① → ∂. FINETUNING

RAG →

chat } Ans:

---

FINETUNING : →

→ COMPLEX.
     REASONING.

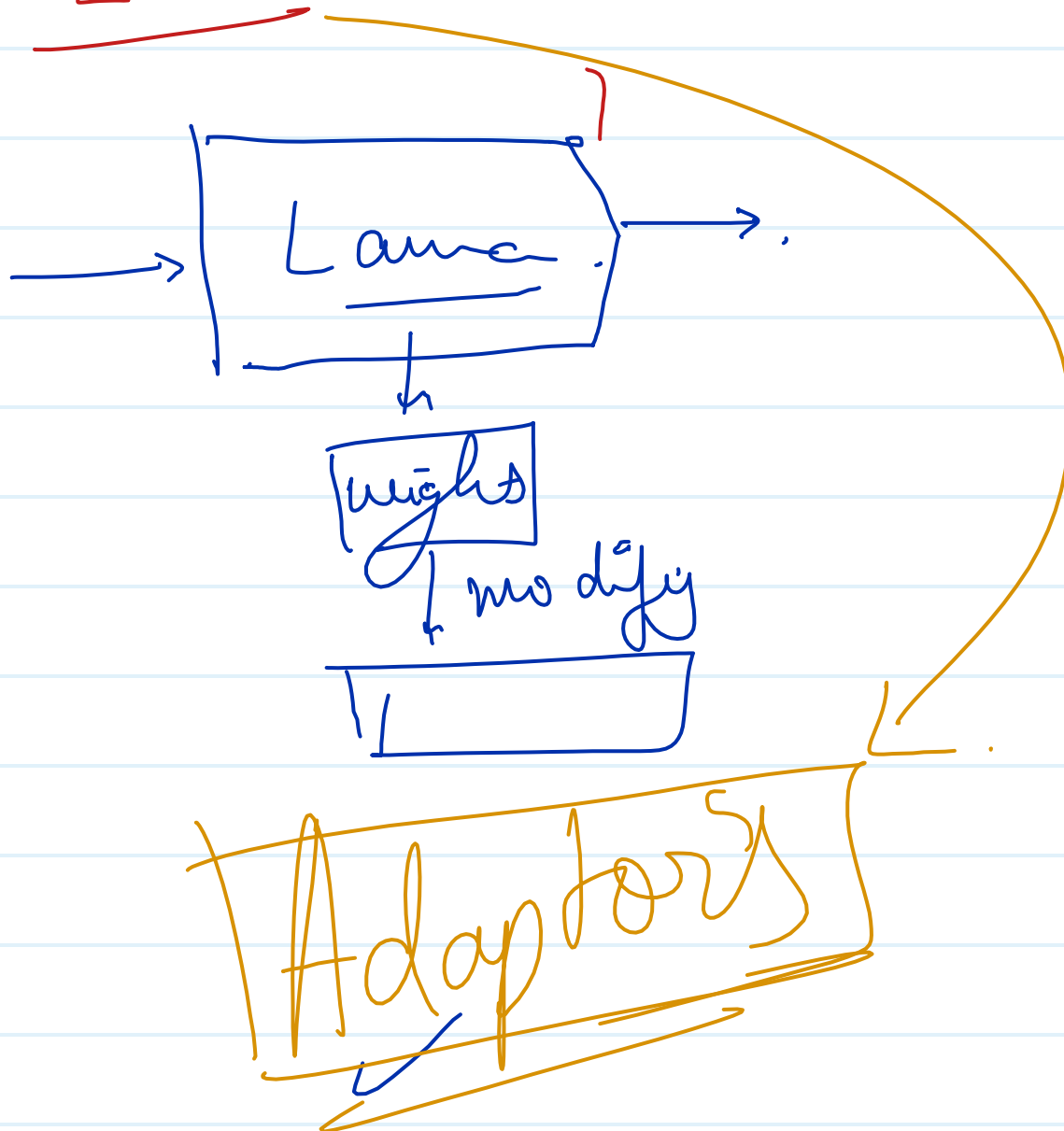→ Specific. output (JSON) <

→ Interview ←

PEFT $\leftarrow$ { Parameter efficient }.
F.T.

✩



90% remains same

FT.

PEFT $\longrightarrow$ LORA
$\longrightarrow$ Low Rank Adaptation

$\longrightarrow$ QLORA.

$\longrightarrow$ Quantized ———

# LORA:

Layer

weights

↱ modify

Adaptors

$$W_{new} = W_{old} + \Delta W$$
$$(n, m)$$

Adaptors $\longrightarrow$ A $\times$ B.

Matrix decompositions.

Noise    zero.

$\square \times \boxed{\phantom{xxxx}}$

$W_{old} \Longrightarrow (4096 \times 4096).$

$A \times B \Rightarrow A \longrightarrow (4096 \times r)$

$B \longrightarrow (r \times 4096).$

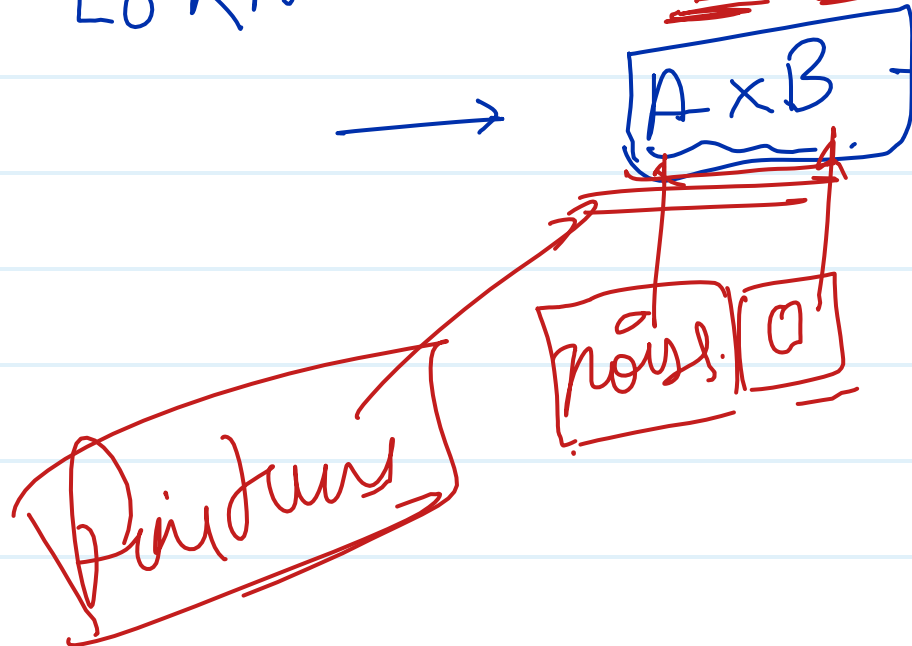Rank $\Longrightarrow r = 1.$

$A \longrightarrow (4096 \times 1)$

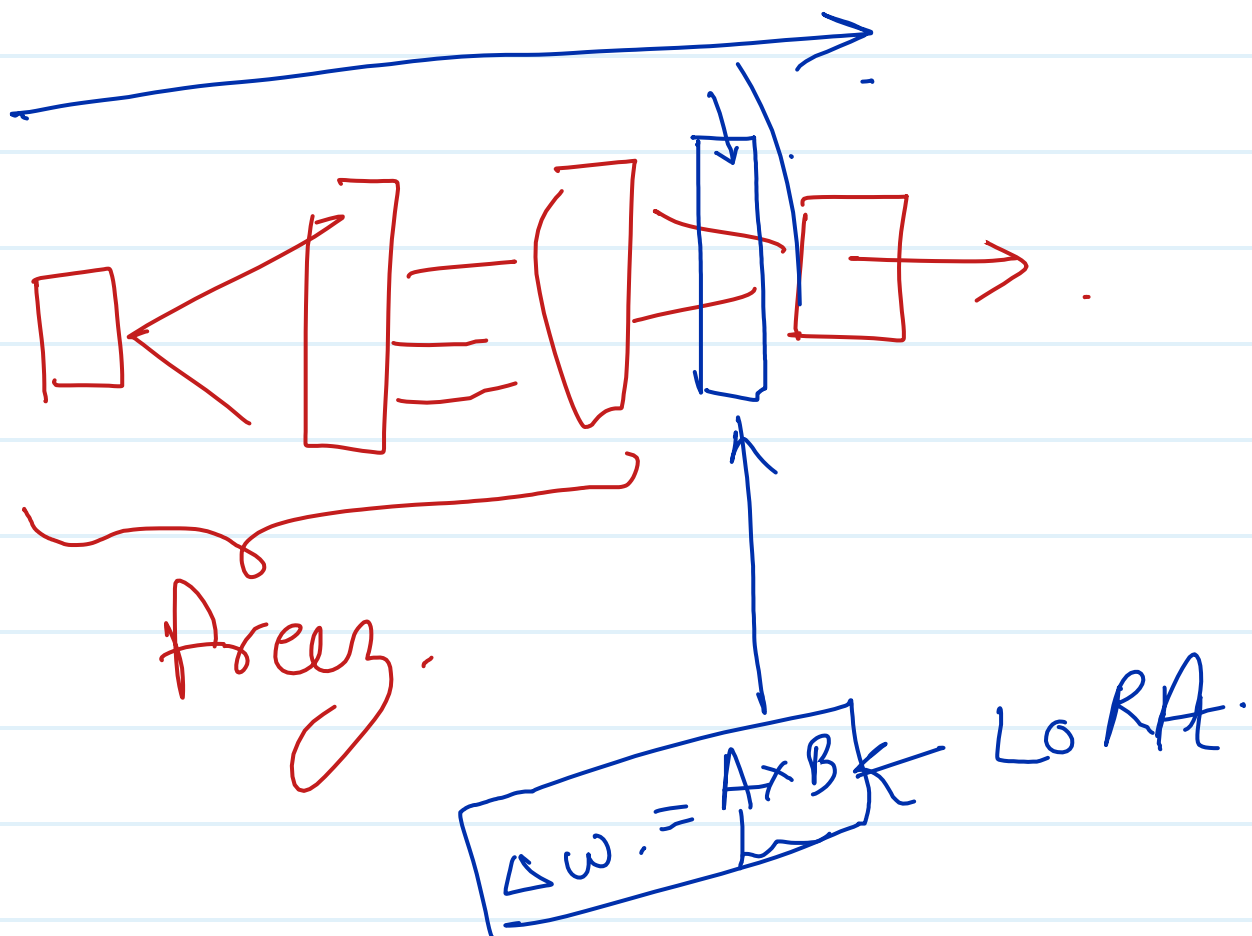$B \longrightarrow (1 \times 4096).$

$$A \times B = 16 \text{ million}$$
$$(n \times m) \quad (m \times p)$$

$A \times B \Rightarrow.$

LoRA $\Rightarrow$ $\times$ Pretrained (Freeze)

$\rightarrow$ $A \times B$ $\uparrow$ $\dagger.$

noise. $\boxed{0}$

Weights

Freez.

$$\Delta w. = A \times B$$ ← LoRA.

$A \rightarrow 4096 \times \boxed{1} = 4096$

$B \rightarrow 1 \times 4096 = + 4096.$

$\sim 9000$

Pretrained $\rightarrow 4096 \times 4096$

$\Rightarrow = 16 \text{ million}.$

ally is $B = $ zero m-1

zero $\rightarrow 0 \rightarrow 0.001$,

$A \times B = 0$.

$W_{new} = W_{old} + \Delta W$

$= 0$ ← no traing

---

$\#\theta + LoRA$

$\longrightarrow$ egc. $\longrightarrow$ 16 bit

Quantize

4 bit

# Quantize → Pretrained (eg. (6 → 4)

→ bits & byts lib.



16 → ④      ⑯

→ Storage.
→ Computation...

GPU ← chunks ← { 10 0 ← pre
                { 10  ④ ← fine
                  (16).

dequantize
4 → 16 ← instance

Formal $\Rightarrow$ $\omega_n = \underbrace{\omega_0}_{} + \overbrace{\left(\dfrac{\alpha}{\gamma}\right)}(\Delta\omega)$

$$\dfrac{}{f} \qquad \underset{(A \times B)}{\underset{\sim}{\uparrow}}.$$

Thumb rule : $\alpha = 2\gamma$.

$$\gamma \underset{\sim}{\sim} 16.$$