

Fine Tuning LLM Models

① Quantization

- ① Full Precision / Half Precision → related to data types / how the data is stored in the memory.
- ② Calibration - Model Quantization
 - ↳ Problems
- ③ Modes of Quantization
 - ↳ Post Training Quantization
 - ↳ Quantization Aware Training

Quantization : Conversion from higher memory format to a lower memory format.)

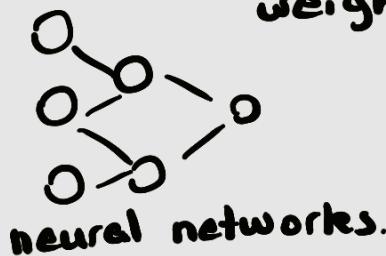
with LLM models, the parameters are a lot and keep increasing. Llama 2 → 70 billion for example.

we can convert this 32 bit to int 8. we can run the model in our resources more easily.

let's say here we have the value 7.23, and this number is stored as 32 bits in the memory. (Full precision/ single precision like floating points.)

weights → matrix

every value are stored in the memory in the form of 32 bits. FP 32



inferencing → giving input and getting output.

quantization helps inferencing.

@ similarly, in Tensorflow data are stored in the TF 32 bit format.

↳ In the scope of LLMs, after Quantization we can also fine tune the models. But this has some problems. When we are converting from 32 bits to lower bits, we may face some data loss and it leads to some accuracy loss. But there are specific techniques to overcome this.

Calibration : How we will be able to convert 32 bit into int 8. It answers "what is the formula required."

* How to perform Quantization? *

① Symmetric Quantization ② Asymmetric Quantization



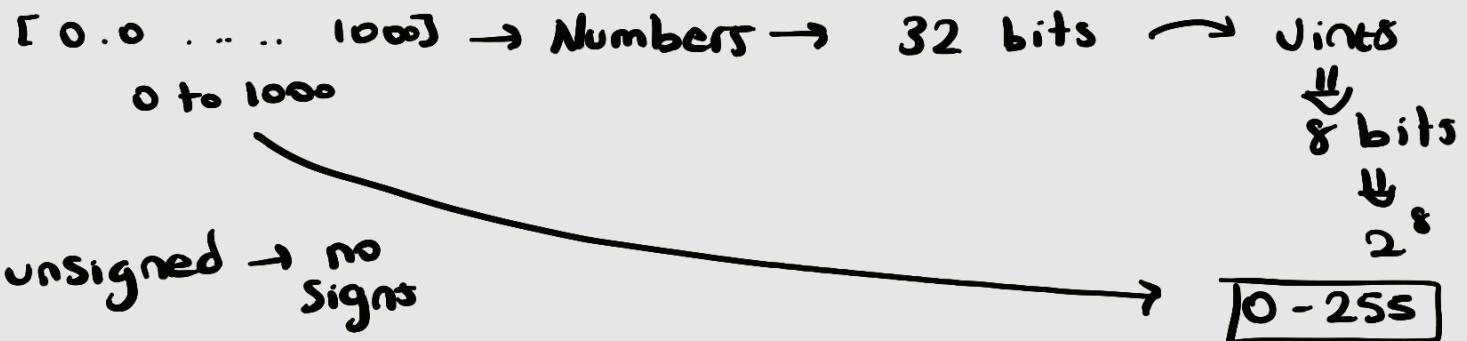
Batch Normalization

From Deep Learning.

technique of Symmetric Quantization.

↳ between all the layers, we apply batch normalization, so that all our weights are zero centered.

(ex) Symmetric Unsigned int 8 quantization.



min max scalar

$0.0 \rightarrow 0$
 $1000 \rightarrow 255$ bits ↓

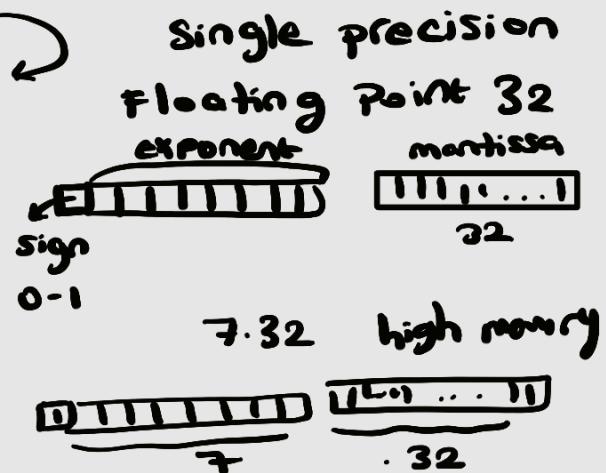
$$\text{scale} = \frac{x_{\max} - x_{\min}}{q_{\max} - q_{\min}}$$

quantization

$$\left(\frac{1000 - 0}{255 - 0} \right) = 3.92 \text{ scale factor.}$$

$$\text{round} \left(\frac{250}{3.92} \right) = 63.77 \rightarrow 64 \text{ symmetric}$$

FP32 → Uint8



② Asymmetric Uint8 quantization

$[-20.0 \dots 1000.0] \rightarrow \text{Numbers} \rightarrow \text{Uint8}$

$[0 \dots 255]$

$$\frac{1000 - (-20)}{255} = 4.0 \text{ scale factor.}$$

but our range is $0 - 255 \rightarrow$ add the same number.
 $\text{round} \left(\frac{-20}{4.0} \right) = -5 \rightarrow -5 + 5 = 0$

add +5 → "zero point"

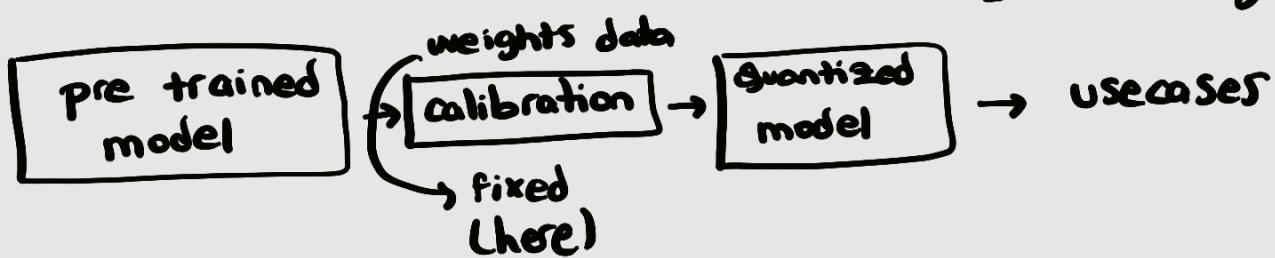
important quantization parameter.

③ if it was signed int8, range would be [-128 ... 127] and then apply the same formula.

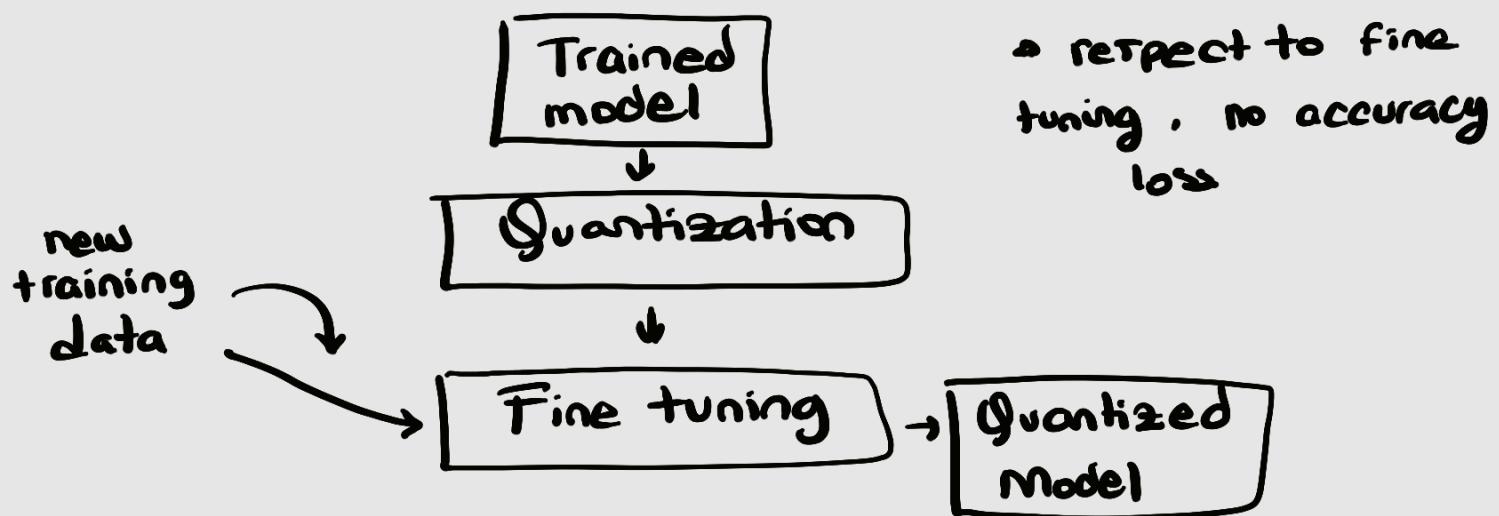
→ the process we apply to the quantization is called calibration.

Modes of Quantization

≈ Post training quantization (PTQ) • there is loss of data
• accuracy decreases



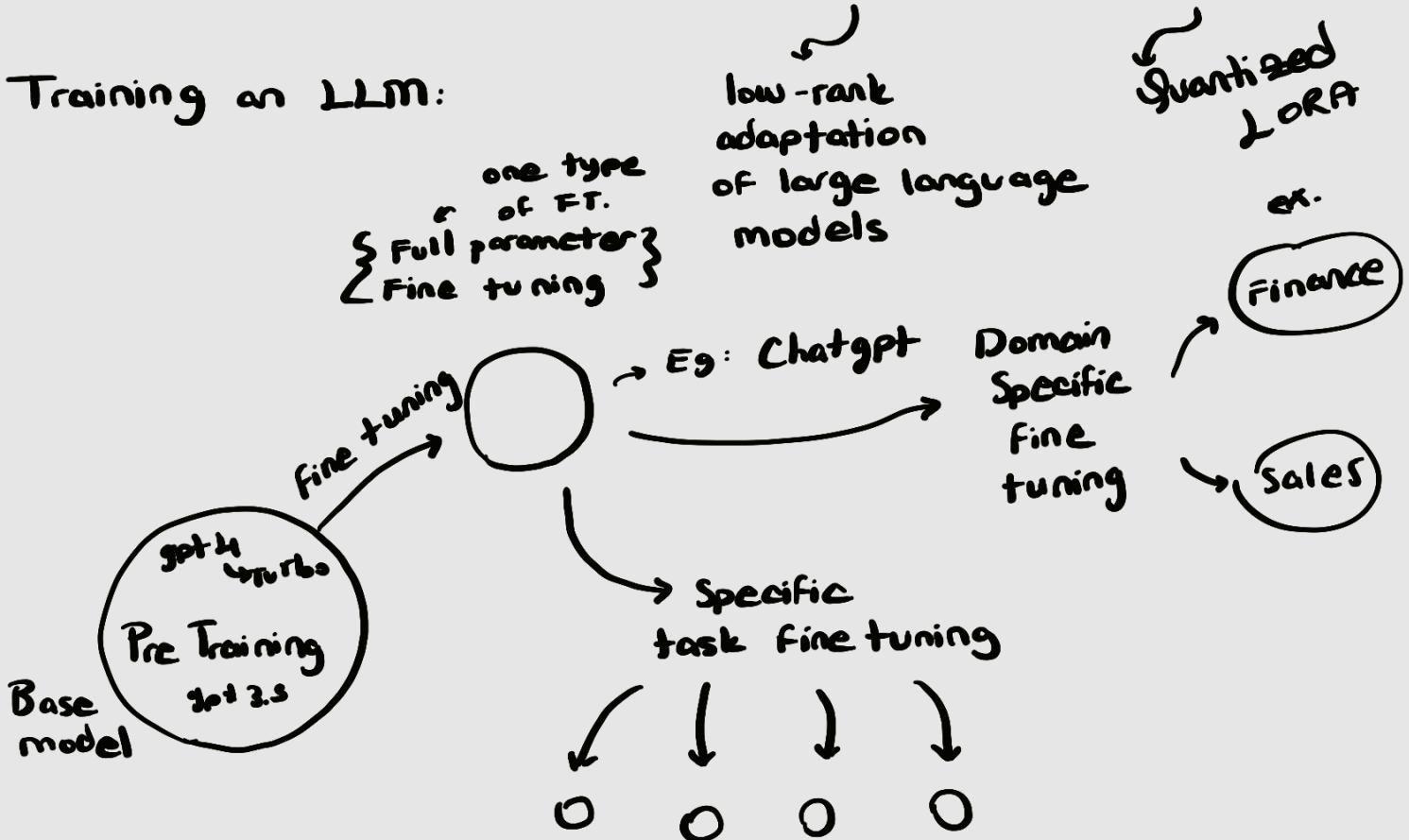
≈ Quantization aware training (QAT)



④ QLora / Lora

Fine Tuning LLM - LORA , QLORA

Training on LLM:



→ pre-trained
llm models (gpt 4 turbo etc)
are base models.

Full Parameter Fine Tuning [challenges]

- ① need to update all model weights.
- ② Hardware Resource Constraint.

Downstream tasks
↳ model monitoring
↳ model inferencing
↳ GPU, RAM constraint
challenges

Because of these challenges
we use LORA, QLORA

What does LORA do?

① Instead of updating all the weights, it will track changes.

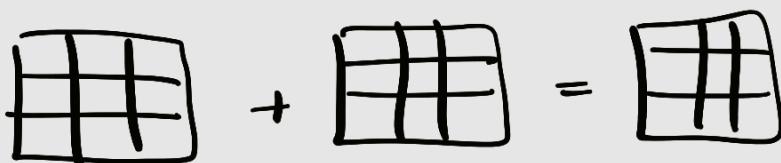


track the changes of the new weights based on fine tuning.

⇒ Lora uses a technique to track changes.

"Matrix Decomposition"

3x3 matrix will be saved on a two smaller matrix.



$$W_0 + \Delta w = W_0 + BA$$

matrix decomposition

$$\begin{matrix} B \\ \times \\ 3 \times 1 \end{matrix} \quad \begin{matrix} A \\ \times \\ 1 \times 3 \end{matrix}$$

Rank = 1

check online
for calculation