



Fine Tuning LLM models

- ① Quantization - is a common technique used to reduce the model size, though it can sometimes result in reduced accuracy.
- ② Quantization-aware training is a method that allows practitioners to apply quantization techniques without sacrificing accuracy.
- ③ It is done in model training process rather than after the fact. The model size can typically be reduced by two to four times & sometimes even more.

④ modes of Quantization → ① Post Training Quantization.

② Quantization aware Training

definition - Quantization: conversion from higher memory format to a lower memory format.

32 bits → int 8 inference

Full precision → half precision, loss of data weight & parameter accuracy

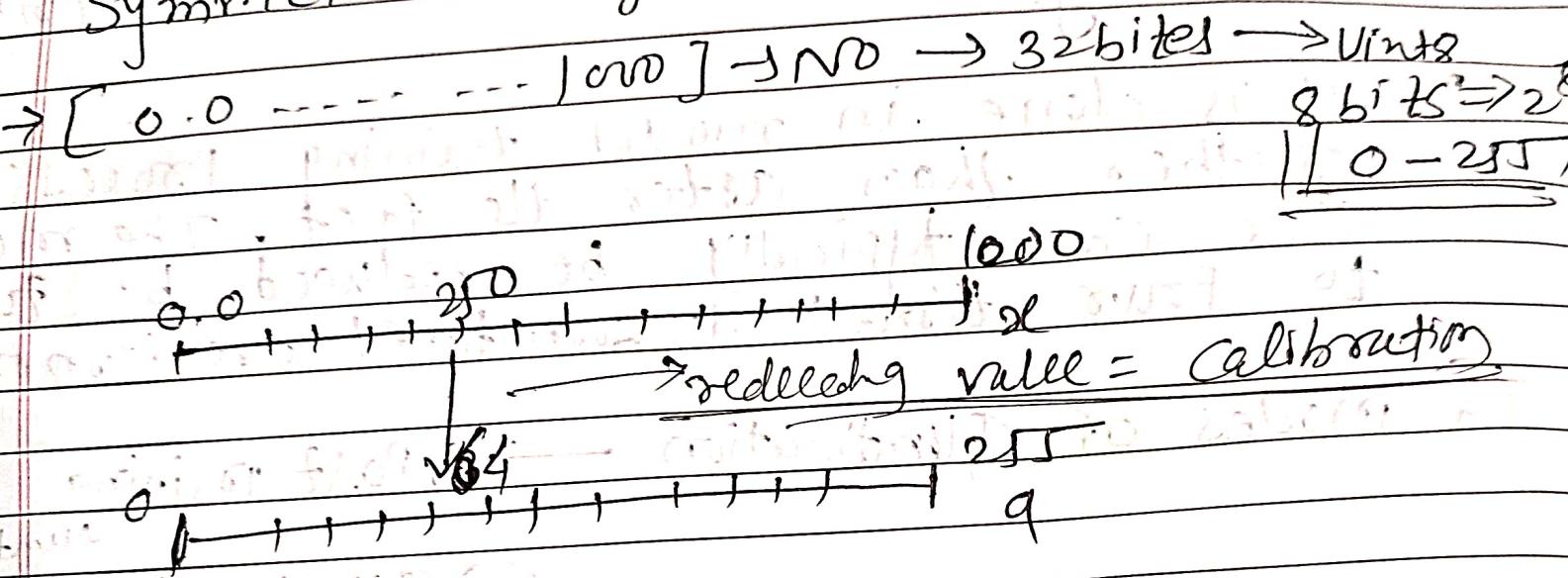
Calibration - model Quantization.

Types of Quantization

① Symmetric quantization ② Asymmetric quantization

1) Batch normalization
mean = 0.46 = 1

Symmetric Unsigned int 8 quantization.



min max scalar

0.0

1000

$$\text{Scale factor} = \frac{x_{\max} - x_{\min}}{q_{\max} - q_{\min}} = \frac{1000 - 0}{255 - 0} = 3.92$$

$$round = \frac{250}{3.92} = 64$$

Higher \rightarrow lower

Zero Point = 0
Scale = 3.92

② Asymmetric Quantization

$$[-20.0 \dots 1000.0]$$

$$[0.0 \dots 255]$$

$$-20.0 \rightarrow 0$$

$$1000 \rightarrow 255$$

$$\frac{1000 - 20}{255} = 4.0$$

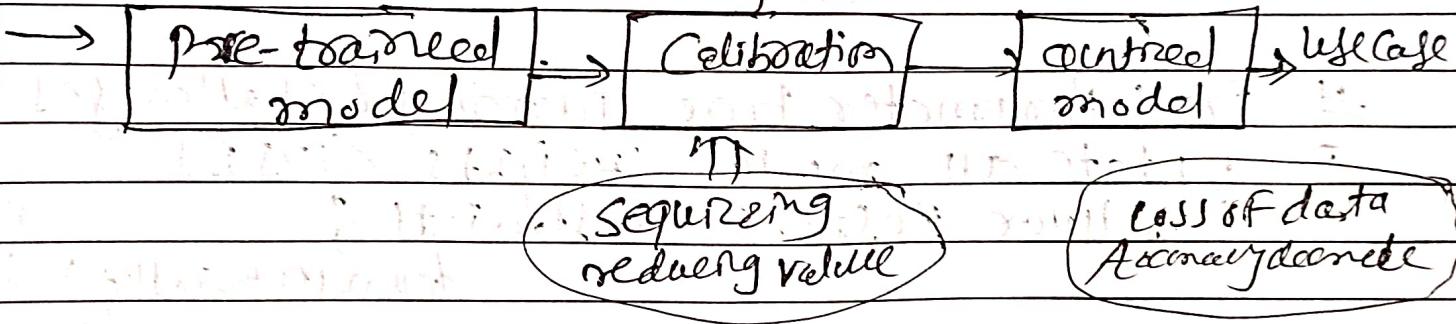
$$\text{mid} = \frac{(-20)}{4} = \underline{\underline{(-5.0)}} \quad -5 + 5 = 0$$

III) Modes of Quantization

zero point = 5
scale factor = 4.0

① Post Training Quantization (PTQ)

: weights + bias



② Quantization Aware Training (QAT)

(fine-tuning training data without sacrificing accuracy)

Trained model

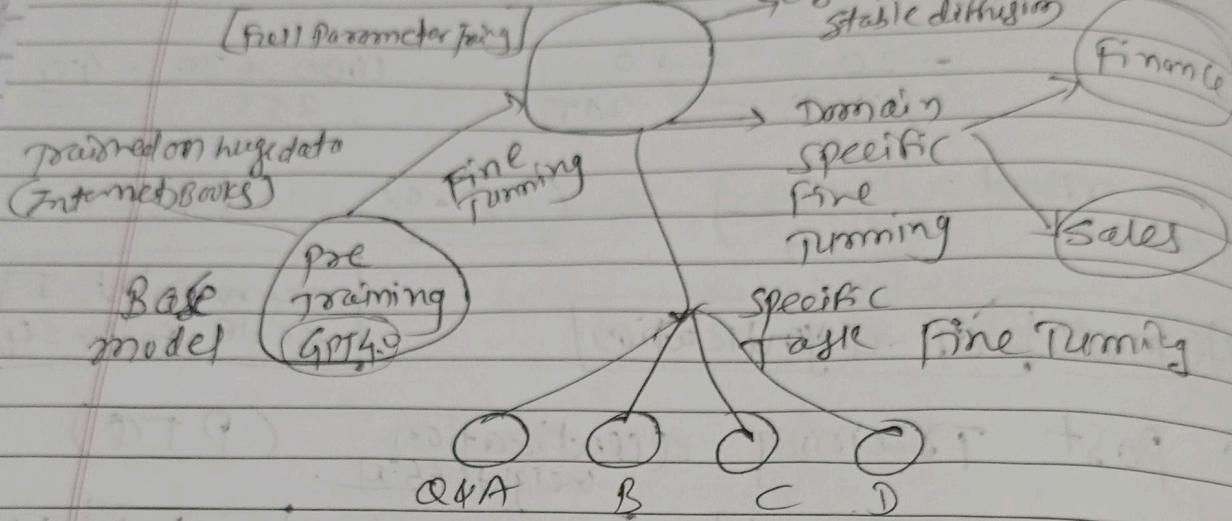
Quantization

Fine Tuning → Quantization model

Fine-tuning LLM model

LORA, QLORA

① LORA - Low-rank Adaptation



Full parameter Fine Tuning [Challenges]

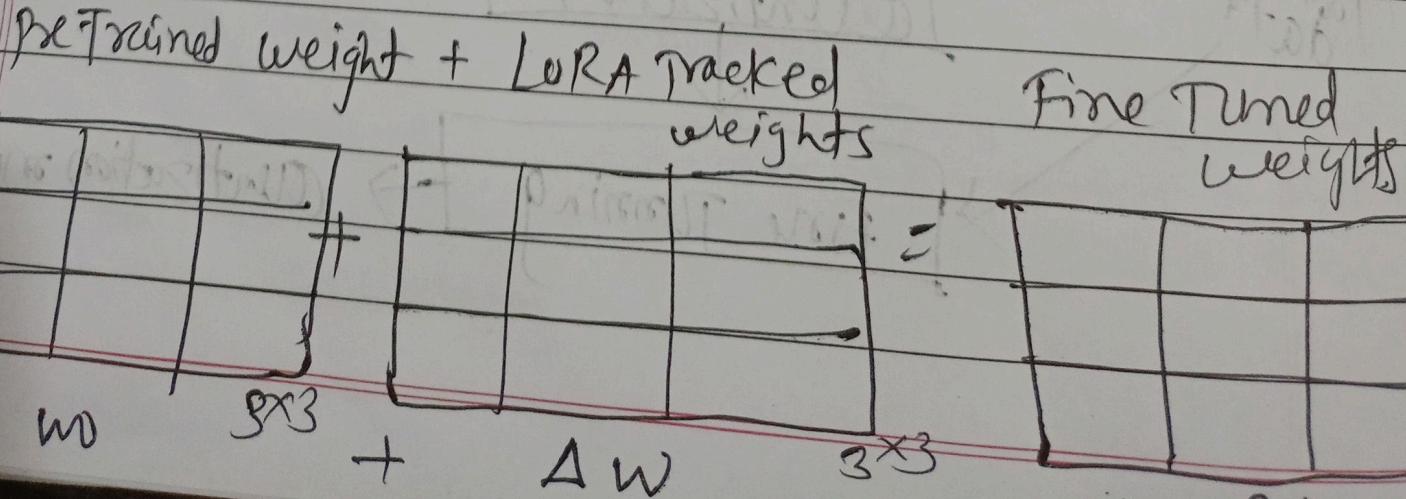
- ① update all model weights (175B)
- ② Hardware Resource Constraints

RAM
GPU
model memory

downstream task

To overcome full parameter fine-tuning we use LORA, QLORA. ($\frac{1}{2}$ Lora)

LORA — Instead of updating weights it tracks changes.



Fine-tuning LLM model

LORA, QLORA

Indepth mathematical Intuition

① LORA - Low-rank Adaptation

[full parameterizing]

trained on huge data
(internet books)

Base
model

pre
training
GPT-3

order

rank

Eg. ChatGPT
stable diffusion

Fine
Tuning

domain
specific
fine
tuning

Finance

specific
task fine tuning

Sales

Q&A B C D

Full Parameter Fine Tuning [Challenges]

① update all model weights (175B)

② Hardware Resource Constraints

RAM
downstream task
GPU
model format

To overcome full parameter fine tuning
we use LORA, QLORA. ($\frac{1}{2} \times 10^9$)

LORA — Instead of updating weights
it tracks changes.

PreTrained weight + LORA tracked
weights

Fine Tuned
weights

$$w_0 \quad 8 \times 3 \quad + \quad \Delta w \quad 3 \times 3 \quad = \quad w_0 + \underline{\Delta w}$$

model complete runway
High rank

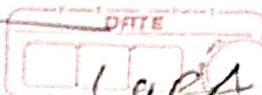
Matrix decomposition

$$\begin{array}{|c|c|} \hline 1 & B \\ \hline 2 & X \\ \hline 3 & \\ \hline \end{array}$$

$\begin{array}{|c|c|c|} \hline 4 & 5 & 6 \\ \hline \end{array}$

$3 \times 1 \quad 1 \times 3$

Linear Algebra



$$W_0 + Aw = w_0 + BA$$

[How to select rank?]

Rank = 1

Rank matrix = The rank of matrix is the number of linearly independent rows.
or
The number of linearly independent columns.

only zero matrix has rank zero.

example = $\begin{bmatrix} 1 & 2 & 3 \\ 0 & 4 & 2 \\ 2 & 6 & 5 \end{bmatrix}$ upper triangle matrix.

Operation: $R_2 \rightarrow R_2 - R_1$

$$\begin{bmatrix} 1 & 2 & 3 \\ 0 & 2 & -1 \\ 0 & 2 & 1 \end{bmatrix}$$

$R_3 \rightarrow R_3 - R_2$

$$\begin{bmatrix} 1 & 2 & 3 \\ 0 & 2 & -1 \\ 0 & 0 & 2 \end{bmatrix}$$

Upper triangle

Rank of matrix = 3

X decomposition
less parameter required, less matrix so, solve
challenges



QLORA = Quantization Low order Rank Adaptation

matrix decomposition

$$\begin{bmatrix} 1 \\ 2 \\ 3 \end{bmatrix}^T \times \begin{bmatrix} 4 & 1 & 5 & 6 \end{bmatrix}$$

Float 16bit

↓ Convert

Float 4bit

Reduce the memory - Convert Higher memory format to lower memory format

bits and bytes = process of quantization

Fine Tuning = parameter efficient Transfer Learning for NLP

Transformers = "Attention - all you need"
"state of the art"

NN Architecture

Transformer Architecture

① Embedding - text input divided into smaller units called tokens.

Tokens \rightarrow numerical vectors
Converting

Called embeddings

~~Bitnet b1.58~~

Ternary { -1, 0, 1 }

12-1-24

~~LLMops & LLM Pipeline~~

VEXT — platform \rightarrow LLMops

Create AI Project

Instead of building any model from scratch
use VEXT LLMops without coding.

Fine tune Gemma Models in Keras
using LORA.
import keras
import keras_NLP

Backend Jax

② Transformer Block —

- ① Attention mechanism - It allows tokens to communicate with other tokens, capturing contextual info & relationship between them
- ② MLP — multilayer perceptron layer — FFNN that operates on each token independently.

- ③ output probabilities — Final linear & softmax layers transform the processed embeddings into probabilities enabling the model to make predictions about the next token in a seq.