

Generated by Midjourney

QUANTIZATION EXPLAINED

Umar Jamil

Downloaded from: <https://github.com/hkproj/quantization-notes>

License: Creative Commons Attribution-NonCommercial 4.0 International (CC BY-NC 4.0):

<https://creativecommons.org/licenses/by-nc/4.0/legalcode>

Not for commercial use

Outline

- What is quantization
 - Why we need quantization
- Numerical representation
 - Integers
 - Floating-point numbers
- Introduction to quantization
 - Review of how neural networks work
 - Quantization of NN layers
- Types of quantization
 - Asymmetric and symmetric quantization
 - Quantization range
 - Quantization granularity
 - Post-Training Quantization
 - Quantization-Aware Training

Prerequisites

- Basic understanding of neural networks
- High school mathematics

Outline

- What is quantization
 - Why we need quantization
- Numerical representation
 - Integers
 - Floating-point numbers
- Introduction to quantization
 - Review of how neural networks work
 - Quantization of NN layers
- Types of quantization
 - Asymmetric and symmetric quantization
 - Quantization range
 - Quantization granularity
 - Post-Training Quantization
 - Quantization-Aware Training

What is quantization?

The problem

- Most modern deep neural networks are made up of billions of parameters. For example, the smallest LLaMA 2 has 7 billion parameters. If every parameter is 32 bit, then we need $\frac{7 \times 10^9 \times 32}{8 \times 10^9} = 28$ GB just to store the parameters on disk.
- When we inference a model, we need to load all its parameters in the memory, this means big models cannot be loaded easily on a standard PC or a smart phone.
- Just like humans, computers are slow at computing floating-point operations compared to integer operations. Try to do 3×6 and compare it to 1.21×2.897 , which one can you compute faster?

The solution

- Quantization aims to reduce the total amount of bits required to represent each parameter, usually by converting floating-point numbers into integers. This way, a model that normally occupies 10 GB can be “compressed” to less than 1 GB (depending on the type of quantization used). **Please note:** quantization doesn’t mean truncating/rounding. We don’t just round up or down all the floating-point numbers! We will see later how it works
- Quantization can also speed up computation, as working with smaller data types is faster (for example multiplying two integers is faster than multiplying two floating point numbers).

Advantages of quantization

- Less memory consumption when loading models (important for devices like smart phones)
- Less inference time due to simpler data types
- Less energy consumption, because inference takes less computation overall.

Outline

- What is quantization
 - Why we need quantization
- Numerical representation
 - Integers
 - Floating-point numbers
- Introduction to quantization
 - Review of how neural networks work
 - Quantization of NN layers
- Types of quantization
 - Asymmetric and symmetric quantization
 - Quantization range
 - Quantization granularity
 - Post-Training Quantization
 - Quantization-Aware Training

How are integers represented in the CPU (or GPU)?

- Computers use a fixed number of bits to represent any piece of data (a number, a character or a pixel's color). A bit string made up of n bits can represent up to 2^N distinct numbers. For example, with 3 bits, we can represent a total of 2^3 distinct numbers. We usually represent numbers in blocks of 8 bits (byte), 16 bits (short), 32 bits (int) or 64 bits (long).

$$110 = 1 \times 2^2 + 1 \times 2^1 + 0 \times 2^0 = 6 \longrightarrow$$

Binary	Number
000	0
001	1
010	2
011	3
100	4
101	5
110	6
111	7

With 3 bits we can represent 2^3 distinct numbers

- In most CPUs integers are represented using the 2's complement: the first bit indicates the sign, while the rest indicate the absolute value of the number (in case it's positive), or its complement in case it's negative. 2's complement also gives a unique representation to the number zero.

But, but... Python can handle anything!

- Python can represent arbitrary big numbers by using the so-called BigNum arithmetic: each number is stored as an array of digits in base 2^{30} .
 - This is a function of CPython (the Python interpreter), not something built-in in the CPU or the GPU. This means that if we want to perform operations fast, using hardware acceleration provided by CUDA (for example), we are forced to use numbers in fixed format (32 bits usually).

```
Type "help", "copyright", "credits" or "license" for more information.
>>> 2**9999
9975315584403791924418710813417925419117484159430962274260044792647194151109733159599808420180972989496655647116045621357782456747068905587968929660481619789278650233
96897263382623275633029947760275043459096655771254304230309052342754537433044812444045244947419004626970816628925310784154736951278456194032612548321937220523379935813
49272661143426908084715788781482038141844038036611426754582073809197819072948473194970542048026813391053231071366669701826278282476530157134011748470016796715832572964
88866398328878030862910157039970990898036891228418811400186514427436259504172322907273252789648007074169608078672940696285476898845596389004134788678372220615310093789
18162751364161894635355186901433196515714066620700812097835845287030709827171162319400624428073652603715996129805898125065496430120854170403802966160080634246144248127
920656422030768369475743557128157555544872757101656910101465820478798232378005202922920783036224814335082575309603155020932111379543354502873032089284759572802753412
56252030037599211309490296185590272223940364531976212741696109913578223265811883804233065168893530199017065985667468273113502815849687277541208904864054916456572017859
3876238425492863846896321661079969993844333040418441891901382164138758613682878637239206147194866905430803711626645987406560809802089140982848737949082265629217067979
93139206506409270314173832454345260523790441307911980992885061203522165291537934519659802301702486578291604336052956650451876411707769872697198857628726525510615547
36608053767374128703876369931741492491703784689778233193109372847496395082860518506822165679086071558956991134919229236672201354820914255025364638741822752893172505504
26493906194736964349770417173079403521979559492907572889588571809849364065729741891601040737491085929005694535614125452913408718110288737960708826857843862807452291452
4962305143150407677916540650509383792811717176947770458781170042244376308132178432441675973186018864662004722812346162717520033901363691887768820336344931812051874570
5483350278523795490501233949408091359629766906412109770141513797042244775073383341948489984431208181566881969516867279007038183709388555276921128697495550932341098482
90825742565247111849738573815345777341088414381001813886288618906826658055984056403963347409436006493218303842758199302673011489357778758973692623184723461543947132974
108504025506161187274814408451786956068416919679587820936692525548513580695771979549579907732720866815582846801556112496898499961339086617901155593132228764956787908750
40999196181423076249405444801161221810868858090431785077342420293111648964269378117432728202684813110094817855144061807837562716691516350145488343252842785787527583637
594495970648556684507495809065758577200386432528659477872546016509265242355690915770366202665951923104201821088185195577531989450037142683609814045173898726660234184
39793429011897610931456004037140977565897407881222414925923075485244401363736078734406579737520486605754024909522790170841347489357065803160534319575584088715239629835
4688
>>>
```

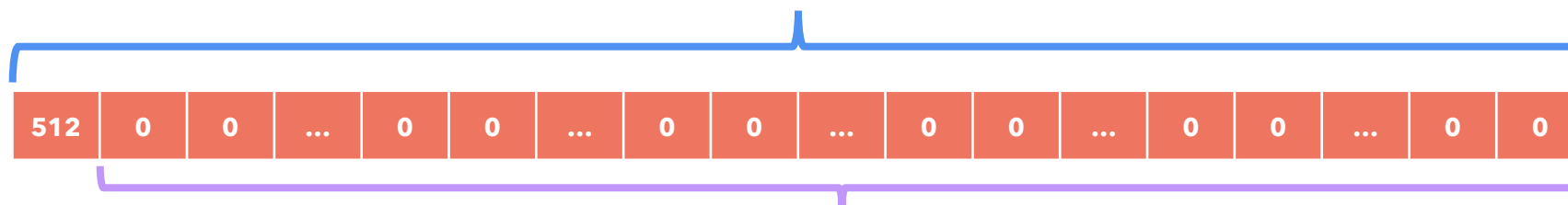


If stored in base 10 digits, it would require an array of 3010 digits

Array of 334 elements, each element is a 32-bit unsigned integer

$$2^{9999} = 512 \times (2^{30})^{333}$$

```
>>> 512*((2**30)**333) == 2**9999
True
```



All elements with index 0... 332 are equal to zero.

How are floating point numbers represented?

Decimal numbers are just numbers that also include negative powers of the base. For example:

$$85.612 = 8 \times 10^1 + 5 \times 10^0 + 6 \times 10^{-1} + 1 \times 10^{-2} + 2 \times 10^{-3}$$

The IEEE-754 standard defines the representation format for floating point numbers in 32 bit.



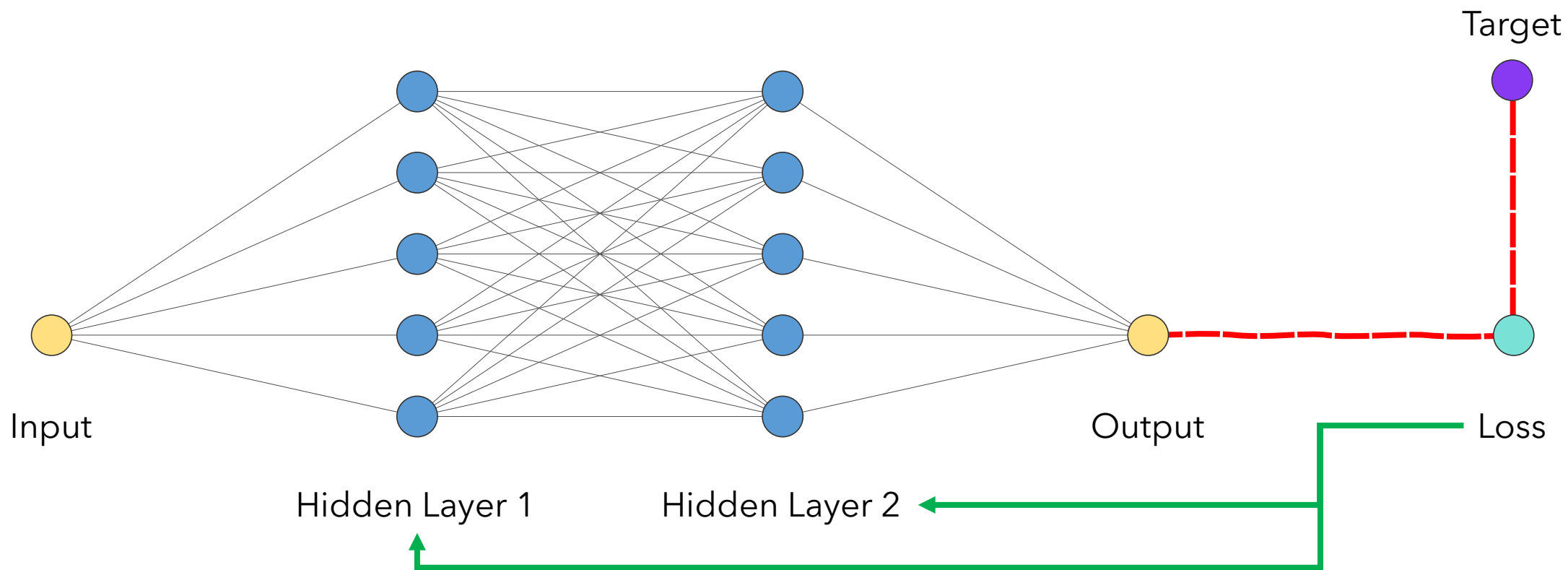
$$\text{Value} = (-1)^{\text{sign}} \times 2^{(E-127)} \times \left(1 + \sum_{i=1}^{23} b_{23-i} 2^{-i}\right) = (+1) \times 2^{-3} \times 1.25 = +0.15625$$

Modern GPUs also support a 16-bit floating point number, with less precision.

Outline

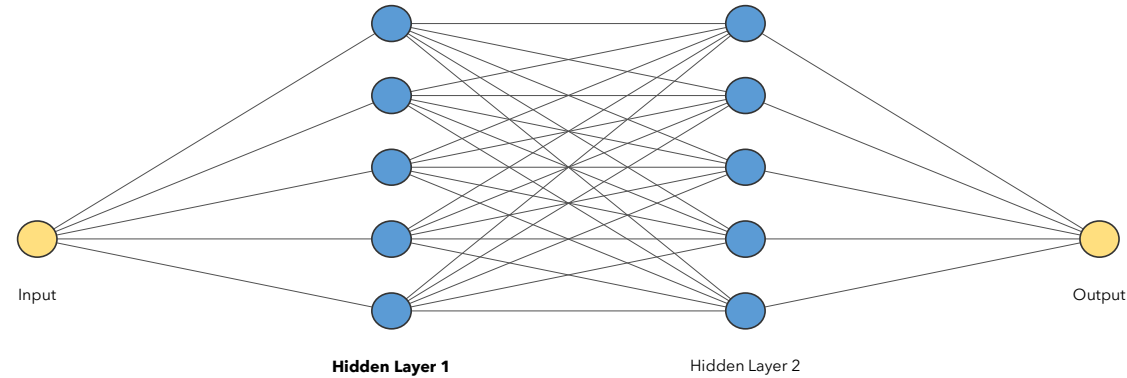
- What is quantization
 - Why we need quantization
- Numerical representation
 - Integers
 - Floating-point numbers
- Introduction to quantization
 - Review of how neural networks work
 - Quantization of NN layers
- Types of quantization
 - Asymmetric and symmetric quantization
 - Quantization range
 - Quantization granularity
 - Post-Training Quantization
 - Quantization-Aware Training

A review of neural networks



- The neural network can be made up of many different layers. For example, linear layers are made up of two matrices, called the **weight** and the **bias**, which are commonly represented using floating-point numbers. Quantization aims to use integer numbers to represent these two matrices, while maintaining the accuracy of the model.

Applying quantization



Dequantize

Output

$$Y = XW + B$$

**Perform all operations
using integer arithmetic**

The main benefit is that integer operations are much faster in most hardware (especially on embedded devices) than floating point operations

Input

Sometimes called "activation"

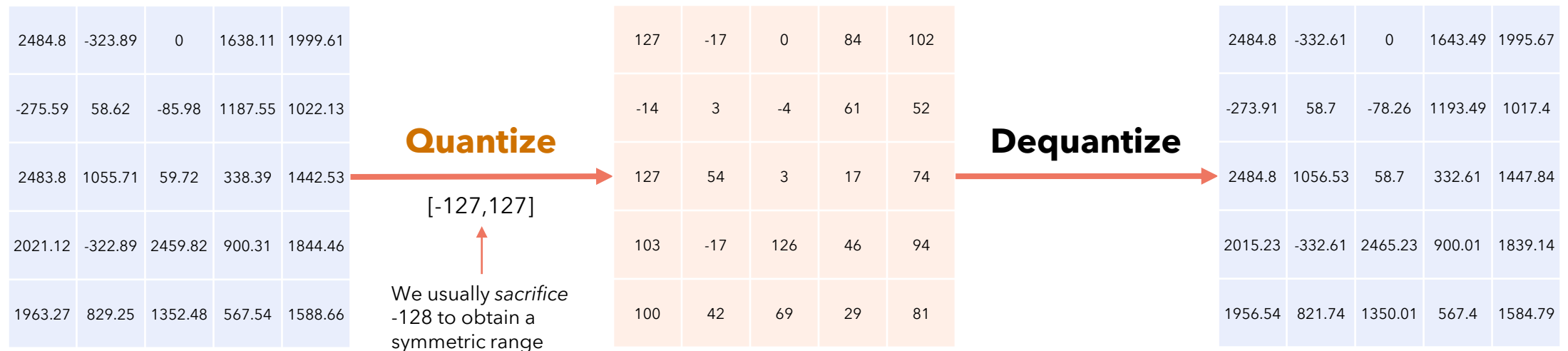
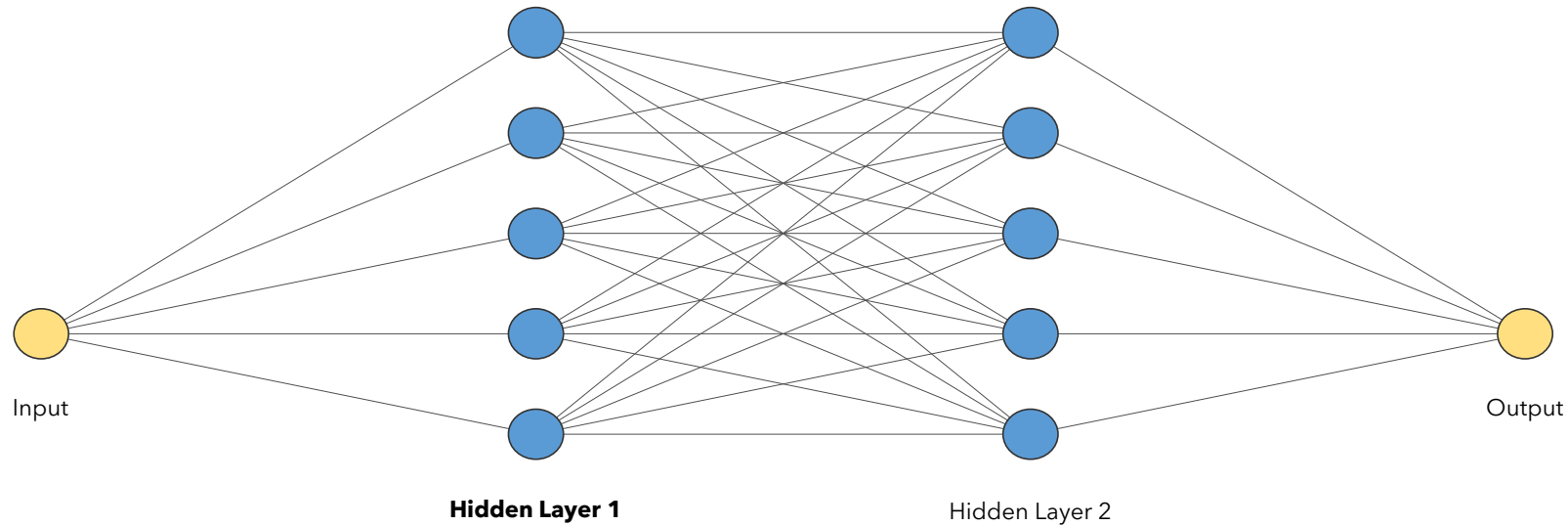
Weight

Bias

Usually quantized as int32

Quantize

Applying quantization



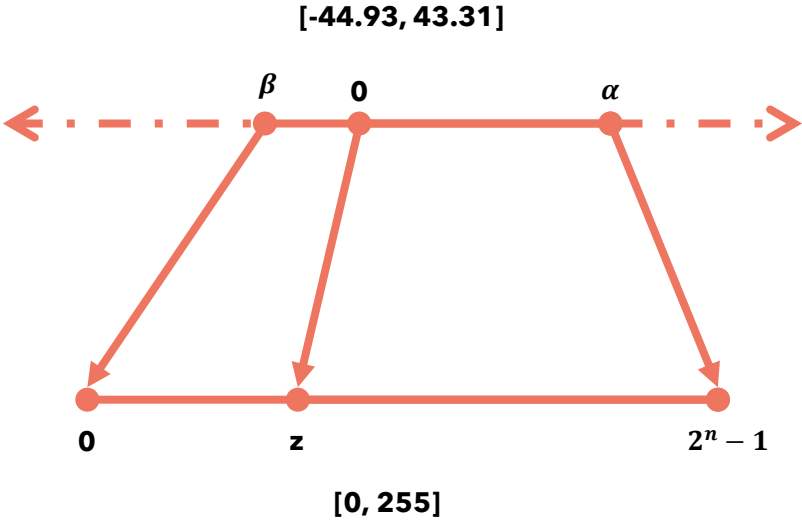
Outline

- What is quantization
 - Why we need quantization
- Numerical representation
 - Integers
 - Floating-point numbers
- Introduction to quantization
 - Review of how neural networks work
 - Quantization of NN layers
- Types of quantization
 - Asymmetric and symmetric quantization
 - Quantization range
 - Quantization granularity
 - Post-Training Quantization
 - Quantization-Aware Training

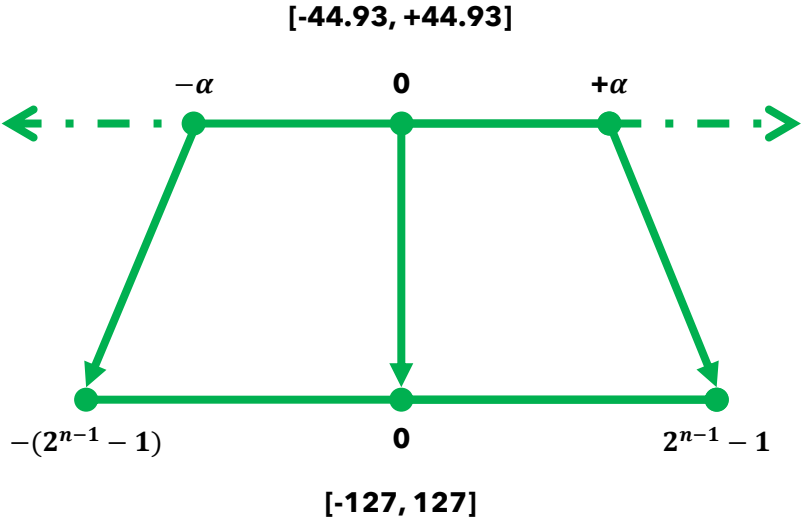
Asymmetric vs Symmetric quantization

43.31	-44.93	0	22.99	-43.93	-11.35	38.48	-20.49	-38.61	-28.02
-------	--------	---	-------	--------	--------	-------	--------	--------	--------

Asymmetric



Symmetric

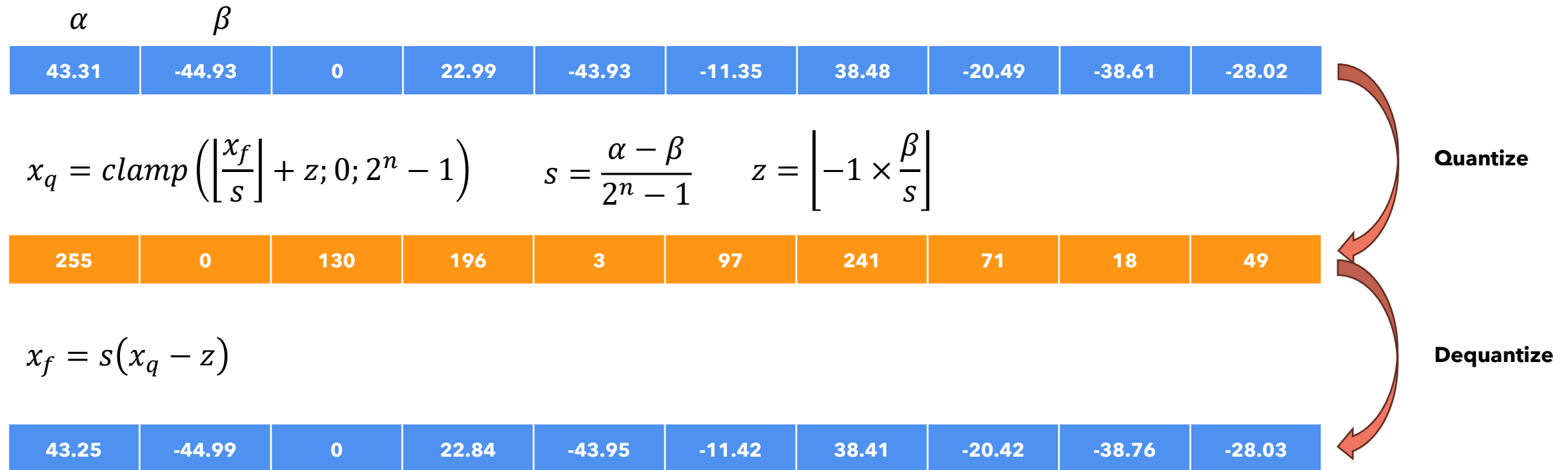


255	0	130	196	3	97	241	71	18	49
-----	---	-----	-----	---	----	-----	----	----	----

122	-127	0	65	-124	-32	109	-58	-109	-79
-----	------	---	----	------	-----	-----	-----	------	-----

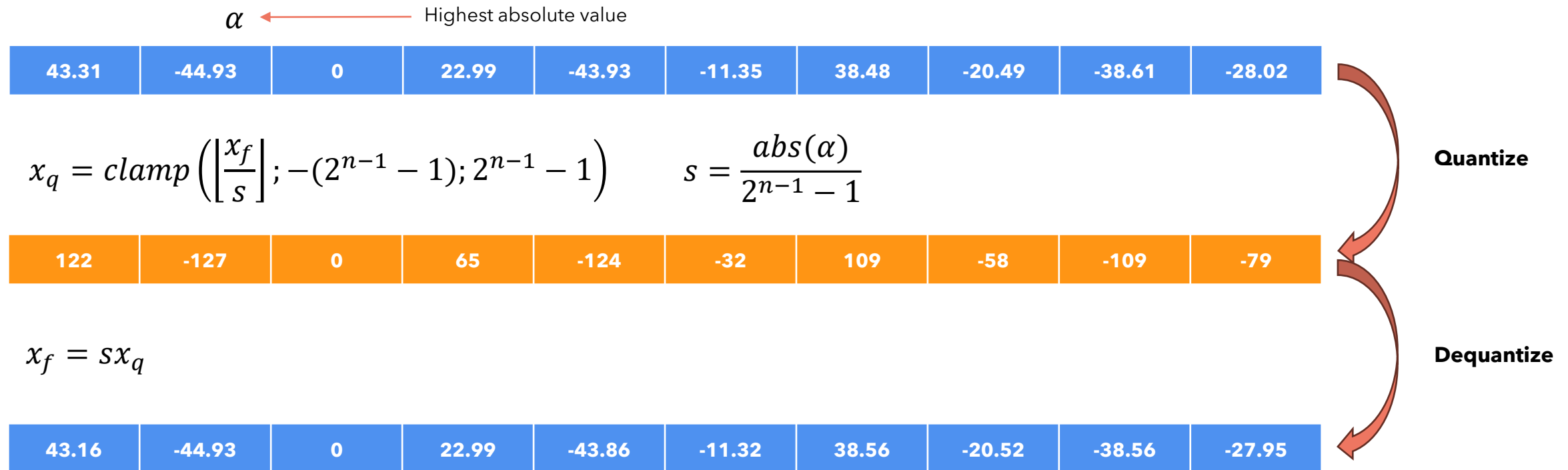
Asymmetric quantization

- It allows to map a series of floating-point numbers in the range $[\beta, \alpha]$ into another in the range $[0, 2^n - 1]$. For example, by using 8 bits, we can represent floating-point numbers in the range $[0, 255]$



Symmetric quantization

- It allows to map a series of floating-point numbers in the range $[-\alpha, \alpha]$ into another in the range $[-(2^{n-1} - 1), 2^{n-1} - 1]$. For example, by using 8 bits, we can represent floating-point numbers in the range $[-127, 127]$

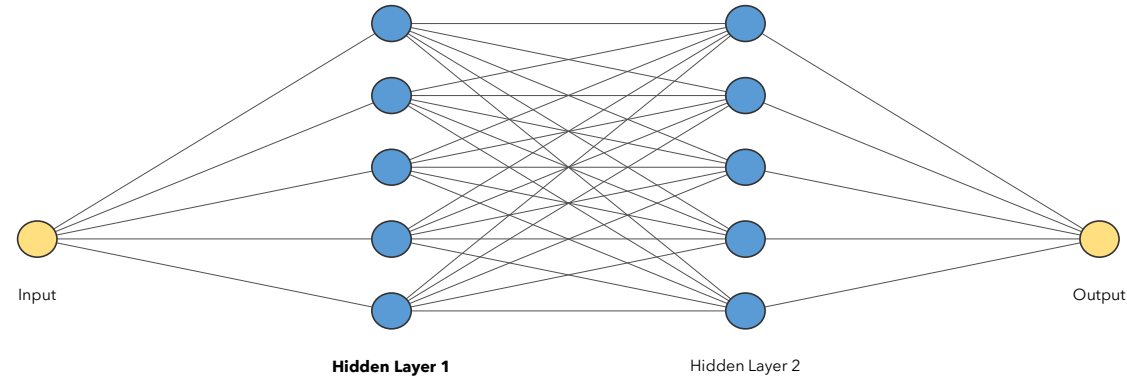




Talk is cheap. Show me the code.

- Linus Torvalds

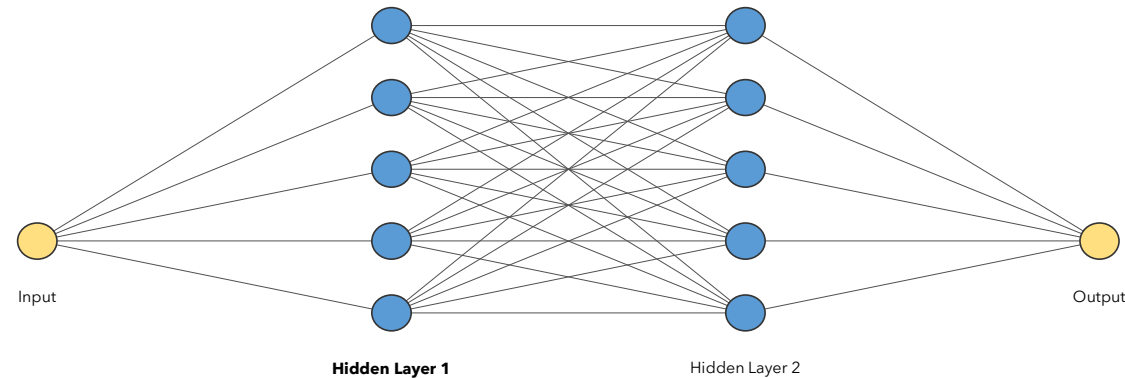
Applying quantization: floating point case



$$\text{Output} \longrightarrow Y_f = XW + B$$

Input \nearrow Weight \nearrow Bias

Applying quantization: integer case



Dequantize

Output

$$Y_q = XW + B$$

**Perform all operations
using integer arithmetic**

How do we dequantize Y?
Y is the result of an operation, how can we dequantize it since we never calculated its scale (s) and zero (z) parameters?

We basically run inference on the model using a few inputs and "observe" the typical output to calculate the scale (s) and the zero (z). This process is called calibration. **We will see it later with Post-Training Quantization.**

Input

Sometimes called "activation".
Can be quantized "on the fly" using a process called (**dynamic quantization**) or with observers.

Weight

Bias

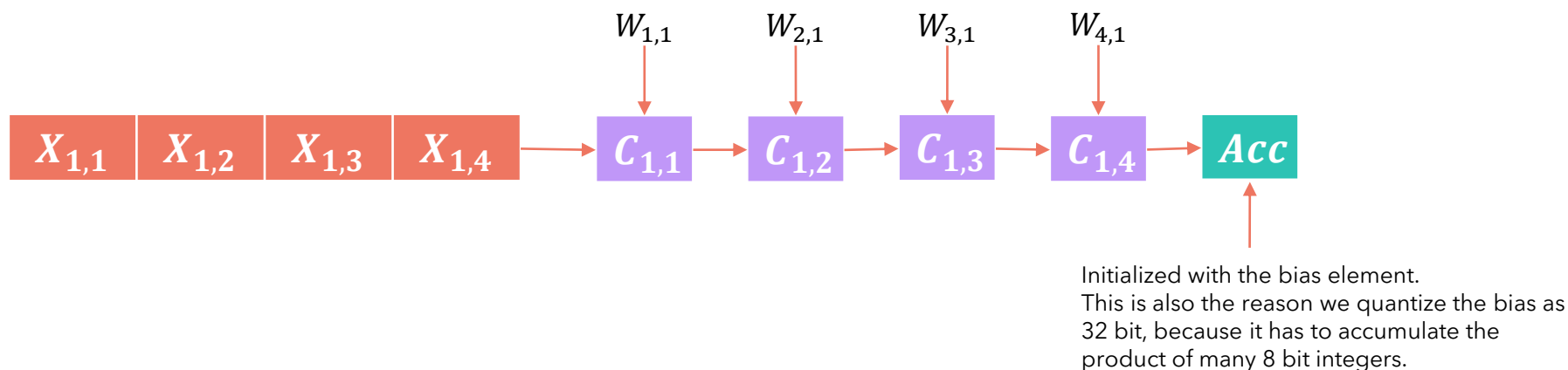
Usually quantized as int32

Quantize

Low-precision matrix multiplication

When we calculate the product $\mathbf{XW} + \mathbf{B}$ in a linear layer, this results in a list of dot products between each row of the matrix \mathbf{X} with each column of the matrix \mathbf{Y} , summing the corresponding element of the bias vector \mathbf{B} .

The GPU can accelerate this operation by using a **Multiply-Accumulate (MAC)** block, which is a physical unit in the GPU that works as follows:

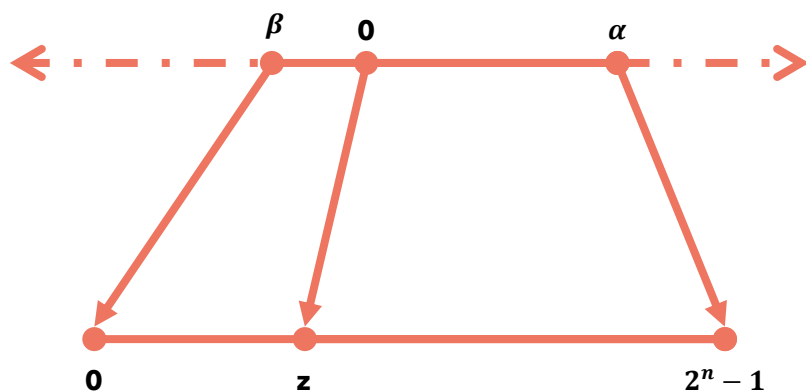


The GPU will perform this operation in parallel for every row and column of the initial matrices using many Multiply-Accumulate (MAC) blocks.

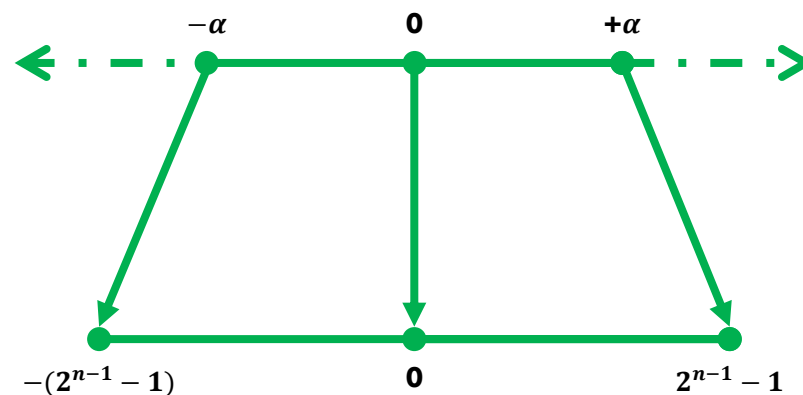
For the full derivation of the math behind low-precision matrix multiplication, Google's GEMM library provides a fantastic explanation: <https://github.com/google/gemmlowp/blob/master/doc/quantization.md>

Asymmetric vs Symmetric quantization

Asymmetric



Symmetric



How do we choose the $[\alpha, \beta]$ range? Any value outside this range will be clipped to the range. We have many strategies to decide it. Let's review them.

Quantization range: how to choose $[\alpha, \beta]$

Min-Max: To cover the whole range of values, we can set

- $\alpha = \max(V)$
- $\beta = \min(V)$
- Sensitive to outliers.

Original										
	43.31	-44.93	0	38.48	-20.49	1000.00	-28.02
Dequantized (Min-Max)	45.08	-45.08	0	24.59	-45.08	-12.29	36.88	-20.49	999.85	-28.68

Outlier



Percentile: Set the range to the percentile of the distribution of V, to reduce sensitivity to outliers

Dequantized (Percentile)	43.38	-44.52	0	38.48	-20.49	50.00	-28.01
--------------------------	-------	--------	---	-----	-----	-----	-------	--------	-------	--------



Only the outlier is quantized with a large error

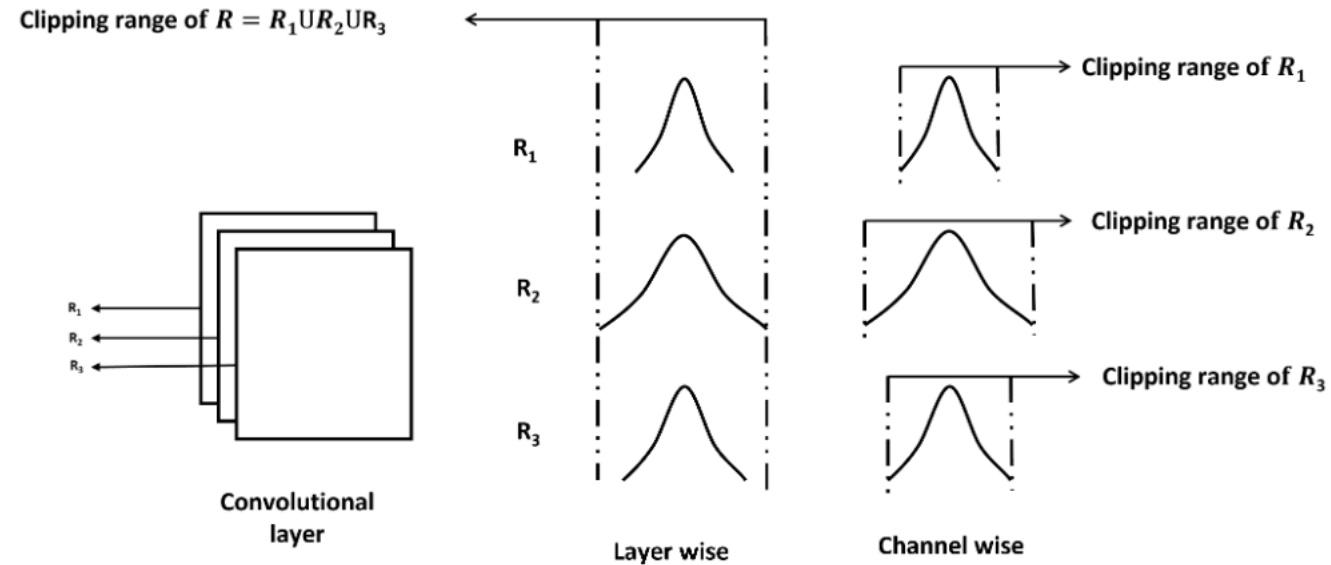
Let's have a look at the code!

Quantization range: how to choose $[\alpha, \beta]$

If the vector V represents the tensor to be quantized, we can choose the $[\alpha, \beta]$ range according to the following strategies:

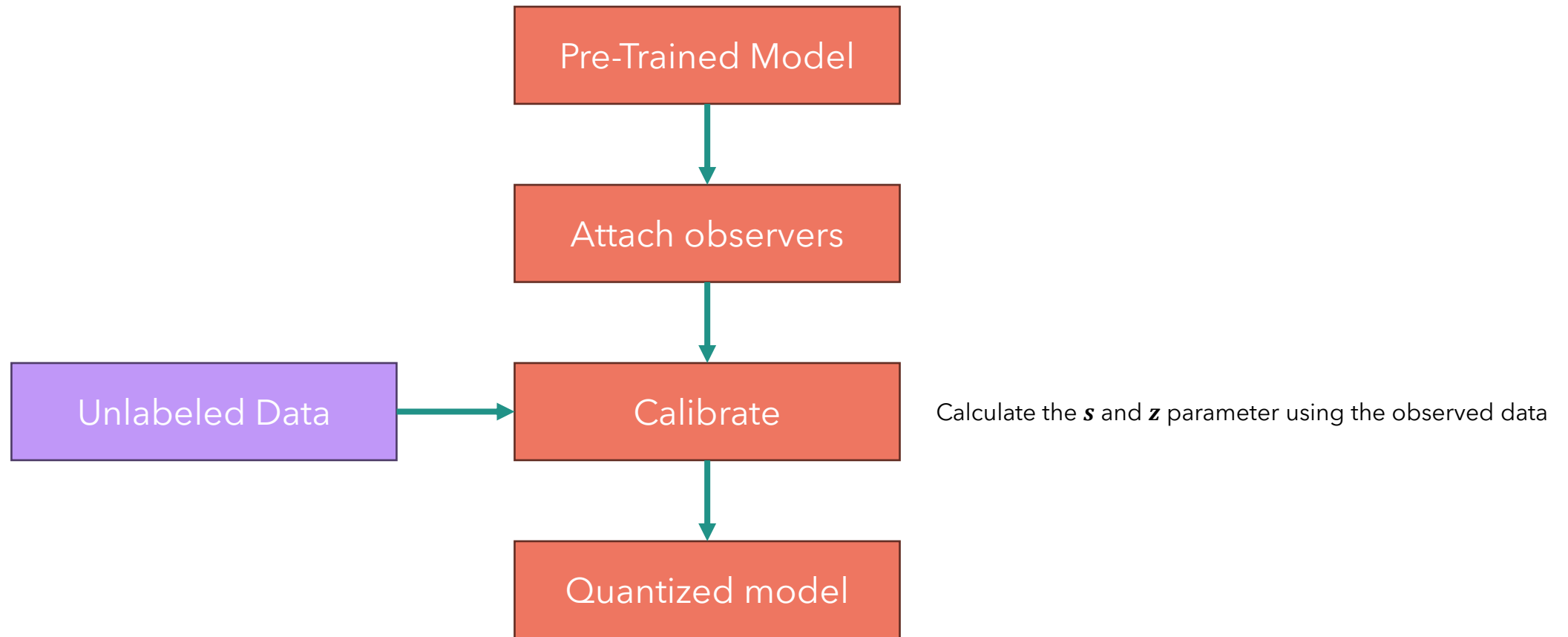
- **Mean-Squared-Error:** choose $[\alpha, \beta]$ such that the MSE error between the original values and the quantized values is minimized.
 - It is usually solved using Grid-Search
- **Cross-Entropy:** used when the values in the tensor being quantized are not equally important. This happens for example in the Softmax layer in Large Language Models. Since most of the inference strategies are Greedy, Top-P or Beam search, it is important to preserve the order of the largest values after quantization.
 - $\operatorname{argmin}_{\alpha, \beta} \operatorname{CrossEntropy}(\operatorname{softmax}(V), \operatorname{softmax}(\hat{V}))$

Quantization granularity



Evaluation of Quantization Methods for Neural Networks, Junaid Mundichipparakkal

Post Training Quantization (PTQ)





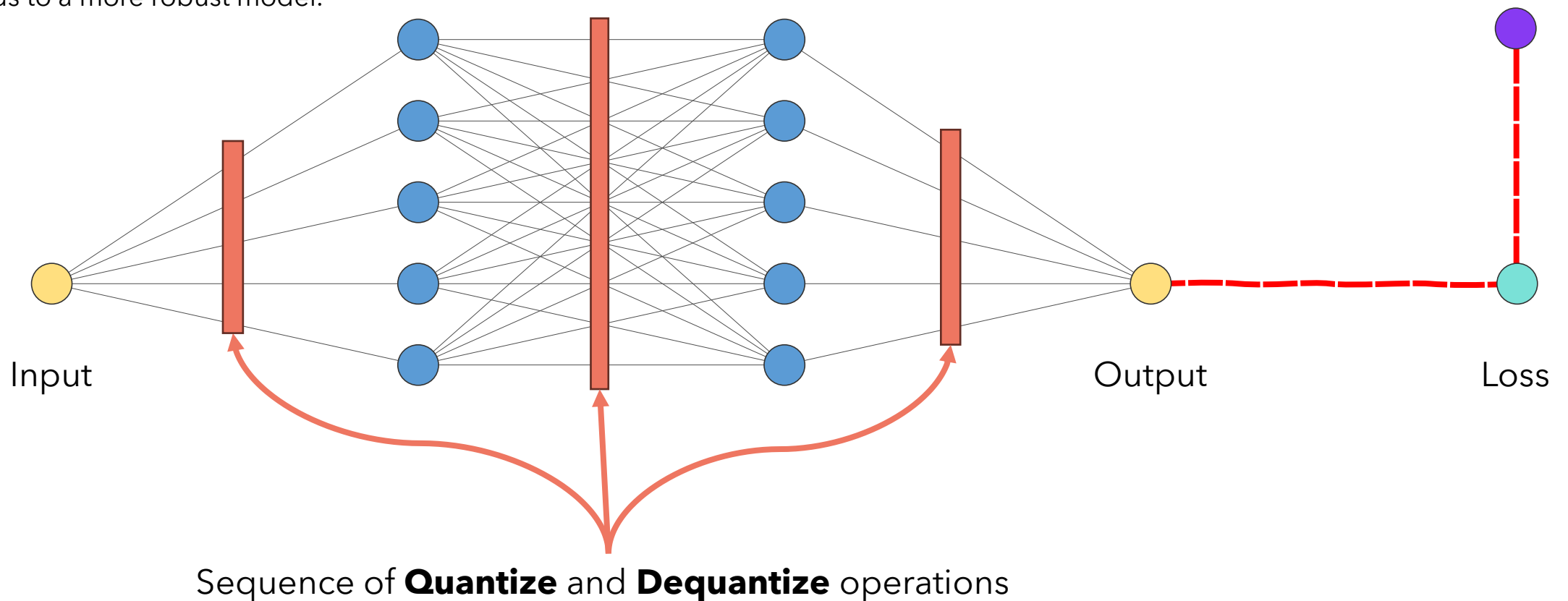
Talk is cheap. Show me the code.

- Linus Torvalds

Quantization Aware Training (QAT)

We insert some fake modules in the computational graph of the model to simulate the effect of the quantization during training.

This way, the loss function gets used to update weights that constantly suffer from the effect of quantization, and it usually leads to a more robust model.





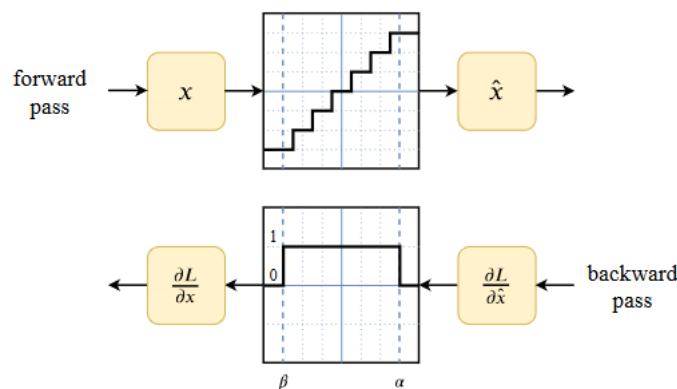
Talk is cheap. Show me the code.

- Linus Torvalds

Quantization Aware Training (QAT): gradient

During backpropagation, the model needs to evaluate the gradient of the loss function w.r.t every weight and input. A problem arises: what is the derivative of the quantization operation we defined before?

A typical solution is to approximate the gradient with the STE (Straight-through Estimator) approximation.

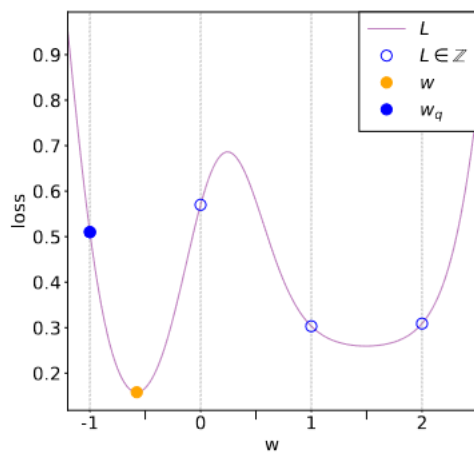


The STE approximation results in 1 if the value being quantized is in the range $[\alpha, \beta]$, otherwise it is 0.

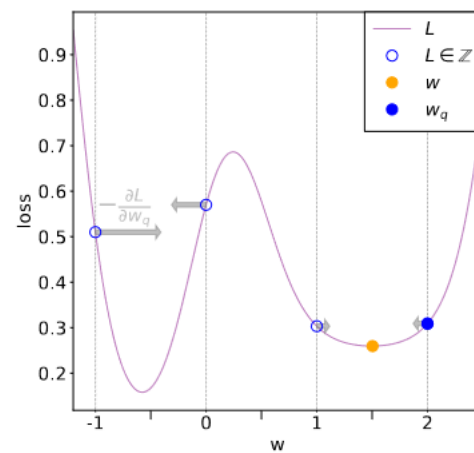
INTEGER QUANTIZATION FOR DEEP LEARNING INFERENCE: PRINCIPLES AND EMPIRICAL EVALUATION, Wu et al.

Quantization Aware Training (QAT): why it works

Why does QAT work? What is the effect of inserting fake quantization operations?



(a) Post training quantization



(b) After quantization aware fine-tuning

INTEGER QUANTIZATION FOR DEEP LEARNING INFERENCE: PRINCIPLES AND EMPIRICAL EVALUATION, Wu et al.

Thanks for watching!
Don't forget to subscribe for
more amazing content on AI
and Machine Learning!