

## Finetuning

GPT → 3.5 turbo ← OPENAI API

Kkama, phi ← microsoft

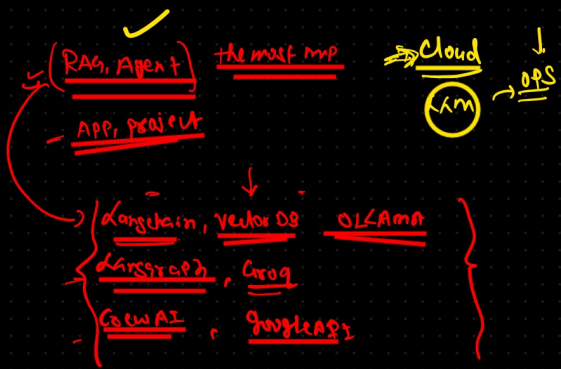
→ modular coding

Some theory points on LoRA, PEFT, Quantization

⇒ Reward learning (RLHF, Preference optimization) ← Sun

Reday made ⇒ Kkama factory, unsloth

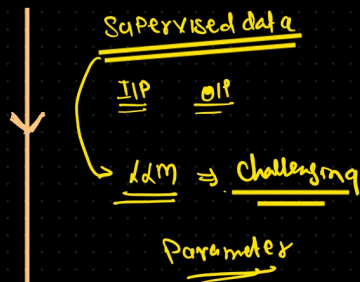
(MARI APP)



① Pretraining

② SFT

③ RLHF



PEFT ⇒ LoRA, Quantization + LoRA ⇒ QLoRA

① PEFT, LoRA ⇒ without modify all parameter we train the model

- 1) Reduce computational cost
- 2) less memory
- 3) Speed up the process
- 4) maintain the model performance

Kkama → transformer ⇒ Self Attention, Neural network

✓ weights

W → huge numbers  
entire weight  
frozen (freeze)

⊂ A  
 ⊂ B

$$W' = W + \Delta W$$

$$\Delta W = A * B$$

$$A (r \times d)$$

$$B (d \times r)$$

⊂ A

$d \Rightarrow$  dimension  
 $r \Rightarrow$  low rank factor

Interview question

LoRA

32 GB RAM

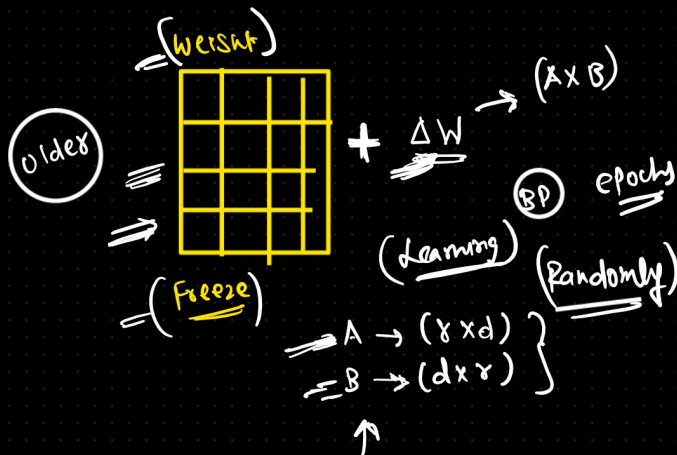
QLoRA

12-24 GB RAM

Transformer  $\Rightarrow$  Self Attention + NN

huge number weight

LoRA subset



$$\text{Rank} \Rightarrow \text{SVD}$$

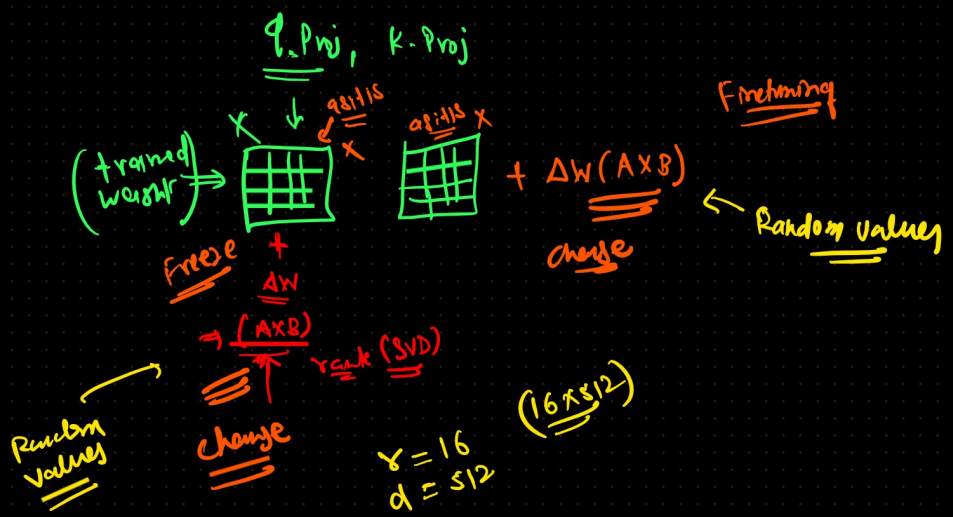
Self attention  $\Rightarrow$  weight matrix

hyperparameter

$\begin{cases} W_q \Rightarrow \text{query} \Rightarrow \text{q-Proj} \\ W_k \Rightarrow \text{key} \Rightarrow \text{k-Proj} \\ W_v \Rightarrow \text{value} \Rightarrow \text{v-Proj} \end{cases}$

Wqkv

out-Proj



⇒ 1 BIT ⇒ 0 or 1

System ⇒ 0.8

1 BYTE ⇒ 8 BIT

1 KB ⇒ 1024 BYTE

1 MB ⇒ 1024 KB

1 GB ⇒ 1024 MB

Datatype ⇒

float 32 = 32 bit = 4 byte  
float 16 = 16 bit = 2 byte  
float 8 = 8 bit = 1 byte  
int 8 = 8 bit = 1 byte

Datatype  
↓  
weight ⇒ 2.87 ⇒ float32 ⇒ 4 byte  
175B × 4 byte  
GPT3.5 ⇒ 175B weight =

Quantization ⇒ reduce precision

FP32 → int8

4 byte → 1 byte

LORA ⇒ 32 RAM, 64 RAM

Q + LORA ⇒ 16 RAM, 24GB RAM

⇒ Quantization ⇒ Reduce the precision

Loss of information

① 4 byte → 2 byte  
4 byte → 1 byte  
4 byte → 0.5 byte

loss ↑ ④    loss ↓ ①

lower precision (loss)



(PTQ)

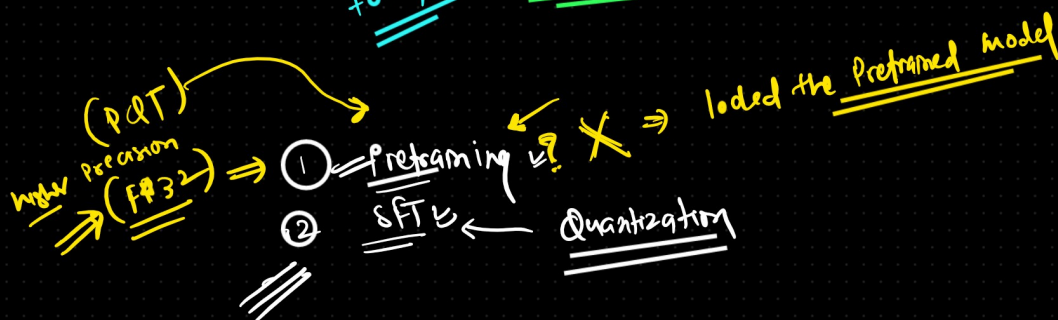
(QAT)

# 1 Post-training quantization

- 1 First trained in full precision (FP32)
- 2 After training the w & b are converted to lower precision

{ FP16  
Int8 }

today's  $\rightarrow$  LoRA + QLoRA  $\Rightarrow$  (PTQ) QAT

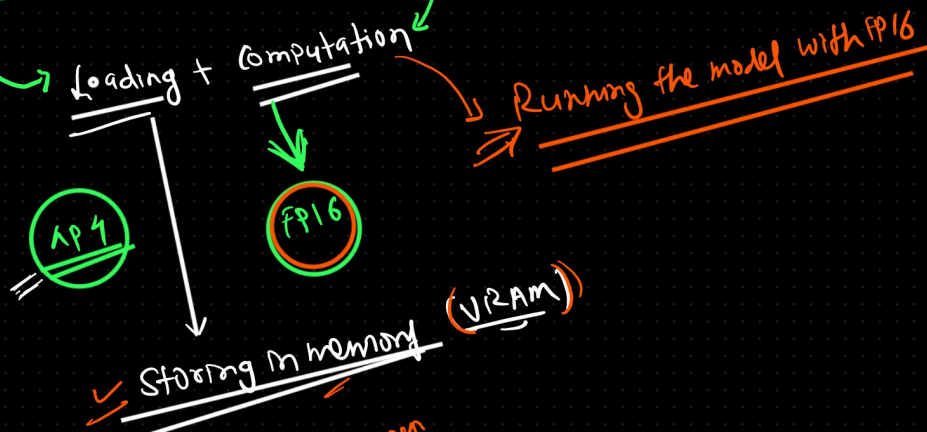


## PTQ $\hookrightarrow$ Post Quantization $\hookrightarrow$ technique

- 1 unsupervised  $\leftarrow$  FP32  $\rightarrow$  full precision
- 2 SFT  $\leftarrow$  Quantization

```
bnb_config=BitsAndBytesConfig(
  load_in_4bit=True,
  bnb_4bit_use_double_quant=True,
  bnb_4bit_quant_type="nf4",
  bnb_4bit_compute_dtype=torch.float16)
```

normal float 4 bit (NP4 bit)



Accuracy  $\leftarrow$  (loss)  $\leftarrow$  (loss)  $\leftarrow$  (loss)  
 (NP4) (FP16) FP32  
 $\rightarrow$  Computation

DPO, PPO