

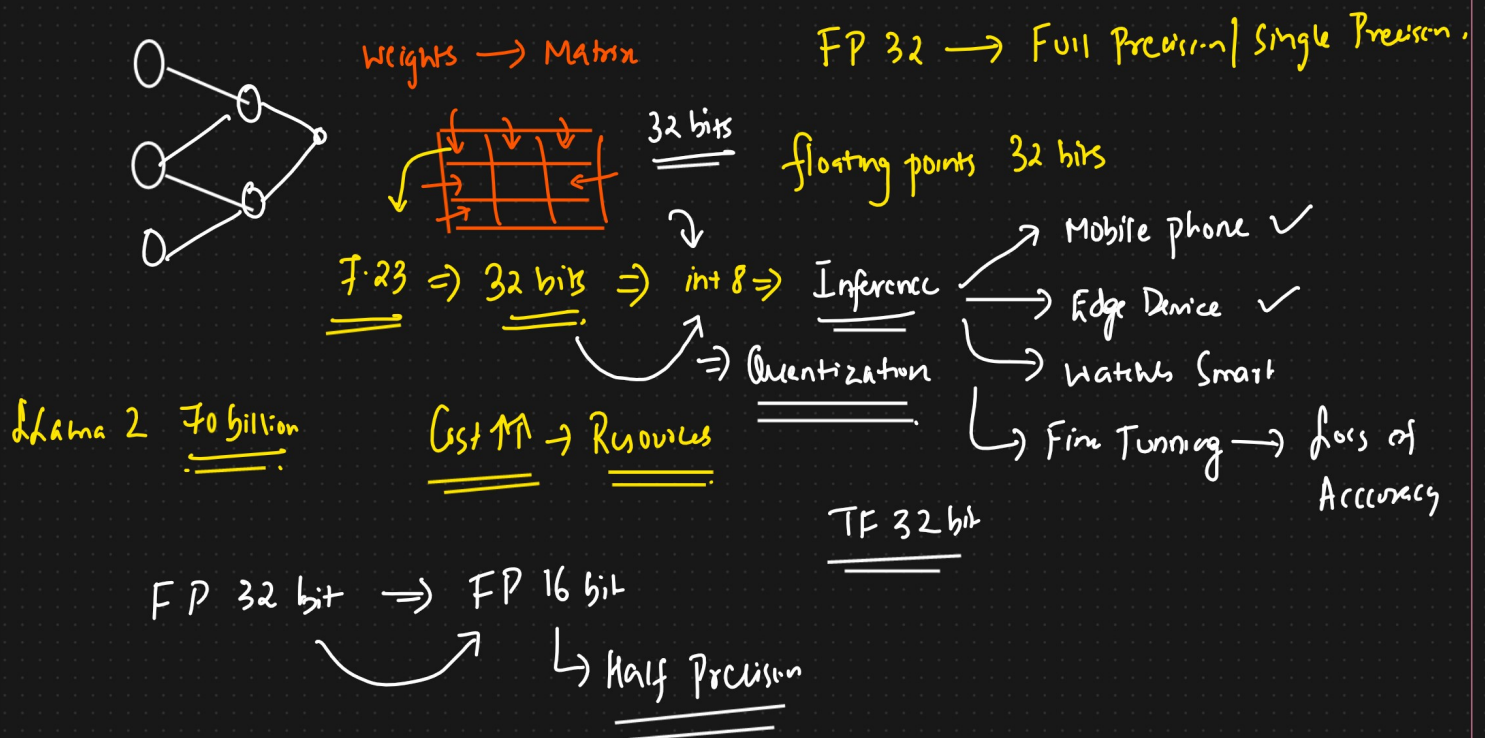
# Fine Tuning LLM Models

QLoRA, LoRA

## ① Quantization

- ① Full Precision / Half Precision → DATA → Weights and parameters ✓
- 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.



## ④ How to perform Quantization

① Symmetric Quantization

↓

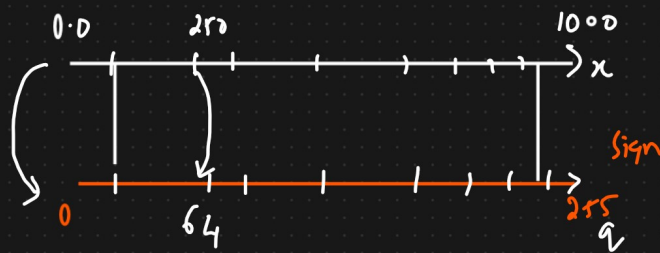
Batch Normalization

② Asymmetric Quantization

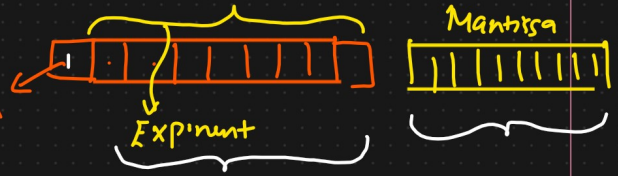
## ① Symmetric uint8 quantization

$$[0.0 \dots 1000] \rightarrow \text{Numbers} \rightarrow 32 \text{ bits} \rightarrow \text{uint8} \Rightarrow 8 \text{ bits} \Rightarrow 2^8$$

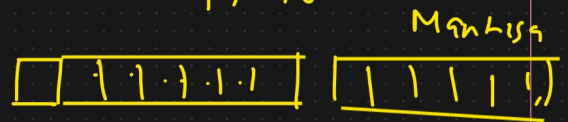
$\hookrightarrow \underline{0-255}$ 
 $\boxed{0-255}$ 
 $\boxed{7.32}$



Single Precision Floating Point 32



FP 16



Min max scalar

$$\begin{aligned} 0.0 &\rightarrow 0 \\ 1000 &\rightarrow 255 \end{aligned}$$

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

$$\text{round} \left( \frac{250}{3.92} \right) = \underline{\underline{64}}$$

$$\text{Zero point} = 0$$

## ② Asymmetric uint8 quantization

$$[-20.0 \dots 1000.0]$$

$$[0 \dots 255]$$

$$-5.0$$

$$\text{Zero point} = \underline{\underline{5}}$$

$$\text{Scale} = \underline{\underline{4.0}}$$

$$\text{round} \left( \frac{-20}{4} \right) = (-5.0) + \boxed{5} = \underline{\underline{0}}$$

$$\frac{1000 + 20}{255} = \underline{\underline{4.0}} \Rightarrow \text{Scale factor}$$

$[-128]$

$[127]$

## ⑧ Post Training Quantization (PTQ)



## ⑨ Quantization Aware Training (QAT)

