-2 precision

FP32 Model	Baseline	FP8 1-5-2
MobileNet_v2 ImageNet	71.81	52.51
ResNet50 ImageNet	76.44	75.31
DensetNet121 ImageNet	74.76	73.64
MaskRCNN COCO [†]	33.58 29.27	32.83 28.65

[†] Box and Mask average precision

<https://proceedings.neurips.cc/paper/2019/file/65fc9fb4897a89789352e211ca2d398f-Paper.pdf>



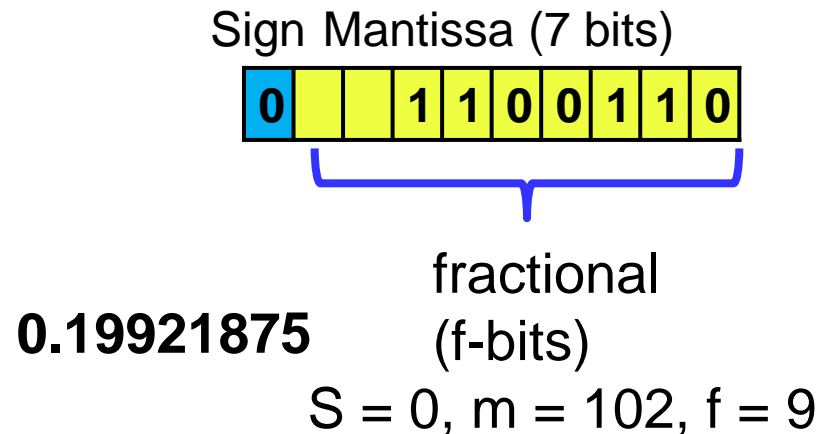
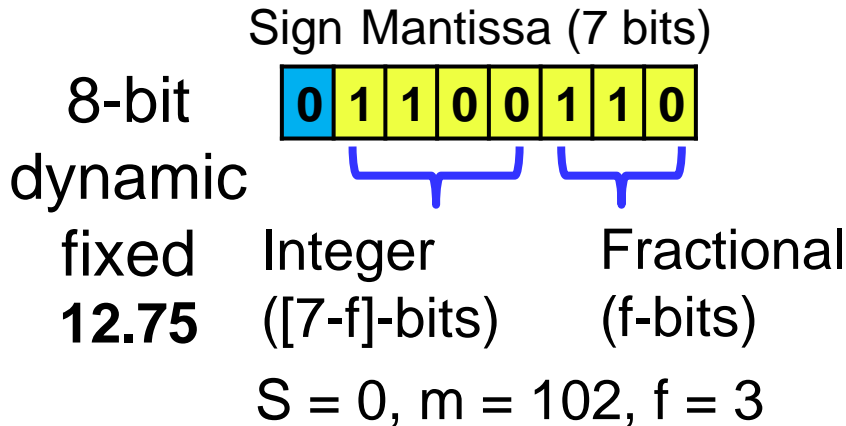
How to Determine Bit Width on DNN ?

- For accuracy, DNN operations decide bit width to achieve sufficient precision
- Which DNN operations affect the accuracy ?
 - **For inference:** weights, activations, and partial sums
 - **For training:** weights, activations, partial sums, gradients, and weight update



Dynamic Fixed Point

- Allow “f” to vary based on data type and layer
- In large layers, the outputs are the result of many accumulations
- The value of network parameters are much smaller than layer output -> **varying bit widths on parameters and outputs**

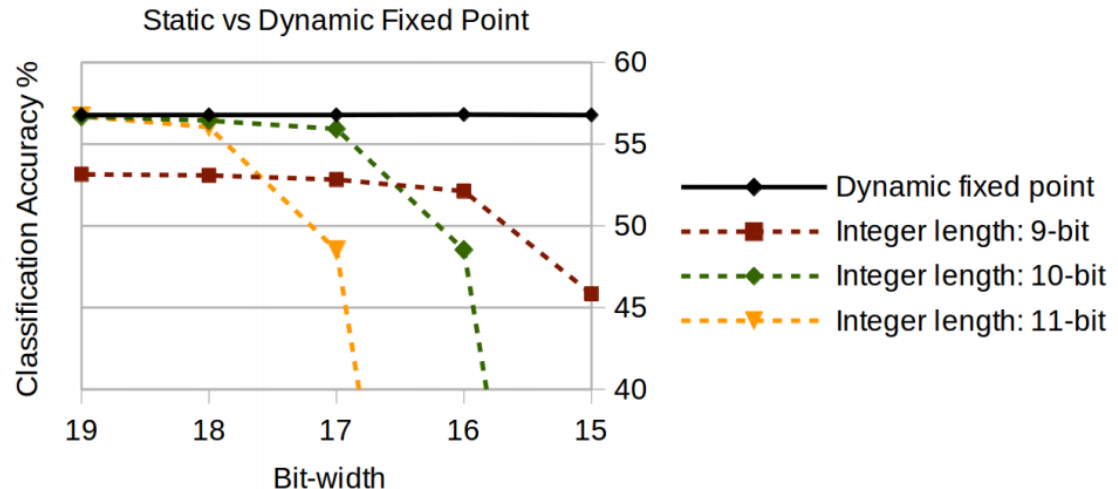




Impact on Accuracy

- The accuracy drops in the small bit width when using static fixed point
- Stable accuracy variation is shown in dynamic fixed point (why ?)

Top-1
accuracy of
CaffeNet on
ImageNet





Impact on Accuracy

- Small bit width cannot adapt to every DNN models very well (training)

	Layer outputs	CONV parameters	FC parameters	Fixed point accuracy
LeNet (Exp 1)	4-bit	4-bit	4-bit	99.0%
LeNet (Exp 2)	4-bit	2-bit	2-bit	98.8%
SqueezeNet	8-bit	8-bit	8-bit	57.1%
CaffeNet	8-bit	8-bit	8-bit	56.0%
GoogleNet	8-bit	8-bit	8-bit	66.6%



Precision Varies from Layer to Layer

- Accuracy varies with the different bit widths in layers
- How to find out the best bit width in each layer while maintaining high accuracy ?

AlexNet

Error rate	Bit per layer
1%	10-8-8-8-8-8-6-4
2%	10-8-8-8-8-8-5-4
5%	10-8-8-8-7-7-5-3
10%	9-8-8-8-7-7-5-3



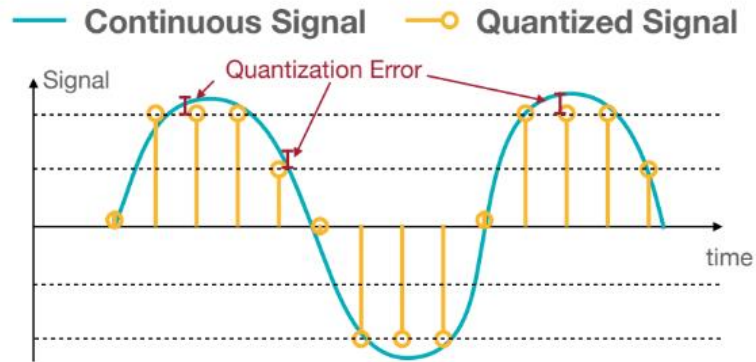
Takeaway Questions

- What are advantages to use BF16 instead of FP16 ?
 - (A) Fast conversion from FP32
 - (B) Get more precise value
 - (C) Represent few different values
- What are benefits to use lower precision data type on neural network ?
 - (A) Reduce the latency of DNN models
 - (B) Save the memory space
 - (C) Lower the power consumption of the accelerator

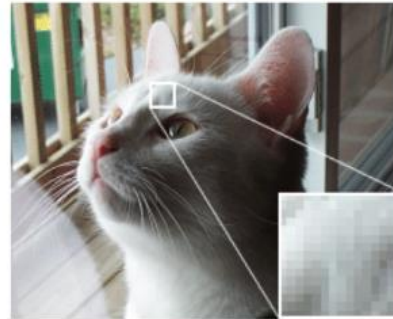


What is Quantization ?

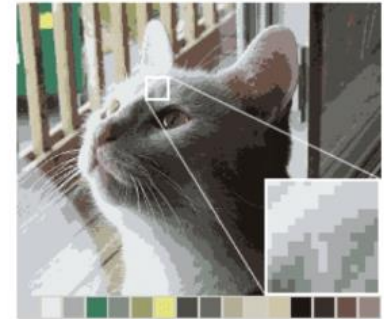
- **Quantization** is the process of **constraining an input from a continuous or large set of values to a discrete set**



Original Image



16-Color Image



Images are in the public domain.

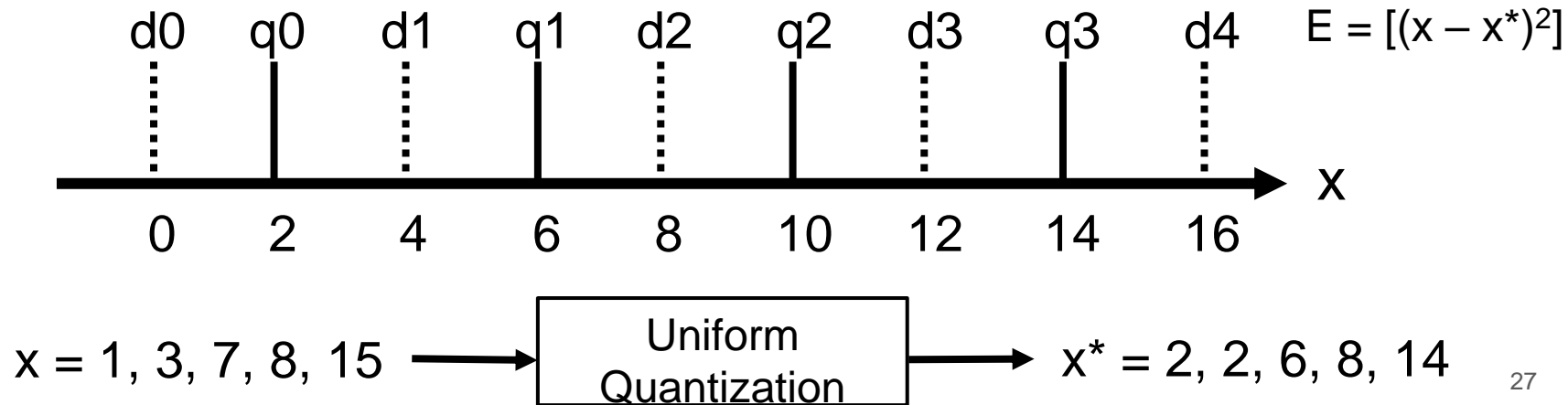
“Palettization”

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



Data Quantization

- **Quantization**
 - Maps data from a full precision to reduced one
- **Quantization error**
 - Measures the average difference between the original full precision and quantized values

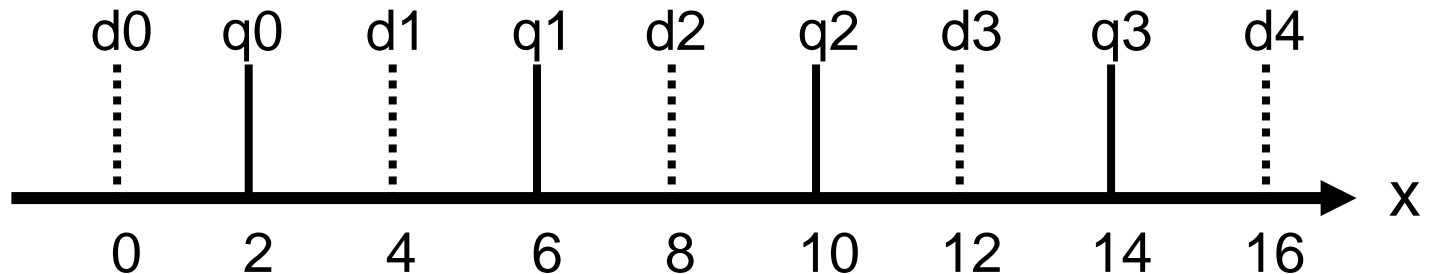




Types of Quantization

- **Uniform Quantization**

- Quantized values are equally spaced out
- x^* can take on are $\{2, 6, 10, 14\}$ with level = 4
- Decision boundaries d_i are used to decide the quantization value that x should be mapped to

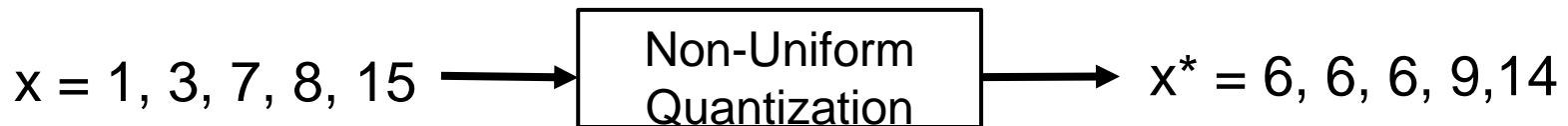
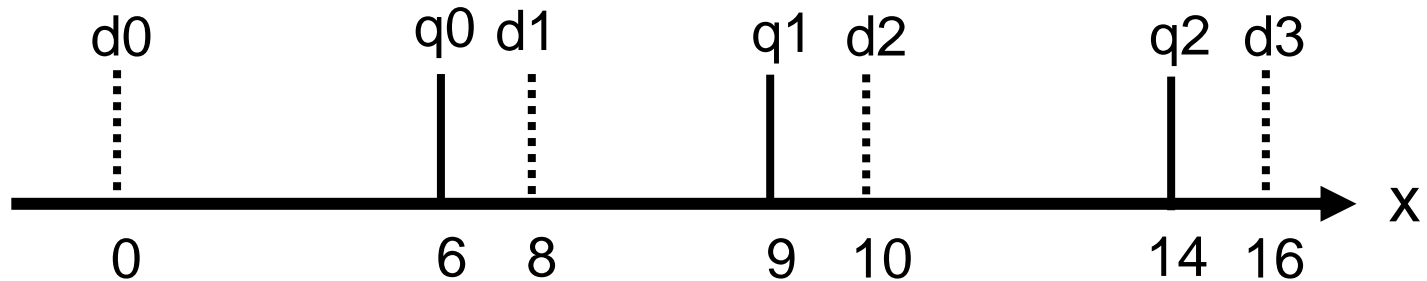




Types of Quantization

- **Non-uniform quantization**

- Spacing can be computed e.g. logarithmic or with look-up-table
- Fewer unique values can make weight sharing and compression





K-Means-based Weight Quantization

- **Storage**
 - Integer Weights; Floating-Point Codebook
- **Computation**
 - Floating-Point Arithmetic

weights
(32-bit float)

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

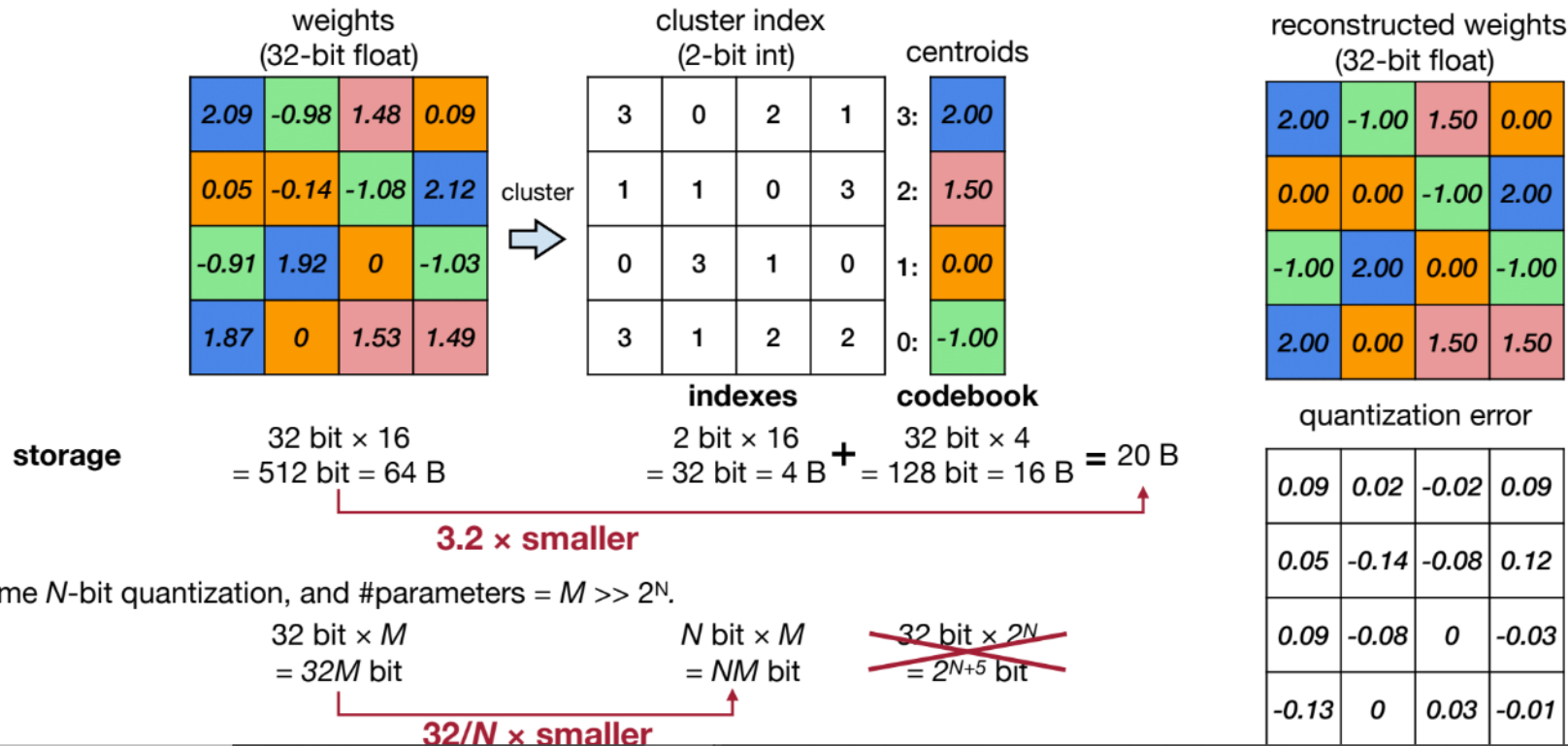
~~2.09, 2.12, 1.92, 1.87~~



2.0



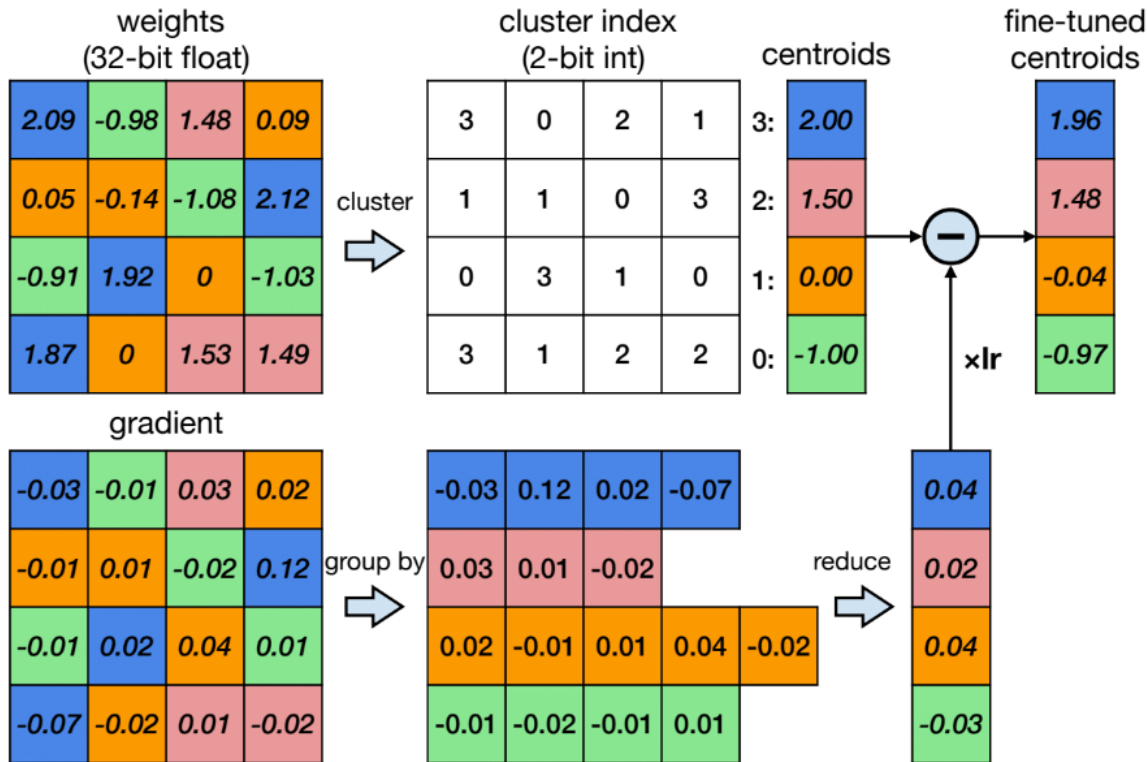
K-Means-based Weight Quantization





K-Means-based Weight Quantization

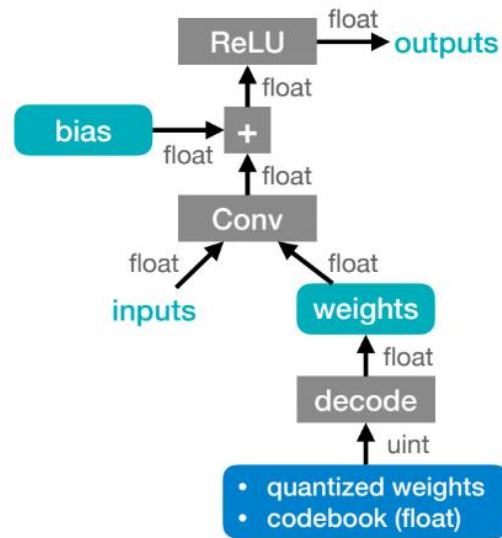
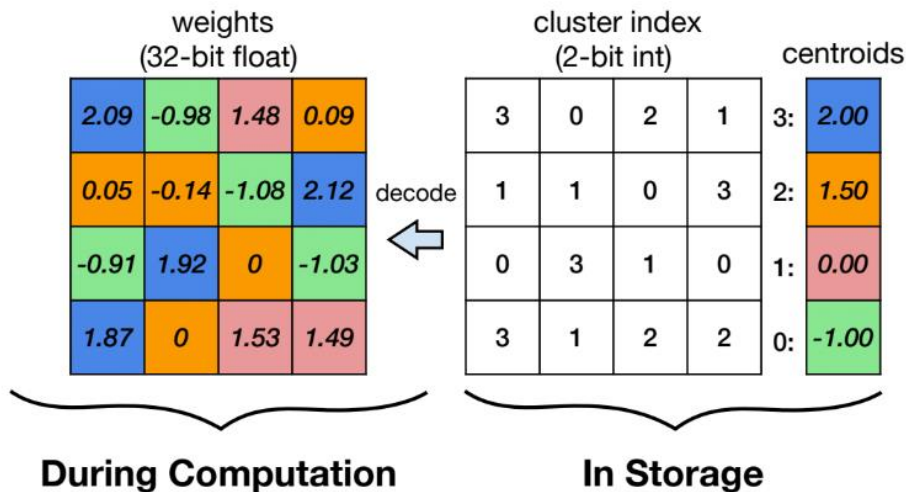
- Fine-tuning Quantized Weights**





K-Means-based Weight Quantization

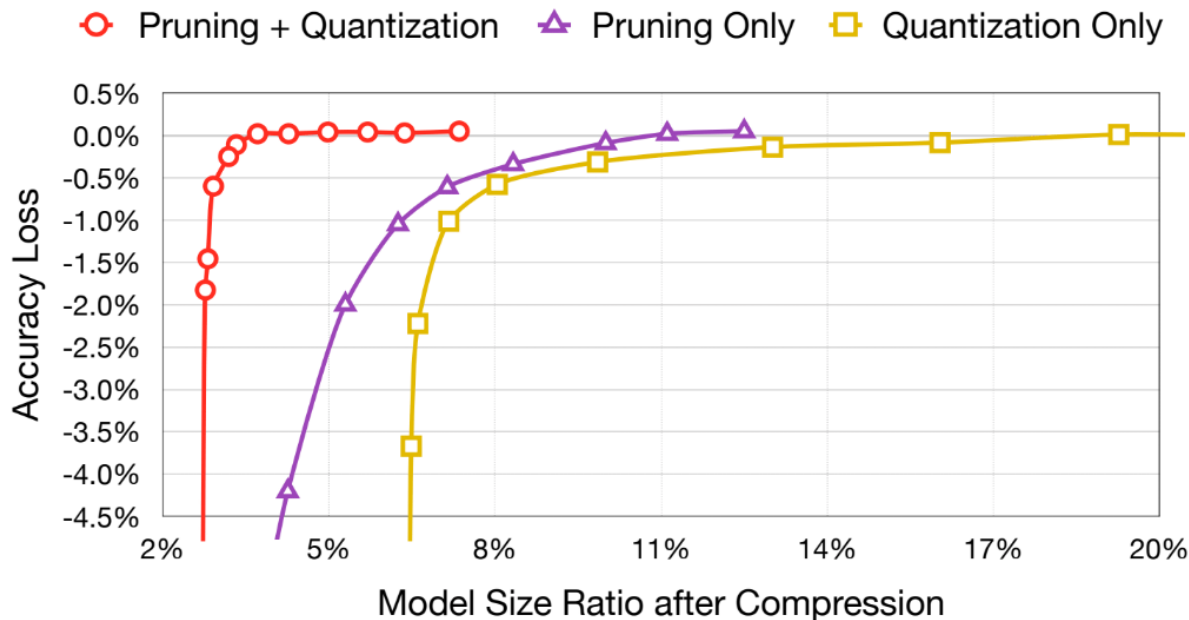
- Weights are decompressed using a lookup table during runtime inference
- Only saves storage cost of a neural network model
- All the computation and memory access are still floating-point





K-Means-based Weight Quantization

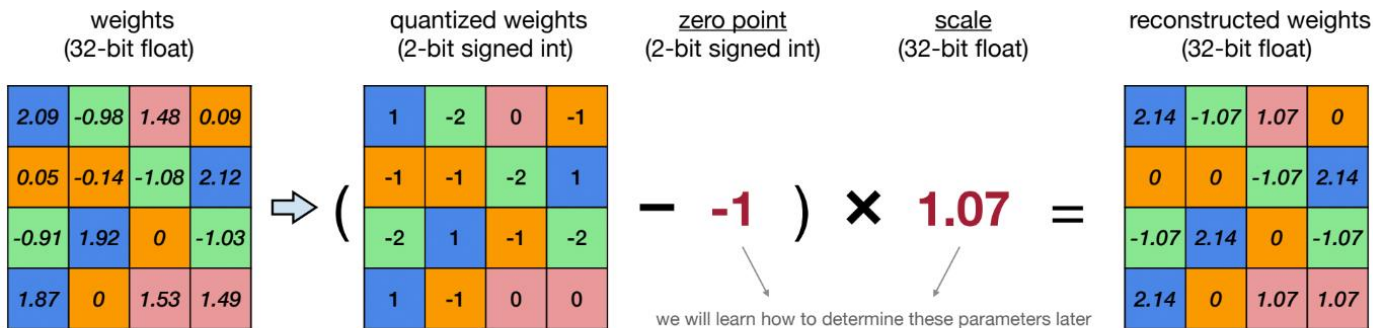
- Accuracy vs. compression rate for AlexNet on ImageNet dataset





What is Linear Quantization ?

- An affine mapping of integers to real numbers
- **Storage:** Integer Weights; **Computation:** Integer Arithmetic



Binary	Decimal
01	1
00	0
11	-1
10	-2

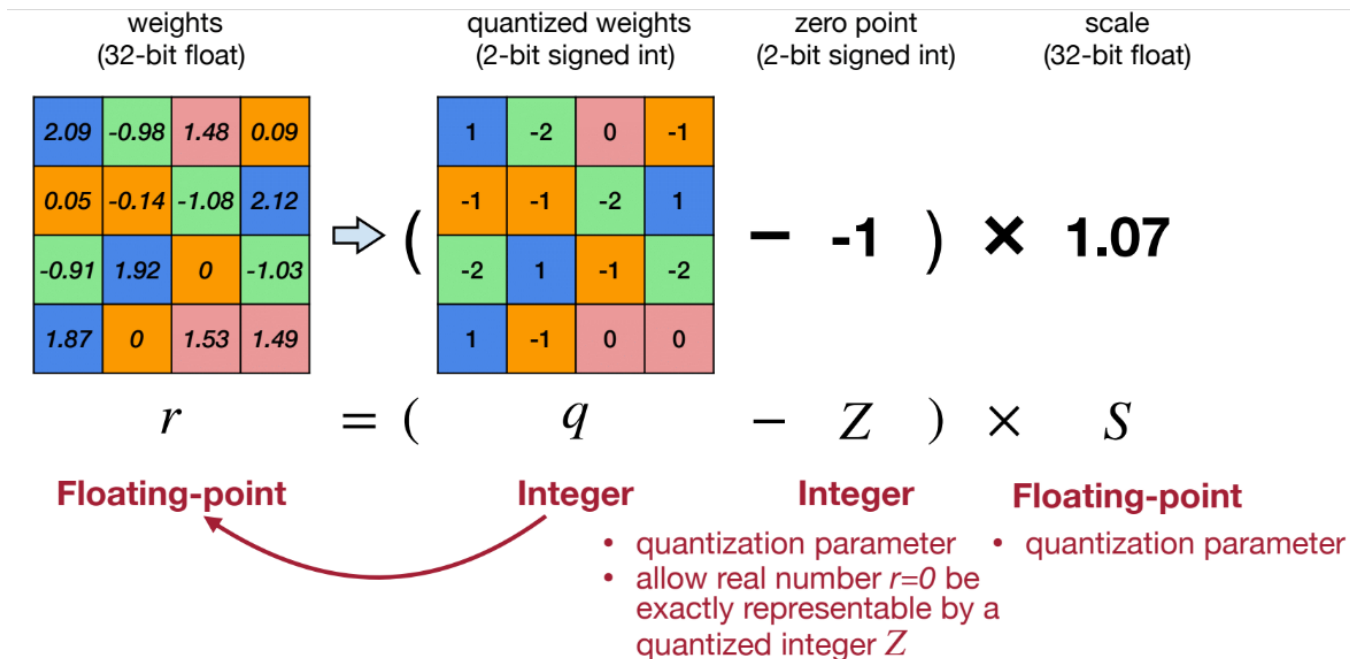
quantization error

-0.05	0.09	0.41	0.09
0.05	-0.14	-0.01	-0.02
0.16	-0.22	0	0.04
-0.27	0	0.46	0.42



Linear Quantization

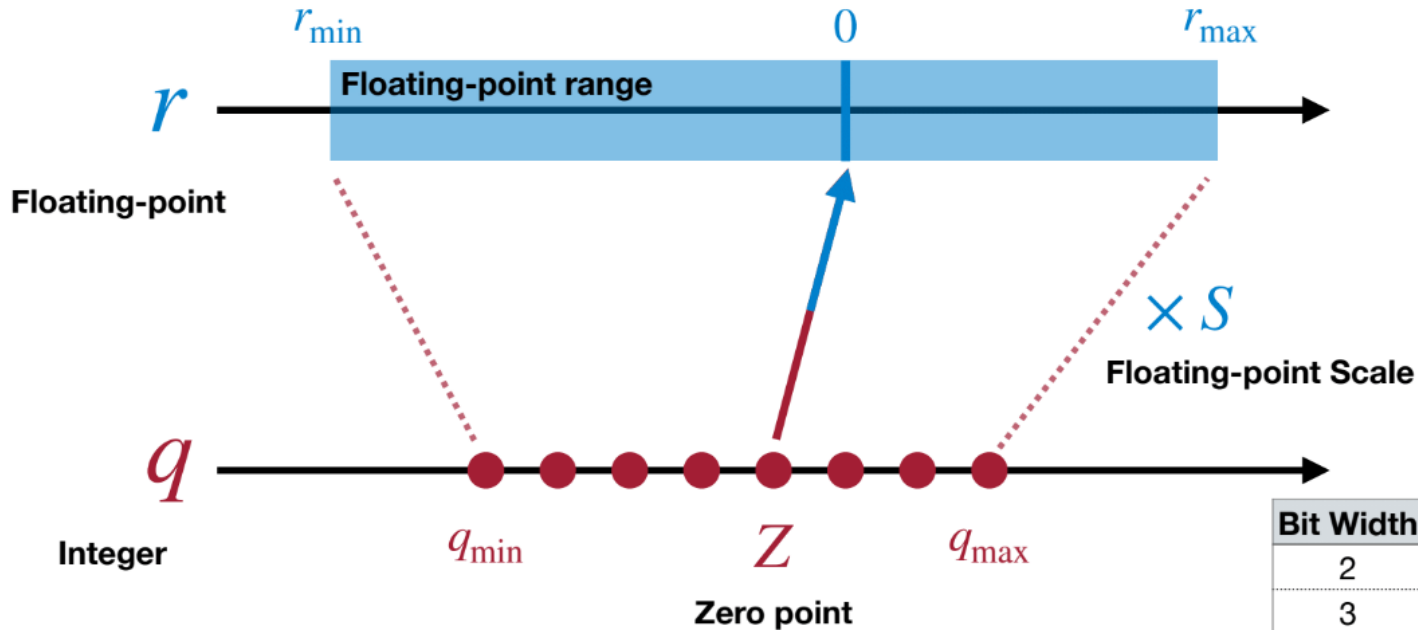
- An affine mapping of integers to real numbers ($r = S(q - Z)$)





Linear Quantization

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

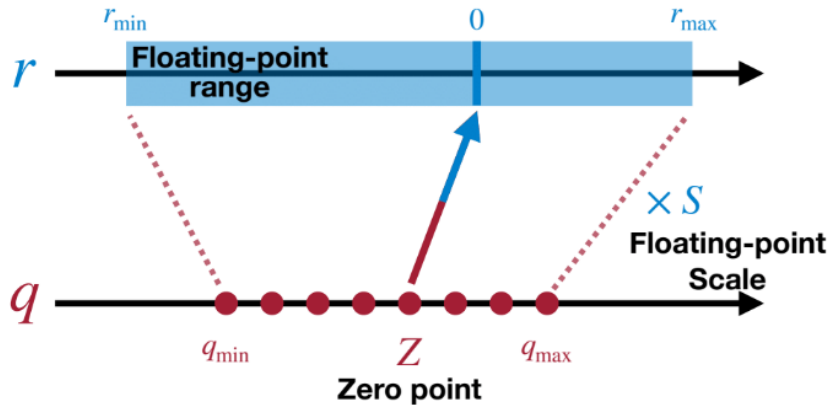


Bit Width	q_{\min}	q_{\max}
2	-2	1
3	-4	3
4	-8	7
N	-2^{N-1}	$2^{N-1}-1$



Scale of Linear Quantization

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

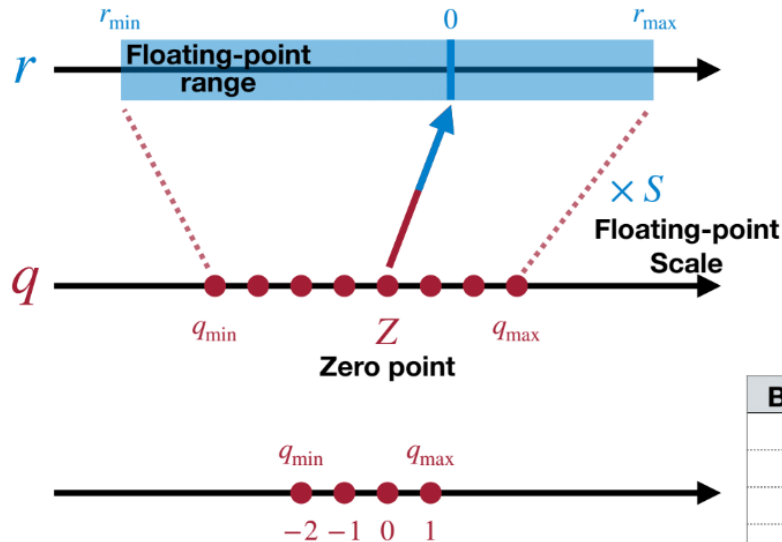


$$\begin{aligned} r_{\max} &= S(q_{\max} - Z) \\ r_{\min} &= S(q_{\min} - Z) \end{aligned} \quad \ominus$$
$$\downarrow$$
$$r_{\max} - r_{\min} = S(q_{\max} - q_{\min})$$
$$S = \frac{r_{\max} - r_{\min}}{q_{\max} - q_{\min}}$$



Scale of Linear Quantization

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



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

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

$$= \frac{2.12 - (-1.08)}{1 - (-2)}$$

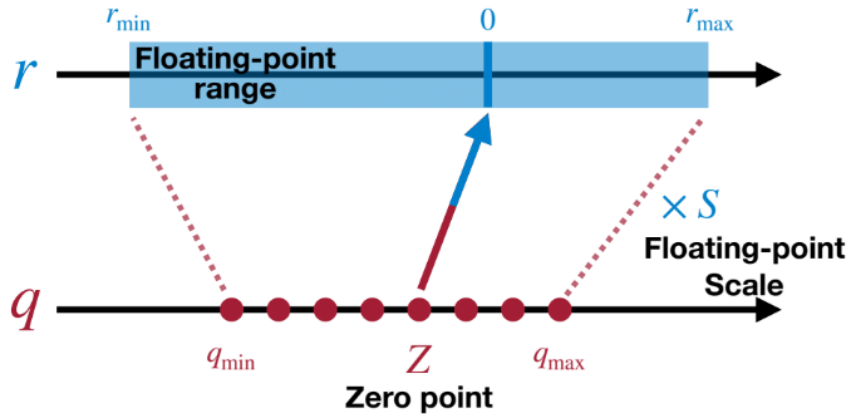
$$= 1.07$$

Binary	Decimal
01	1
00	0
11	-1
10	-2



Zero Point of Linear Quantization

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



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

↓

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

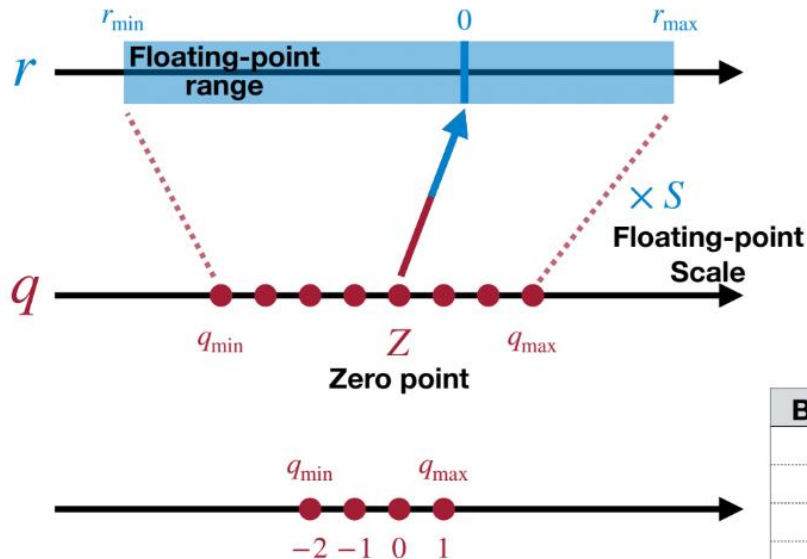
↓

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



Zero Point of Linear Quantization

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



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

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

$$= \text{round}\left(-2 - \frac{-1.08}{1.07}\right)$$

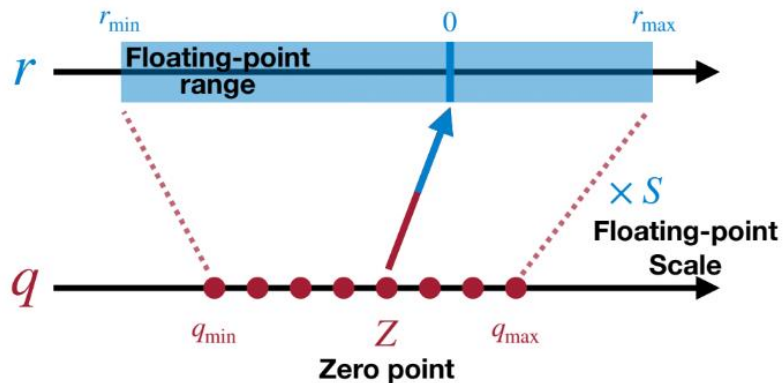
$$= -1$$

Binary	Decimal
01	1
00	0
11	-1
10	-2



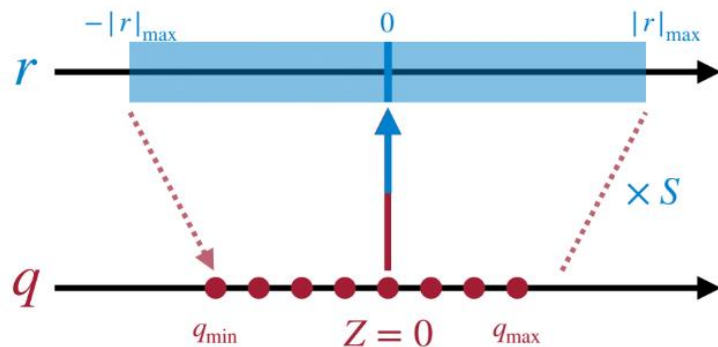
Symmetric Linear Quantization

- Full range mode



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

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



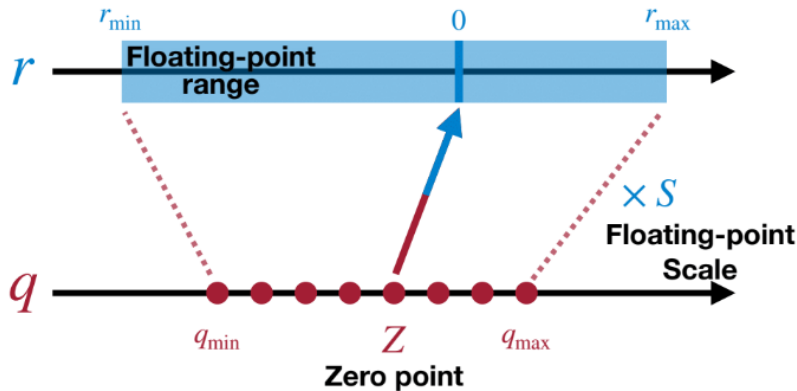
$$S = \frac{r_{\min}}{q_{\min} - Z} = \frac{-|r|_{\max}}{q_{\min}} = \frac{|r|_{\max}}{2^{N-1}}$$

- use full range of quantized integers
- example: PyTorch's native quantization, ONNX



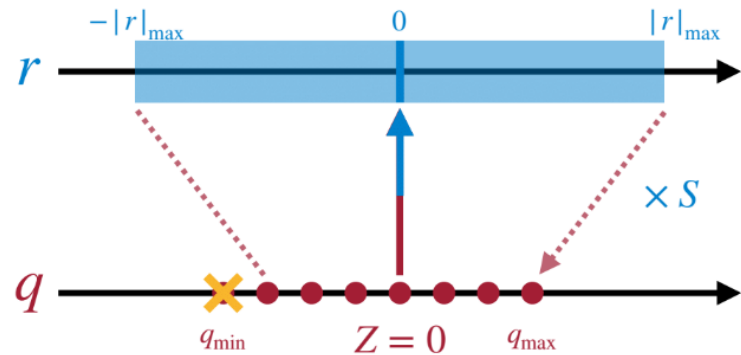
Symmetric Linear Quantization

- Restricted range mode



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

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

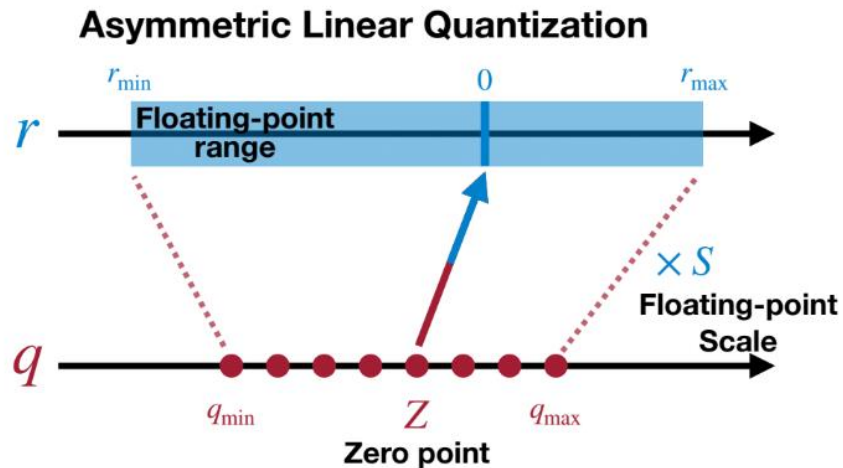


$$S = \frac{r_{\max}}{q_{\max} - Z} = \frac{|r|_{\max}}{q_{\max}} = \frac{|r|_{\max}}{2^{N-1} - 1}$$

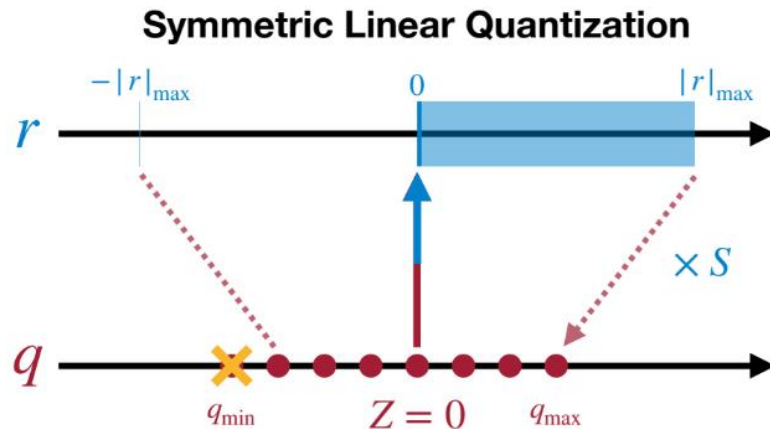
- example: TensorFlow, NVIDIA TensorRT, Intel DNNL



Asymmetric vs. Symmetric



- The quantized range is fully used.
- The implementation is more complex, and zero points require additional logic in hardware.



- The quantized range will be wasted for biased float range.
 - Activation tensor is non-negative after ReLU, and thus symmetric quantization will lose 1 bit effectively.
- The implementation is much simpler.



Linear Quantized Matrix Multiplication

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

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

$$S_{\mathbf{Y}} (\mathbf{q}_{\mathbf{Y}} - Z_{\mathbf{Y}}) = S_{\mathbf{W}} (\mathbf{q}_{\mathbf{W}} - Z_{\mathbf{W}}) \cdot S_{\mathbf{X}} (\mathbf{q}_{\mathbf{X}} - Z_{\mathbf{X}})$$

$$\mathbf{q}_{\mathbf{Y}} = \frac{S_{\mathbf{W}}S_{\mathbf{X}}}{S_{\mathbf{Y}}} (\mathbf{q}_{\mathbf{W}} - Z_{\mathbf{W}}) (\mathbf{q}_{\mathbf{X}} - Z_{\mathbf{X}}) + Z_{\mathbf{Y}}$$

$$\mathbf{q}_{\mathbf{Y}} = \frac{S_{\mathbf{W}}S_{\mathbf{X}}}{S_{\mathbf{Y}}} (\mathbf{q}_{\mathbf{W}}\mathbf{q}_{\mathbf{X}} - Z_{\mathbf{W}}\mathbf{q}_{\mathbf{X}} - Z_{\mathbf{X}}\mathbf{q}_{\mathbf{W}} + Z_{\mathbf{W}}Z_{\mathbf{X}}) + Z_{\mathbf{Y}}$$



Linear Quantized Matrix Multiplication

- An affine mapping of integers to real numbers ($r = S(q - Z)$)
 - Consider the following matrix multiplication

$$Y = WX$$

$$q_Y = \frac{S_W S_X}{S_Y} \left(q_W q_X - Z_W q_X - Z_X q_W + Z_W Z_X \right) + Z_Y$$

Rescale to
N-bit Integer
N-bit Integer Multiplication
32-bit Integer Addition/Subtraction
N-bit Integer
Addition



Linear Quantized Matrix Multiplication

- Consider the following matrix multiplication.

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

$$\mathbf{q}_Y = \frac{S_W S_X}{S_Y} \left(\mathbf{q}_W \mathbf{q}_X - Z_W \mathbf{q}_X - Z_X \mathbf{q}_W + Z_W Z_X \right) + Z_Y$$

- Empirically, the scale $\frac{S_W S_X}{S_Y}$ is always in the interval (0, 1).

Fixed-point Multiplication

$$\frac{S_W S_X}{S_Y} = 2^{-n} M_0, \text{ where } M_0 \in [0.5, 1)$$

Bit Shift



Linear Quantized Matrix Multiplication

$$Y = WX$$

$$q_Y = \frac{S_W S_X}{S_Y} \left(q_W q_X - Z_W q_X - Z_X q_W + Z_W Z_X \right) + Z_Y$$

Rescale to N-bit Integer
**N-bit Integer Multiplication
32-bit Integer Addition/Subtraction**
N-bit Integer Addition

$$q_Y = \frac{S_W S_X}{S_Y} \left(q_W q_X - Z_X q_W \right) + Z_Y$$

$Z_W = 0$



Linear Quantized Fully-Connected Layer

- An affine mapping of integers to real numbers ($r = S(q - Z)$)
 - Now, we consider the following fully-connected layer with bias

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

$$S_Y (\mathbf{q}_Y - Z_Y) = S_W (\mathbf{q}_W - Z_W) \cdot S_X (\mathbf{q}_X - Z_X) + S_b (\mathbf{q}_b - Z_b)$$

↓ $Z_W = 0$

$$S_Y (\mathbf{q}_Y - Z_Y) = \underbrace{S_W S_X}_{\text{blue box}} (\mathbf{q}_W \mathbf{q}_X - Z_X \mathbf{q}_W) + \underbrace{S_b}_{\text{blue box}} (\mathbf{q}_b - \underbrace{Z_b}_{\text{cyan box}})$$



Linear Quantized Fully-Connected Layer

- An affine mapping of integers to real numbers ($r = S(q - Z)$)
 - Now, we consider the following fully-connected layer with bias

$$\mathbf{Y} = \mathbf{WX} + \mathbf{b}$$

$$S_Y (\mathbf{q}_Y - Z_Y) = S_W (\mathbf{q}_W - Z_W) \cdot S_X (\mathbf{q}_X - Z_X) + S_b (\mathbf{q}_b - Z_b)$$

$\downarrow Z_W = 0$

$$S_Y (\mathbf{q}_Y - Z_Y) = S_W S_X (\mathbf{q}_W \mathbf{q}_X - Z_X \mathbf{q}_W) + S_b (\mathbf{q}_b - Z_b)$$

$\downarrow Z_b = 0, S_b = S_W S_X$

$$S_Y (\mathbf{q}_Y - Z_Y) = S_W S_X (\mathbf{q}_W \mathbf{q}_X - Z_X \mathbf{q}_W + \mathbf{q}_b)$$



Linear Quantized Fully-Connected Layer

- An affine mapping of integers to real numbers ($r = S(q - Z)$)
 - Now, we consider the following fully-connected layer with bias

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

$$Z_{\mathbf{W}} = 0 \downarrow Z_{\mathbf{b}} = 0, S_{\mathbf{b}} = S_{\mathbf{W}}S_{\mathbf{X}}$$

$$S_{\mathbf{Y}} (\mathbf{q}_{\mathbf{Y}} - Z_{\mathbf{Y}}) = S_{\mathbf{W}}S_{\mathbf{X}} (\mathbf{q}_{\mathbf{W}}\mathbf{q}_{\mathbf{X}} - Z_{\mathbf{X}}\mathbf{q}_{\mathbf{W}} + \mathbf{q}_{\mathbf{b}})$$

$$\mathbf{q}_{\mathbf{Y}} = \frac{S_{\mathbf{W}}S_{\mathbf{X}}}{S_{\mathbf{Y}}} (\mathbf{q}_{\mathbf{W}}\mathbf{q}_{\mathbf{X}} + \mathbf{q}_{\mathbf{b}} - Z_{\mathbf{X}}\mathbf{q}_{\mathbf{W}}) + Z_{\mathbf{Y}}$$

$$\downarrow \mathbf{q}_{bias} = \mathbf{q}_{\mathbf{b}} - Z_{\mathbf{X}}\mathbf{q}_{\mathbf{W}}$$

$$\mathbf{q}_{\mathbf{Y}} = \frac{S_{\mathbf{W}}S_{\mathbf{X}}}{S_{\mathbf{Y}}} (\mathbf{q}_{\mathbf{W}}\mathbf{q}_{\mathbf{X}} + \mathbf{q}_{bias}) + Z_{\mathbf{Y}}$$



Linear Quantized Fully-Connected Layer

- An affine mapping of integers to real numbers ($r = S(q - Z)$)
 - Now, we consider the following fully-connected layer with bias

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

$$\begin{aligned} Z_{\mathbf{W}} &= 0 \\ Z_{\mathbf{b}} &= 0, \quad S_{\mathbf{b}} = S_{\mathbf{W}}S_{\mathbf{X}} \\ \mathbf{q}_{bias} &= \mathbf{q}_{\mathbf{b}} - Z_{\mathbf{X}}\mathbf{q}_{\mathbf{W}} \end{aligned}$$

$$\mathbf{q}_{\mathbf{Y}} = \frac{S_{\mathbf{W}}S_{\mathbf{X}}}{S_{\mathbf{Y}}} \left(\mathbf{q}_{\mathbf{W}}\mathbf{q}_{\mathbf{X}} + \mathbf{q}_{bias} \right) + Z_{\mathbf{Y}}$$

Rescale to
N-bit Int Mult.
N-bit Int
N-bit Int
32-bit Int Add.
Add

Note: both $q_{\mathbf{b}}$ and q_{bias} are 32 bits.

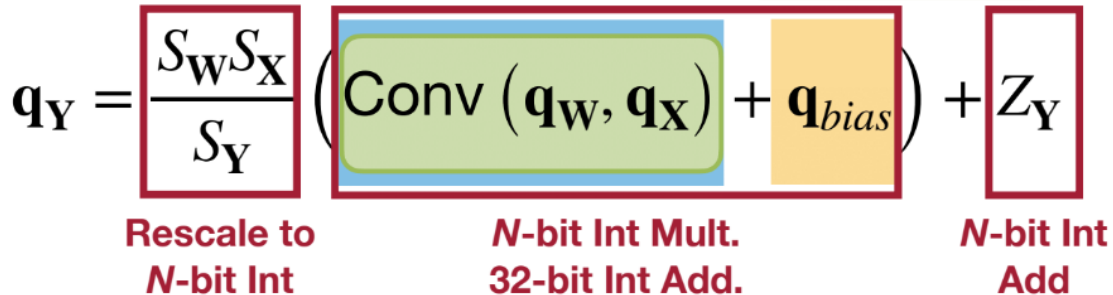


Linear Quantized Convolution Layer

- An affine mapping of integers to real numbers ($r = S(q - Z)$)
 - Now, we consider the following convolution layer

$$\mathbf{Y} = \text{Conv}(\mathbf{W}, \mathbf{X}) + \mathbf{b}$$

$$\begin{aligned} Z_{\mathbf{W}} &= 0 \\ Z_{\mathbf{b}} &= 0, \quad S_{\mathbf{b}} = S_{\mathbf{W}}S_{\mathbf{X}} \\ \mathbf{q}_{bias} &= \mathbf{q}_{\mathbf{b}} - \text{Conv}(\mathbf{q}_{\mathbf{W}}, Z_{\mathbf{X}}) \end{aligned}$$



Note: both $\mathbf{q}_{\mathbf{b}}$ and \mathbf{q}_{bias} are 32 bits.



Linear Quantized Convolution Layer

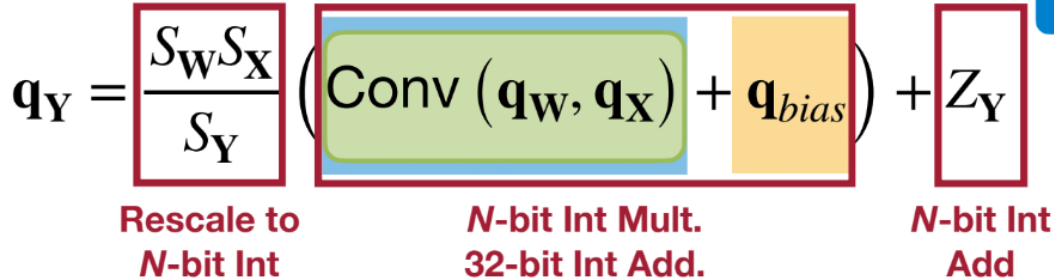
- An affine mapping of integers to real numbers ($r = S(q - Z)$)
 - Now, we consider the following convolution layer

$$Y = \text{Conv}(W, X) + b$$

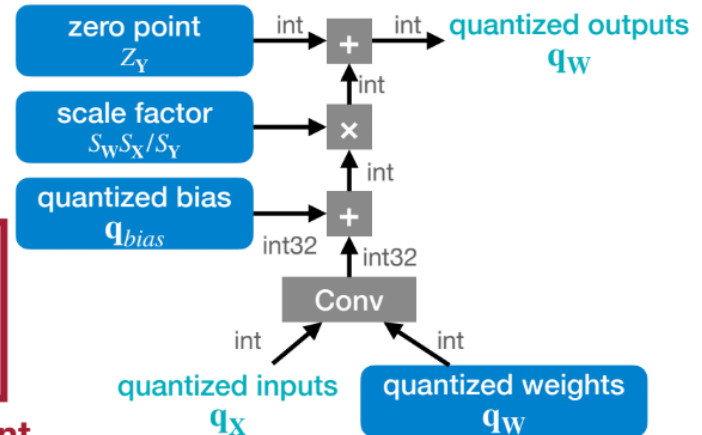
$$Z_W = 0$$

$$Z_b = 0, \quad S_b = S_W S_X$$

$$q_{bias} = q_b - \text{Conv}(q_W, Z_X)$$



Note: both q_b and q_{bias} are 32 bits.



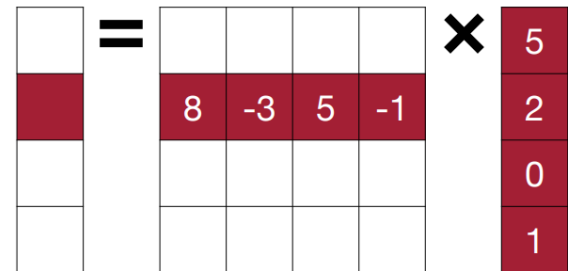
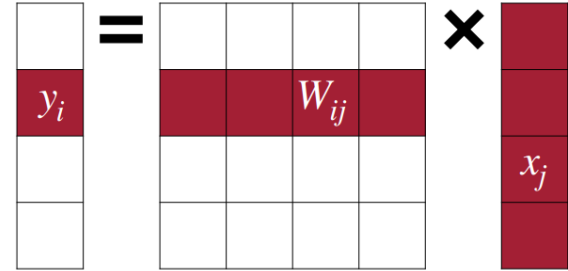


Binary/Ternary Quantization

- Could we push the quantization precision to 1 bit?

$$y_i = \sum_j W_{ij} \cdot x_j$$

$$= 8 \times 5 + (-3) \times 2 + 5 \times 0 + (-1) \times 1$$



input	weight	operations	memory	computation
\mathbb{R}	\mathbb{R}	$+ \times$	1x	1x

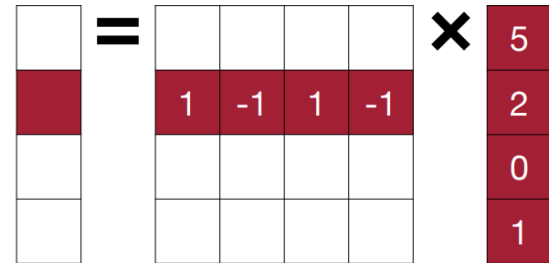
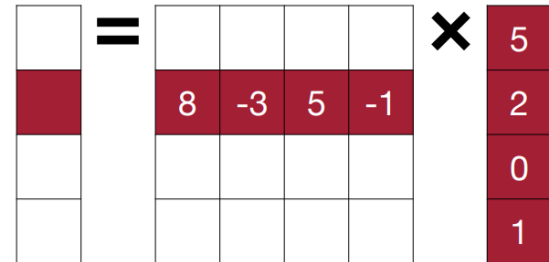


Binary/Ternary Quantization

- If weights are quantized to +1 and -1

$$y_i = \sum_j W_{ij} \cdot x_j$$

$$= 5 - 2 + 0 - 1$$



input	weight	operations	memory	computation
\mathbb{R}	\mathbb{R}	+ ×	1×	1×
\mathbb{R}	\mathbb{B}	+ -	~32× less	~2× less



Binarization

- **Deterministic Binarization**

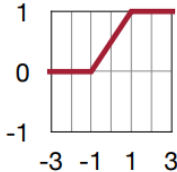
- directly computes the bit value based on a threshold, usually 0, resulting in a sign function.

$$q = \text{sign}(r) = \begin{cases} +1, & r \geq 0 \\ -1, & r < 0 \end{cases}$$

- **Stochastic Binarization**

- use global statistics or the value of input data to determine the probability of being -1 or +1
 - e.g., in Binary Connect (BC), probability is determined by hard sigmoid function $\sigma(r)$

$$q = \begin{cases} +1, & \text{with probability } p = \sigma(r) \\ -1, & \text{with probability } 1 - p \end{cases}, \quad \text{where } \sigma(r) = \min(\max(\frac{r+1}{2}, 0), 1)$$



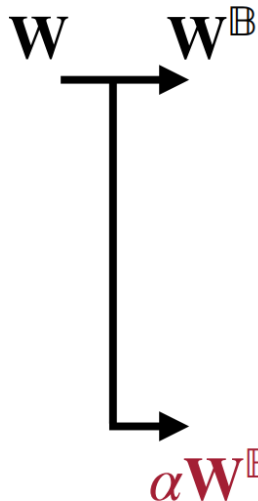
- harder to implement as it requires the hardware to generate random bits when quantizing.



Minimizing Quantization Error in Binarization

weights
(32-bit float)

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



binary weights
(1-bit)

1	-1	1	1
1	-1	-1	1
-1	1	1	-1
1	1	1	1

1	-1	1	1
1	-1	-1	1
-1	1	1	-1
1	1	1	1

$$\mathbf{W}^{\mathbb{B}} = \text{sign}(\mathbf{W})$$

$$\alpha = \frac{1}{n} \|\mathbf{W}\|_1$$

$$\|\mathbf{W} - \mathbf{W}^{\mathbb{B}}\|_F^2 = 9.28$$

scale
(32-bit float)

$$\times 1.05 = \frac{1}{16} \|\mathbf{W}\|_1$$

$$\|\mathbf{W} - \alpha \mathbf{W}^{\mathbb{B}}\|_F^2 = 9.24$$

AlexNet-based Network	ImageNet Top-1 Accuracy Delta
BinaryConnect	-21.2%
Binary Weight Network (BWN)	0.2%



Binary Net

- **Binary Connect**

- Weights $\{-1, 1\}$ (Bipolar binary),
Activation 32-bit float
- Accuracy loss: 19 % on AlexNet

- **Binarized Neural Networks**

- Weights $\{-1, 1\}$, Activations $\{-1, 1\}$
- Both of operands are binary, the multiplication turns into an XNOR
- Accuracy loss: 29.8 % on AlexNet

XNOR

A	B	Out
0	0	1
1	0	0
0	1	0
1	1	1

Popcount (110010001) = 4

for each i in width:

$C += A[\text{row}][i] * B[i][\text{col}]$



for each i in width:

$C += \text{popcount}(\text{XNOR}(A[\text{row}][i] * B[i][\text{col}])))$



Case Study: Binary Multiplication

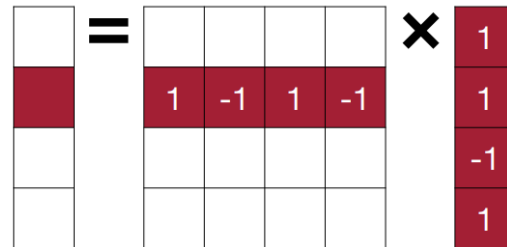
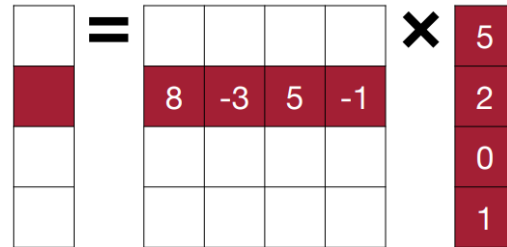
- $A = 10010$, $B = 01111$ (0 is really -1 here)
- **Dot product:**
 - $A * B = (1 * -1) + (-1 * 1) + (-1 * 1) + (1 * 1) + (-1 * 1) = -3$
- $P = \text{XNOR}(A, B) = 00010$, $\text{popcount}(P) = 1$
- $\text{Result} = 2 * P - N$, where N is the total number of bits
- $2 * P - N = 2 * 1 - 5 = -3$



XNOR-Net

- If both activations and weights are binarized

$$\begin{aligned}y_i &= \sum_j W_{ij} \cdot x_j \\ &= 1 \times 1 + (-1) \times 1 + 1 \times (-1) + (-1) \times 1 \\ &= 1 + (-1) + (-1) + (-1) = -2\end{aligned}$$





XNOR-Net

- If both activations and weights are binarized

$$y_i = \sum_j W_{ij} \cdot x_j$$
$$= 1 \times 1 + (-1) \times 1 + 1 \times (-1) + (-1) \times 1$$
$$= 1 + (-1) + (-1) + (-1) = -2$$

W	X	Y=W _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



XNOR-Net

- If both activations and weights are binarized

$$y_i = \sum_j W_{ij} \cdot x_j$$

$$= 1 \times 1 + (-1) \times 1 + 1 \times (-1) + (-1) \times 1$$

$$= 1 + (-1) + (-1) + (-1) = -2$$

$$= 1 \text{ xnor } 1 + 0 \text{ xnor } 1 + 1 \text{ xnor } 0 + 0 \text{ xnor } 1$$

$$= 1 + 0 + 0 + 0 = 1$$



W	X	Y=W _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



XNOR-Net

- If both activations and weights are binarized

$$y_i = \sum_j W_{ij} \cdot x_j$$

$$= 1 \times 1 + (-1) \times 1 + 1 \times (-1) + (-1) \times 1$$

$$= 1 + (-1) + (-1) + (-1) = -2$$

$$y_i = -n + 2 \cdot \sum_j W_{ij} \text{ xnor } x_j$$

$$= 1 \text{ xnor } 1 + 0 \text{ xnor } 1 + 1 \text{ xnor } 0 + 0 \text{ xnor } 1$$

$$= 1 + 0 + 0 + 0 = \boxed{1 \times 2 + (-4)} = -2$$

↑+2
Assuming -1 -1 -1 -1 →

W	X	Y=W _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



XNOR-Net

- If both activations and weights are binarized

$$y_i = -n + 2 \cdot \sum_j W_{ij} \text{xnor } x_j \quad \rightarrow \quad y_i = -n + \text{popcount}(W_i \text{ xnor } x) \ll 1$$
$$= -4 + 2 \times (1 \text{ xnor } 1 + 0 \text{ xnor } 1 + 1 \text{ xnor } 0 + 0 \text{ xnor } 1)$$
$$= -4 + 2 \times (1 + 0 + 0 + 0) = -2$$

→ popcount: return the number of 1

W	X	Y=WX
1	1	1
1	-1	-1
-1	-1	1
-1	1	-1

bw	bx	XNOR(bw, bx)
1	1	1
1	0	0
0	0	1
0	1	0



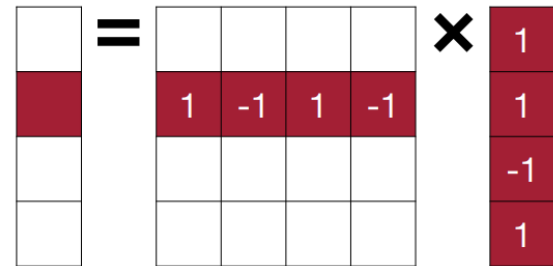
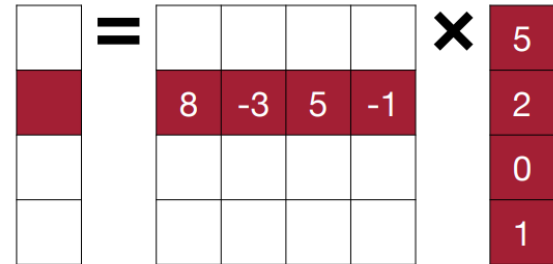
XNOR-Net

- If both activations and weights are binarized

$$y_i = -n + \text{popcount}(W_i \text{ xnor } x) \ll 1$$

$$= -4 + \text{popcount}(1010 \text{ xnor } 1101) \ll 1$$

$$= -4 + \text{popcount}(1000) \ll 1 = -4 + 2 = -2$$

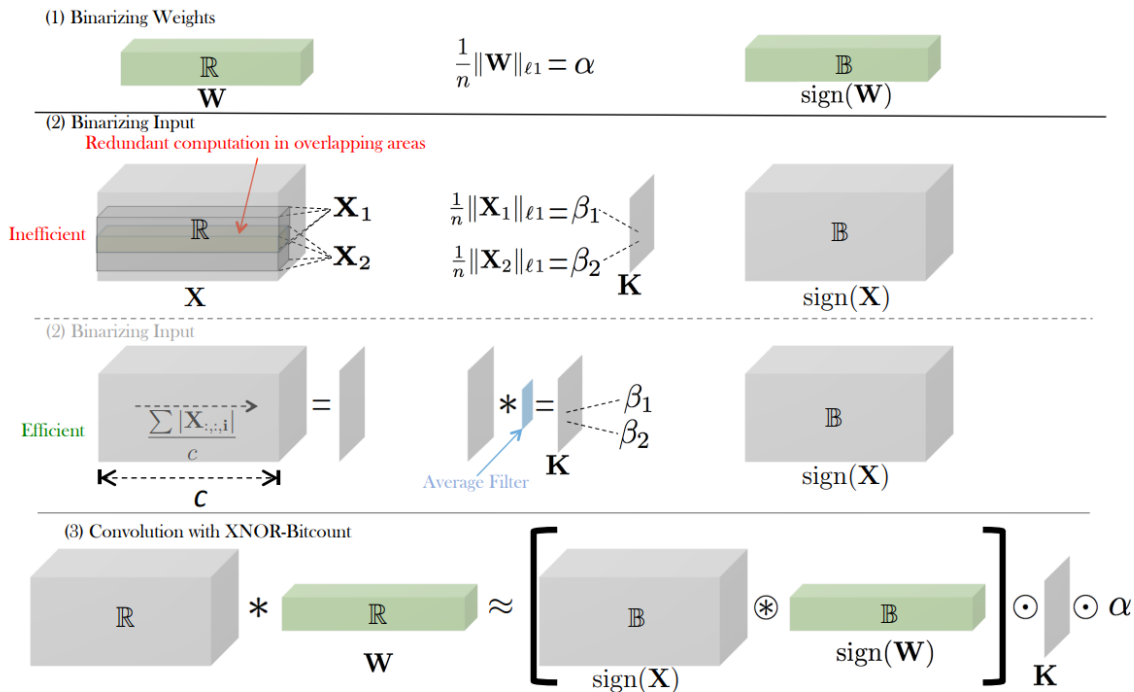


input	weight	operations	memory	computation
R	R	+ ×	1×	1×
R	B	+ -	~32× less	~2× less
B	B	xnor, popcount	~32× less	~58× less



XNOR-Net

- Minimizing quantization error in binarization





XNOR-Net

Neural Network	Quantization	Bit-Width		ImageNet Top-1 Accuracy Delta
		W	A	
AlexNet	BWN	1	32	0.2%
	BNN	1	1	-28.7%
	XNOR-Net	1	1	-12.4%
GoogleNet	BWN	1	32	-5.80%
	BNN	1	1	-24.20%
ResNet-18	BWN	1	32	-8.5%
	XNOR-Net	1	1	-18.1%

- * BWN: Binary Weight Network with scale for weight binarization
- * BNN: Binarized Neural Network without scale factors
- * XNOR-Net: scale factors for both activation and weight binarization



Ternary Weight Networks (TWN)

Weights are quantized to +1, -1 and 0

$$q = \begin{cases} r_t, & r > \Delta \\ 0, & |r| \leq \Delta, \\ -r_t, & r < -\Delta \end{cases} \text{ where } \Delta = 0.7 \times \mathbb{E}(|r|), r_t = \mathbb{E}_{|r| > \Delta}(|r|)$$

weights \mathbf{W}
(32-bit float)

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



ternary weights \mathbf{W}^\top
(2-bit)

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

$$\Delta = 0.7 \times \frac{1}{16} \|\mathbf{W}\|_1 = 0.73$$

$$\times 1.5 = \frac{1}{11} \|\mathbf{W}_{\mathbf{W}^\top \neq 0}\|_1$$

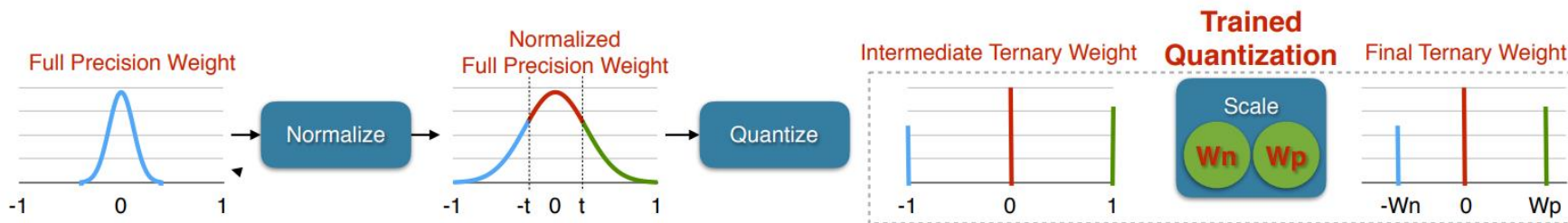
ImageNet Top-1 Accuracy	Full Precision	1 bit (BWN)	2 bit (TWN)
ResNet-18	69.6	60.8	65.3



Ternary Weight Networks (TWN)

- Instead of using fixed scale r_p , TTQ introduces two *trainable* parameters w_p and w_n to represent the positive and negative scales in the quantization.

$$q = \begin{cases} w_p, & r > \Delta \\ 0, & |r| \leq \Delta \\ -w_n, & r < -\Delta \end{cases}$$



ImageNet Top-1 Accuracy	Full Precision	1 bit (BWN)	2 bit (TWN)	TTQ
ResNet-18	69.6	60.8	65.3	66.6



What do we Learn from Quantization?

- **Quantization** can improve DNN computational throughput while maintaining accuracy
- Layers on DNN models can be offered with **different bit widths**
- Varying bit width requires **the support of the hardware**
- **No systematic approach** to figure out the proper bit width in layers of DNN models
- What else ?



Takeaway Questions

- What are purposes of data quantization ?
 - (A) Constrain the value of inputs to a set of discrete values
 - (B) Create more values
 - (C) Improve the degree of parallelism on DNN training
- Why training requires large bit width ?
 - (A) The training needs to compute more data
 - (B) Avoid the value underflow and overflow
 - (C) Gradient and weight update have a larger range