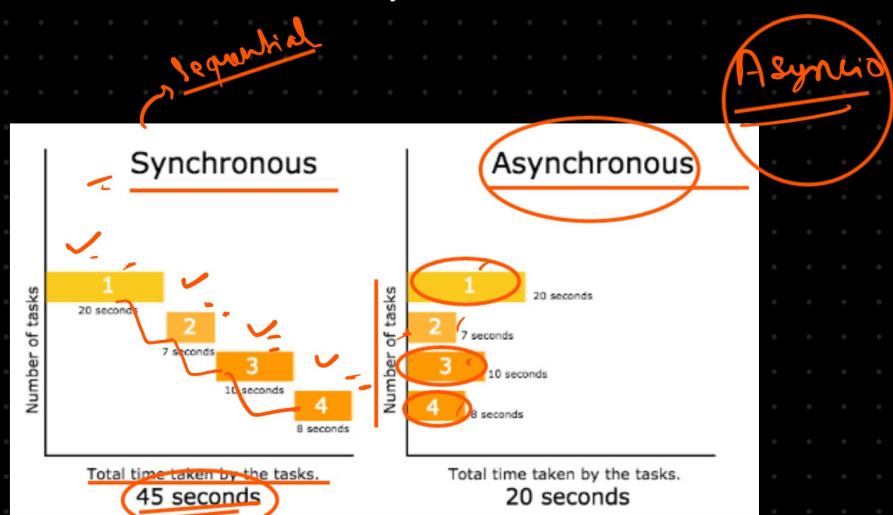
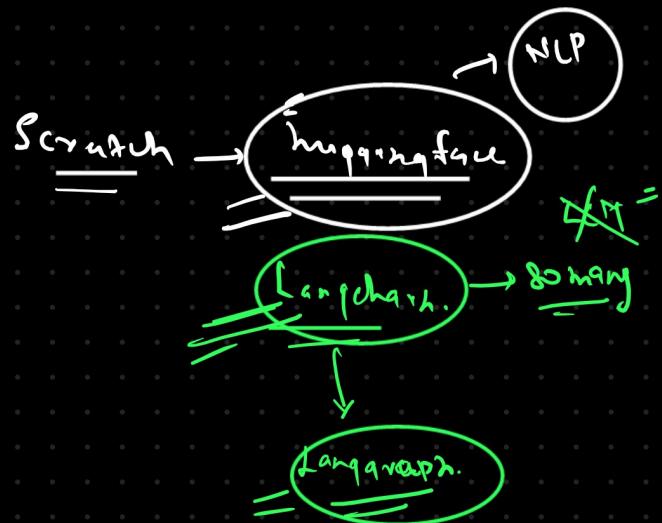
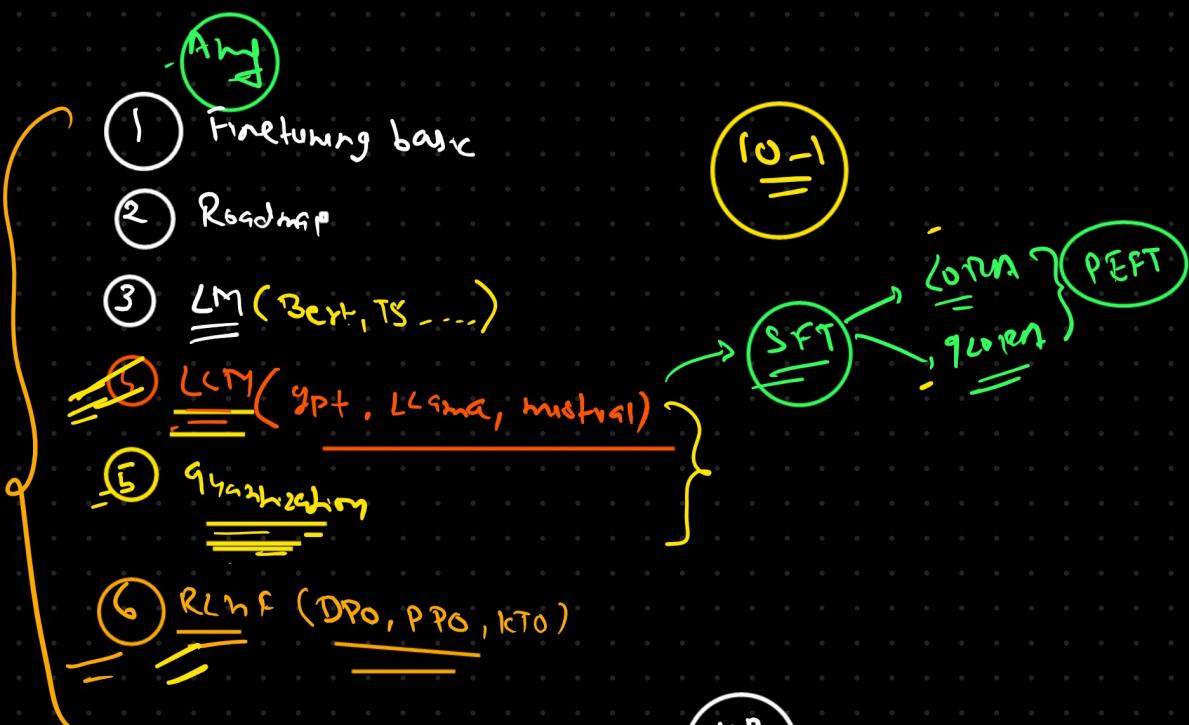


$$\underline{\text{Prompt}} + \underline{\text{LLM}} = \underline{\underline{\text{Output}}}$$

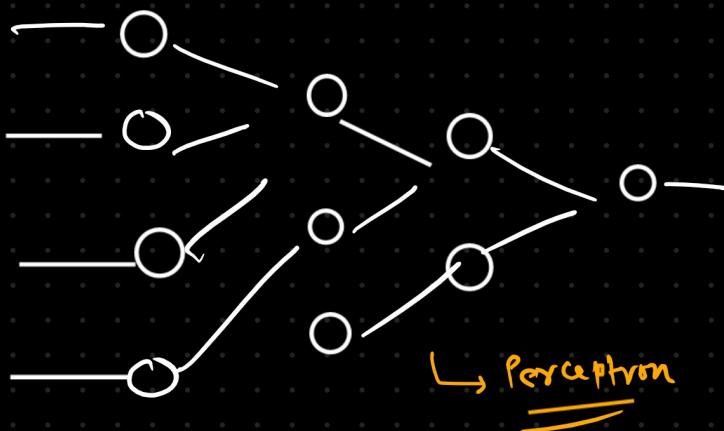
Prompt + LLM + StructParser = Output





LLM → NN

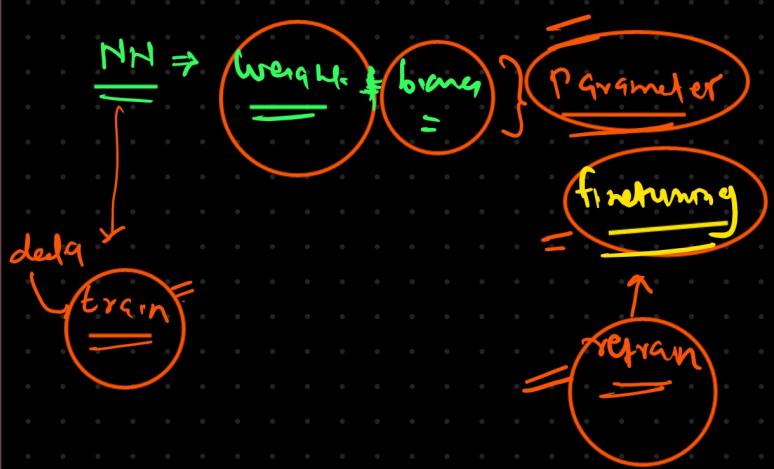
ANN



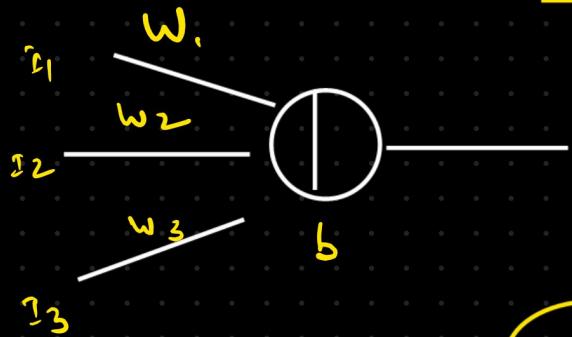
CNN → Convolution

RNN → Recurrent



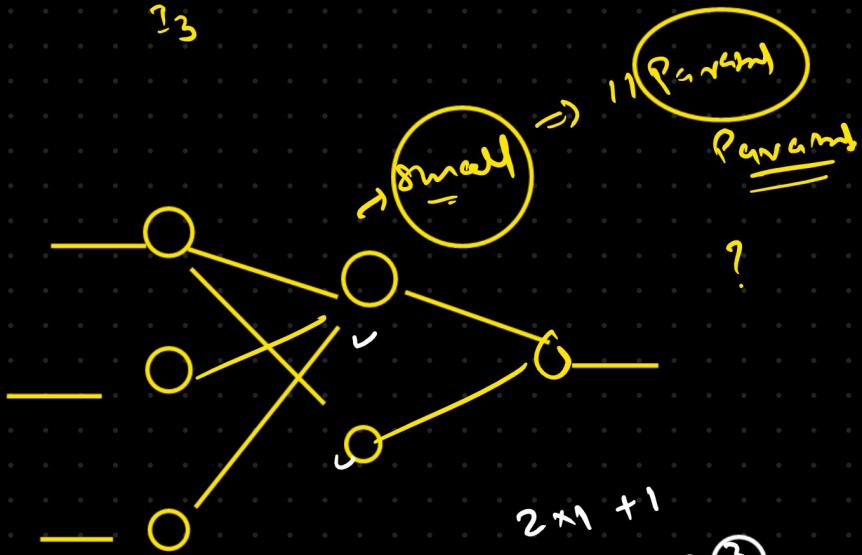


Neural net -



w<sub>1</sub>, w<sub>2</sub>, w<sub>3</sub>, b

parameter = 4



$$3 \times 2 = 6 + 2$$

$$2 \times 1 + 1$$

$$2 + 1 = 3$$

$$\textcircled{8} \quad 8 \times 3 = \textcircled{11} \xrightarrow{\text{Parameter}} \text{variable param}$$

variable param

NN  $\Rightarrow$  X

Parameter  $\rightarrow$  variable param

(w & b)

GPT3 → 175B =

Billion → Size of the model

Llama → 2B, 7B, 40B, 100B

+ Data

Mistral → 2, 1.5, 7B, 3B, 207B

Gemm → Gemm 1.5 pwo → 175B.

Umade → Billion

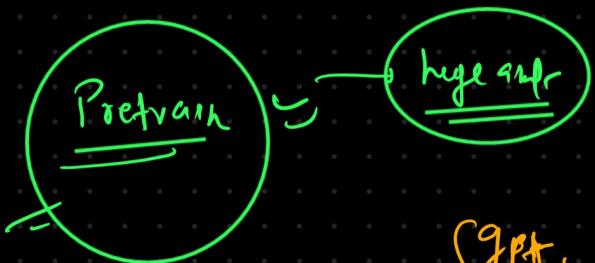
Training ↗ Start

You, ↑  
small org.



LM → training → huge Param

(Billion)



= (GPT, Llama, Mistral)

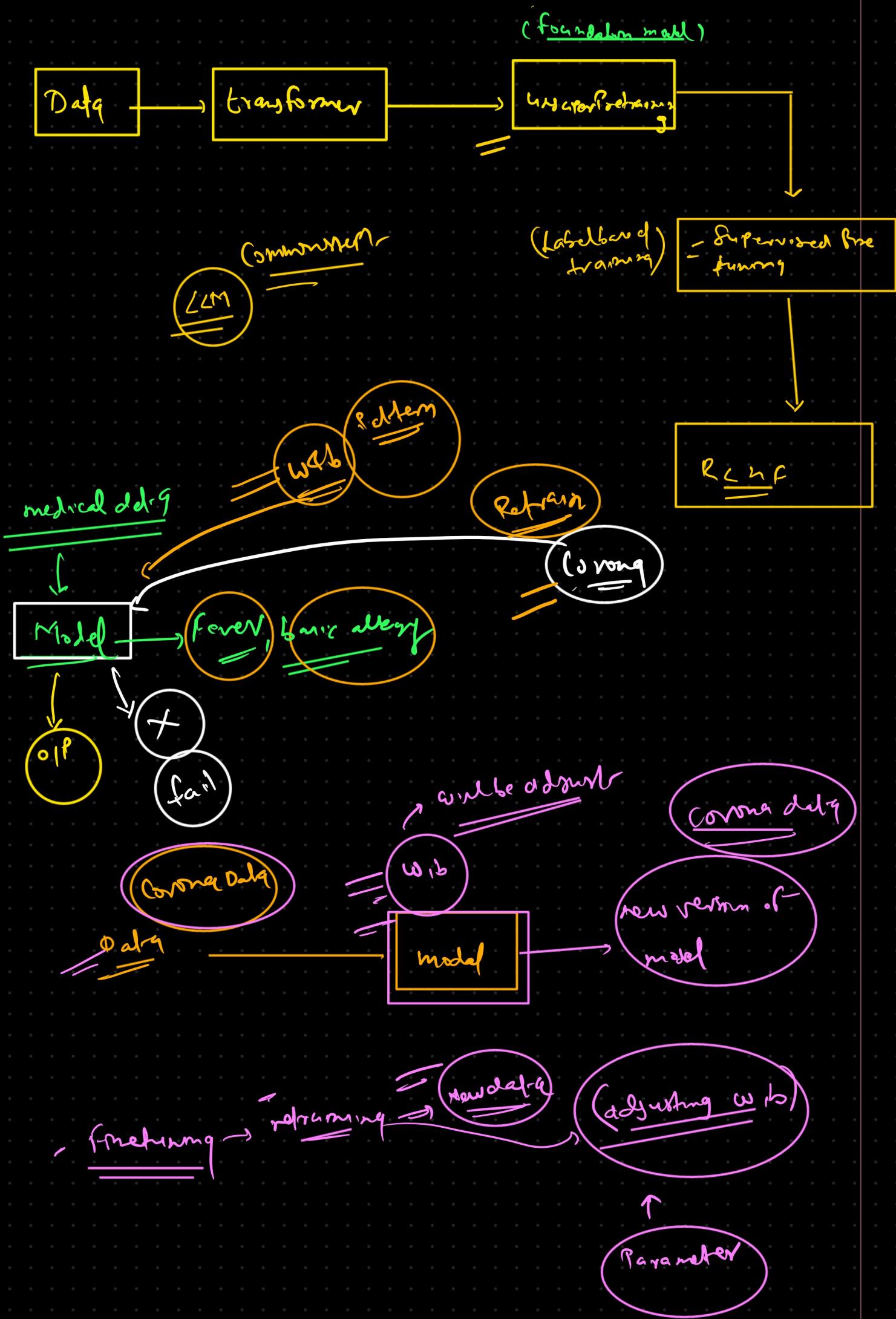
transformer

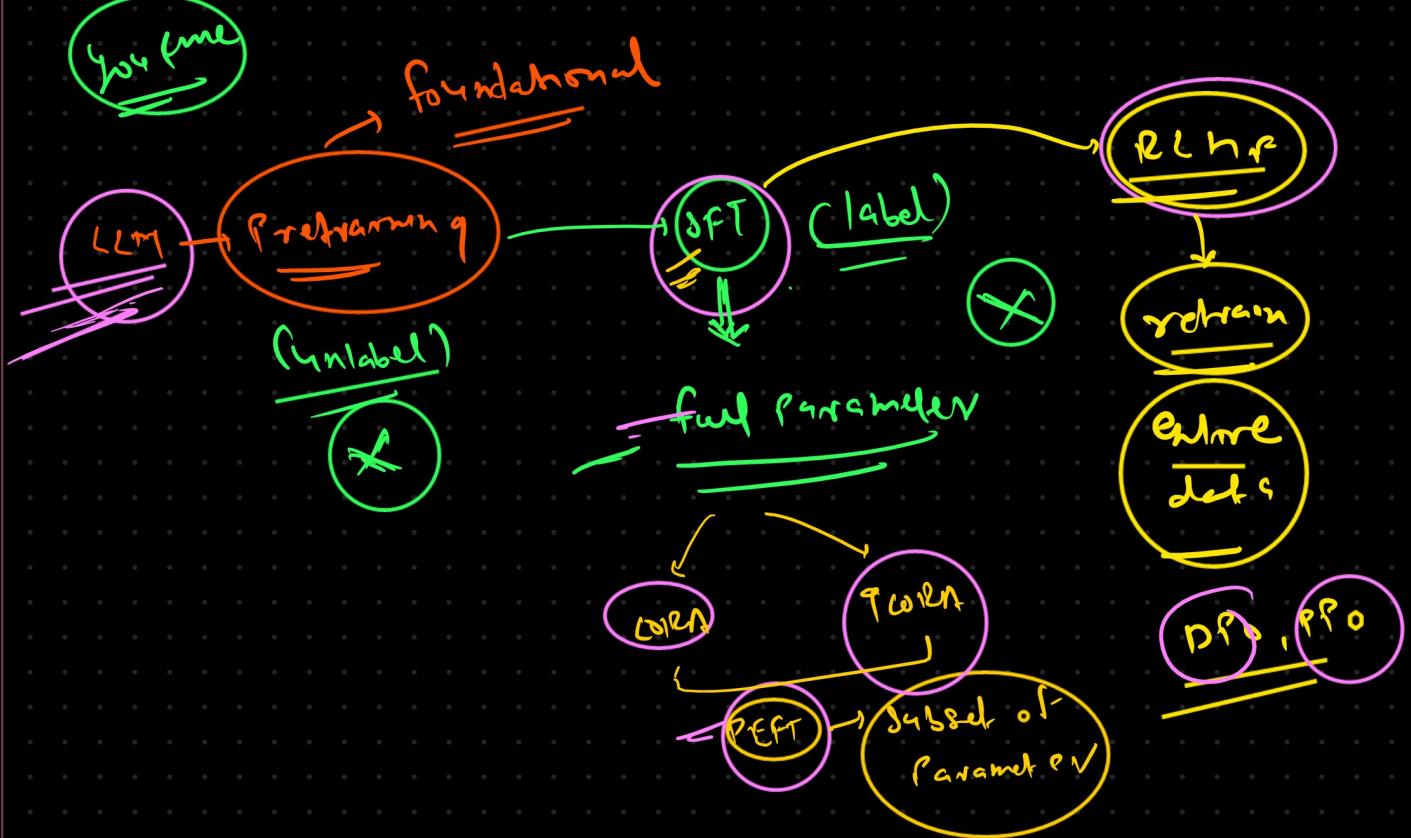
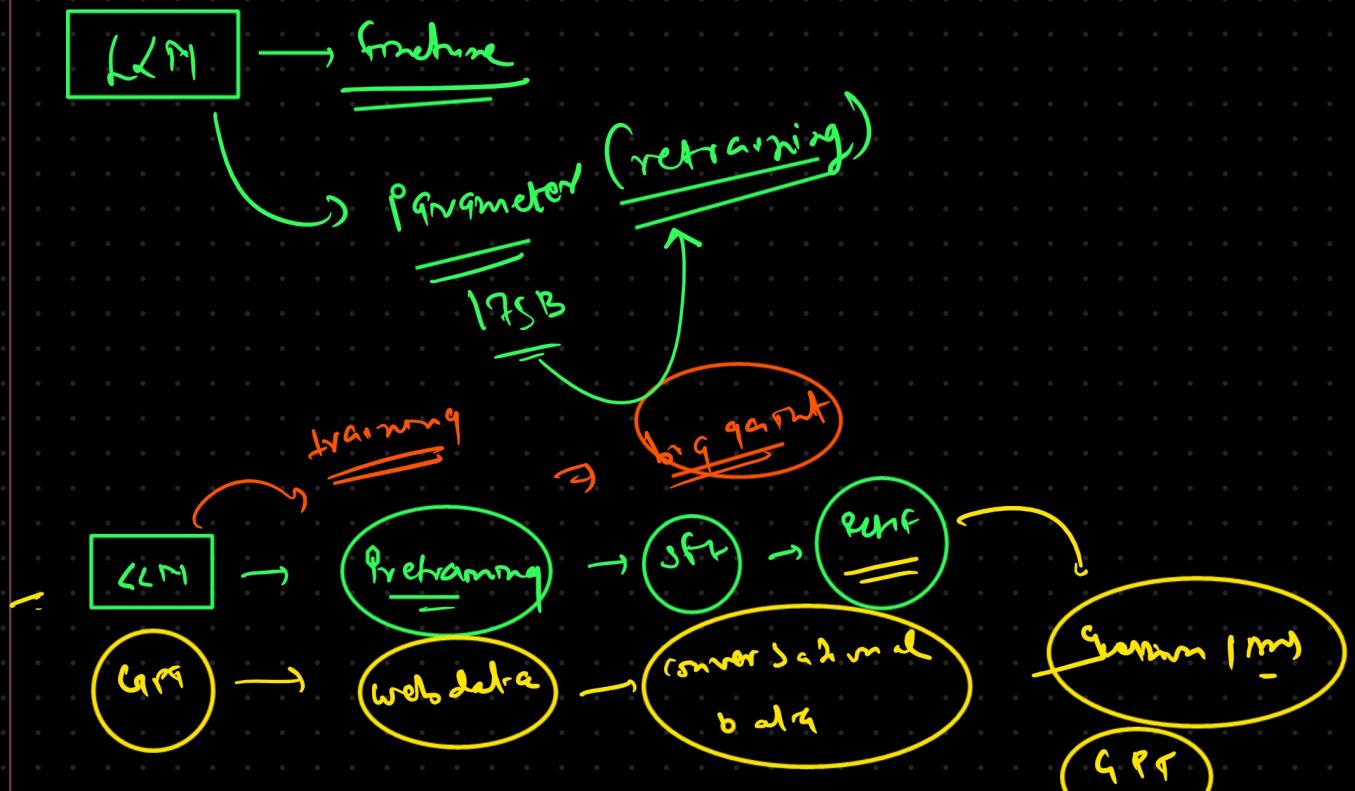
- Attention → weight matrix

wk

rows

Base unit

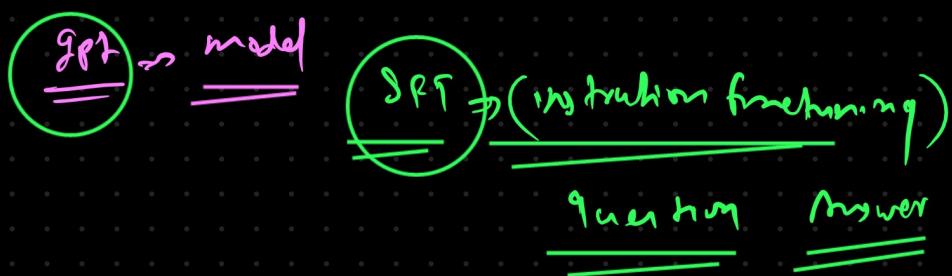




finetuning → retraining → param → Data

↑  
domain specific

Lama, mpt6

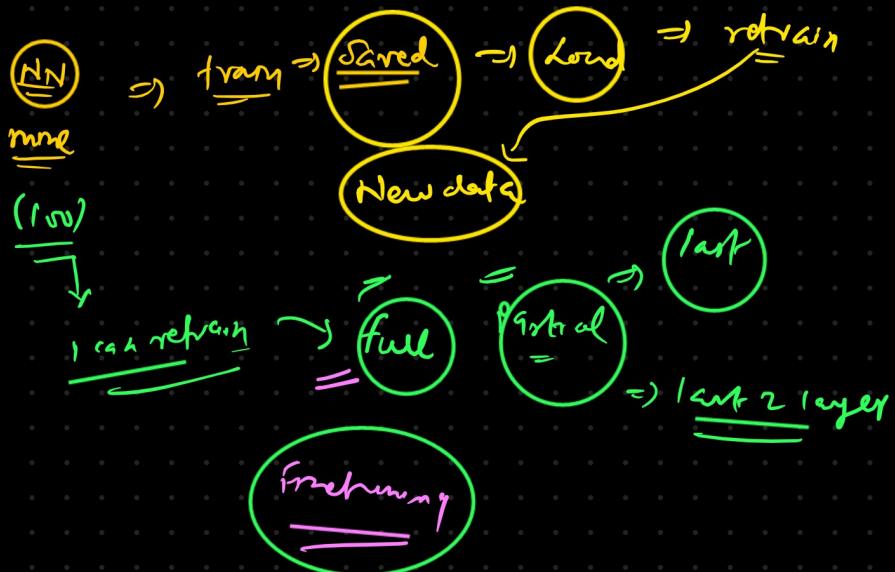


Finetuning

parameters ( $w \& b$ )

Retraining

Domain specific data



LM vs LLM

Raw  
CNN  
GRU

transforming

huge data

Millions  
Ts  
billions

LLM (Bert)  
Brown, GPT

GPT-3, LLaMA, Qwen

Scaling

foundational

fine-tuning

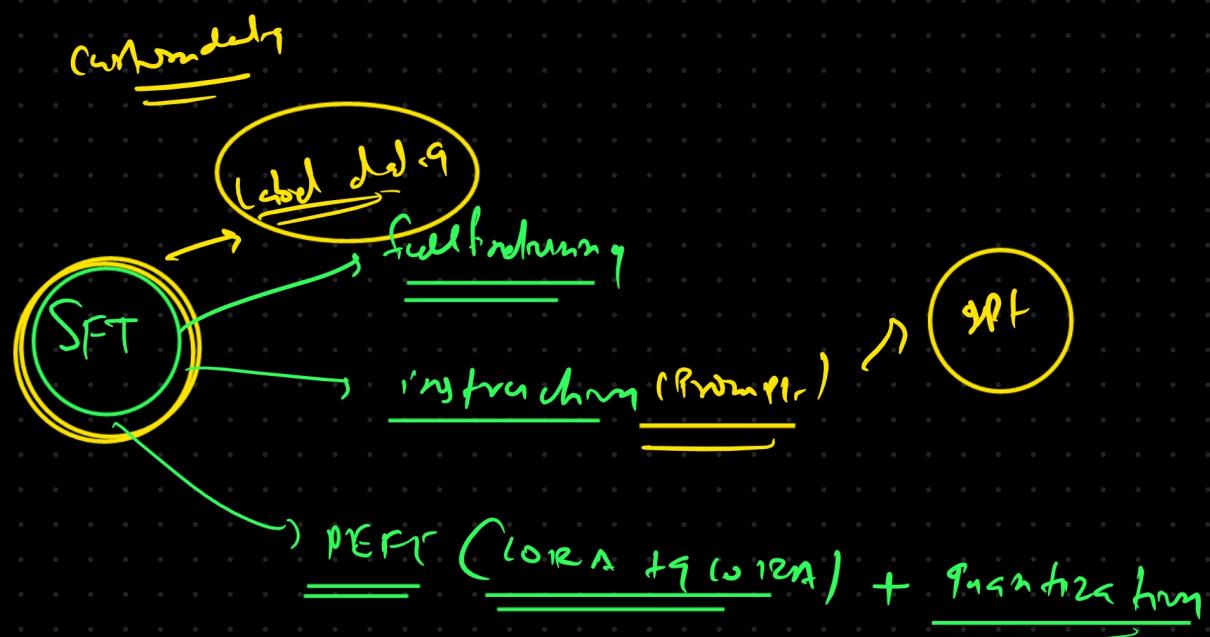
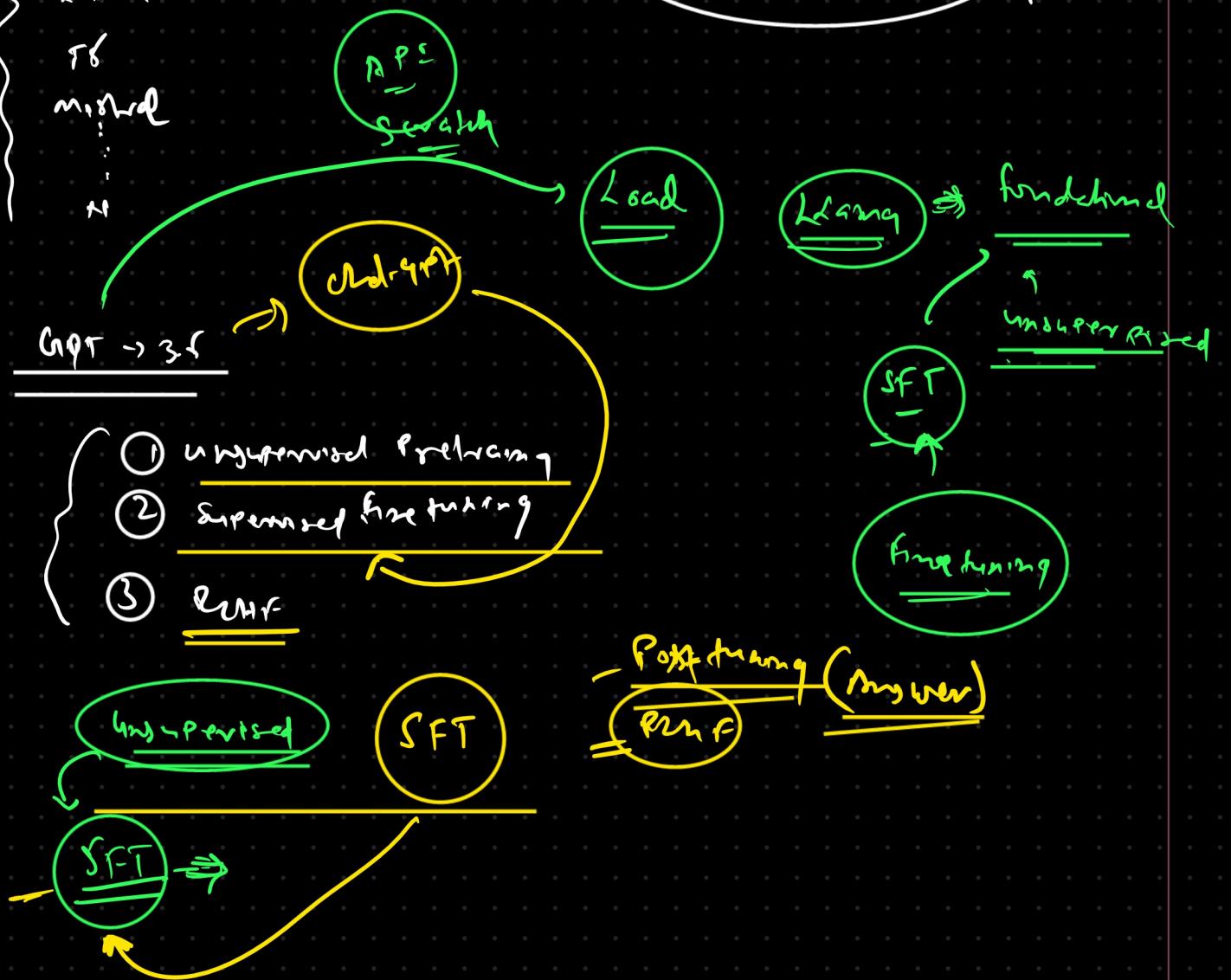
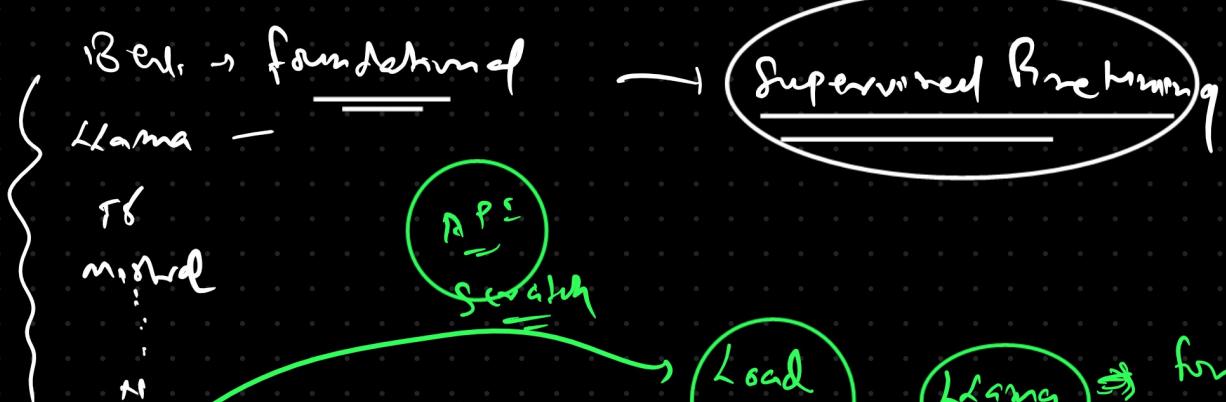
unsupervised training

LLM  $\Rightarrow$  X

Data

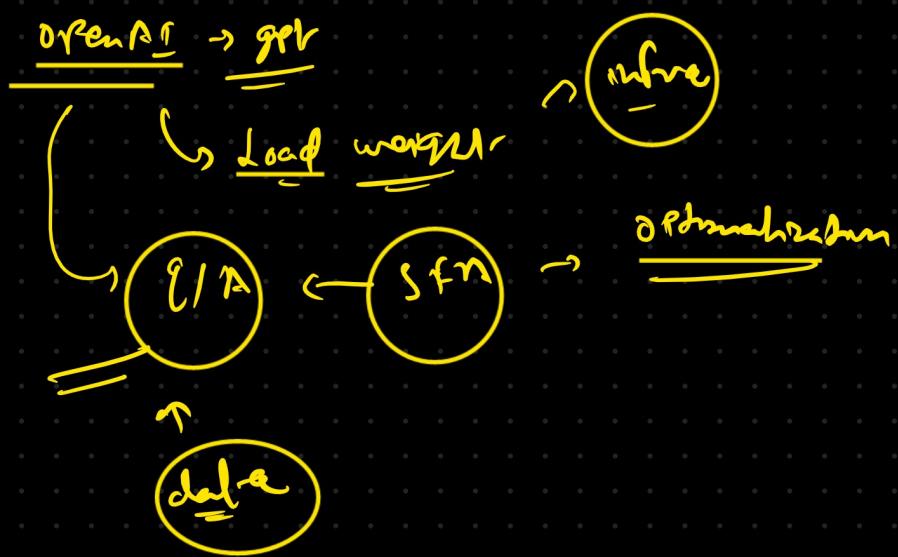
Not label

training



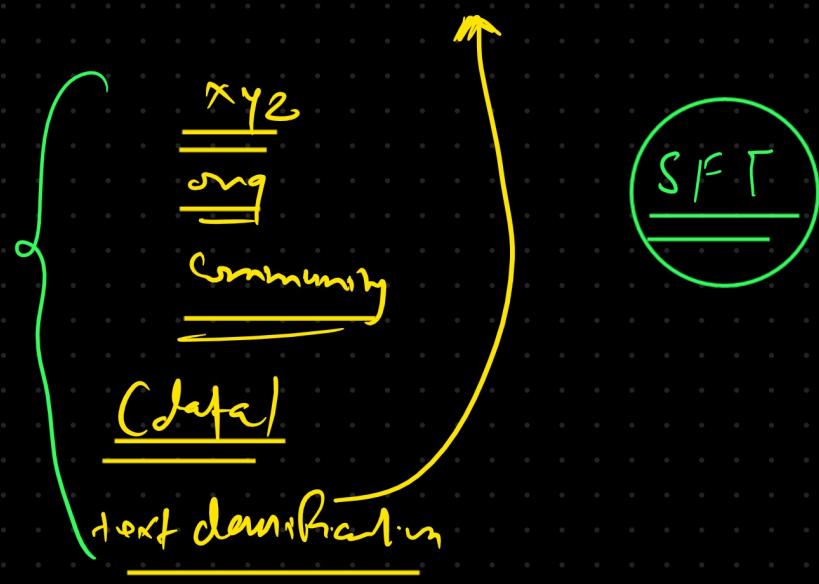
SFT  $\Rightarrow$  Supervised Learning

Label data



Foundational  $\Rightarrow$  Unsupervised model

Pretrain mode

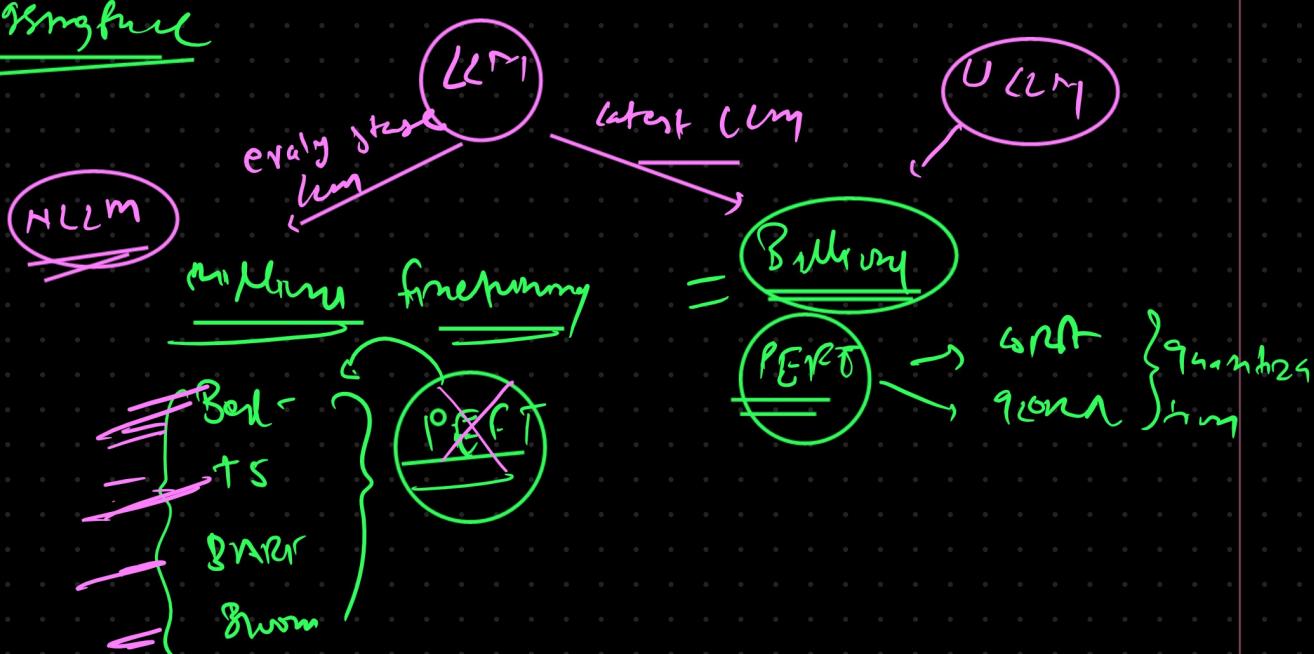


<u>Question</u>	<u>Answer</u>
-----------------	---------------

PLNF

Warren, Michael, Zephyrus or any model

logging rule



\*

NN → fram → save → load → file naming

Model  $\rightarrow$  Pretarm (unsupervised Pretraining)

$\downarrow$   
load model (W&B) Parm  
 $\downarrow$

Data  $\xrightarrow{\text{Input}}$  (Model)  $\rightarrow$  output  
Inferencing

Prediction  $\xrightarrow{\text{Pretarm}}$  load  $\rightarrow$  Path  $\hookrightarrow$  Prediction  
outPut

Inferencing  $\rightarrow$  difficult  $\rightarrow$  huge  
Llama 7.3.1  $\Rightarrow$  Infra  
quantization

quantization  $\xrightarrow{(?)} \rightarrow$  fair inferencing

Inferencing  
 $\downarrow$   
Prediction

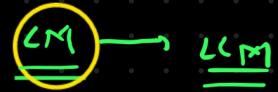
Model  $\rightarrow$  head  
 $\xrightarrow{\text{W&B}}$  quantizer

inference  $\rightarrow$  APF  $\leftarrow$  whq  
exm

BFR

Done (hdd, sdd)

1 inference



BERT

empty

B

$$\text{Learn} \xrightarrow{\text{3.1}} \text{out} + B \xrightarrow{\text{Inputize}}$$

2 finetuning

Next



BERT, (GPT-3.5, LLaMA, mT5)



Mistakes

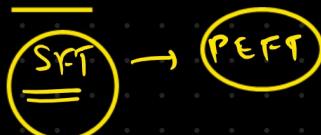
~~PEFT~~

Pretrain

SFT

(LLM)

LoRA, QLoRA

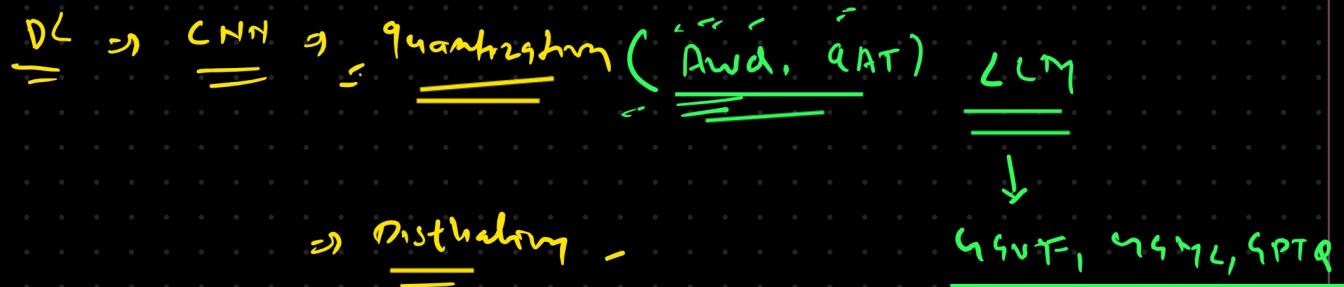
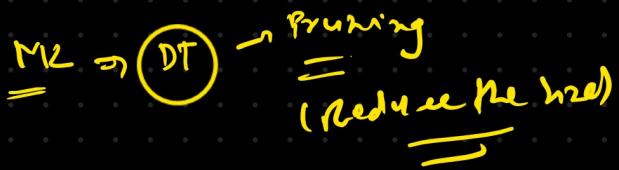


① unsupervised pretraining (web) (Self-Supervised)

② SFT

③ RHf → reference → OpU }  
→ PPo

Reduce the size.

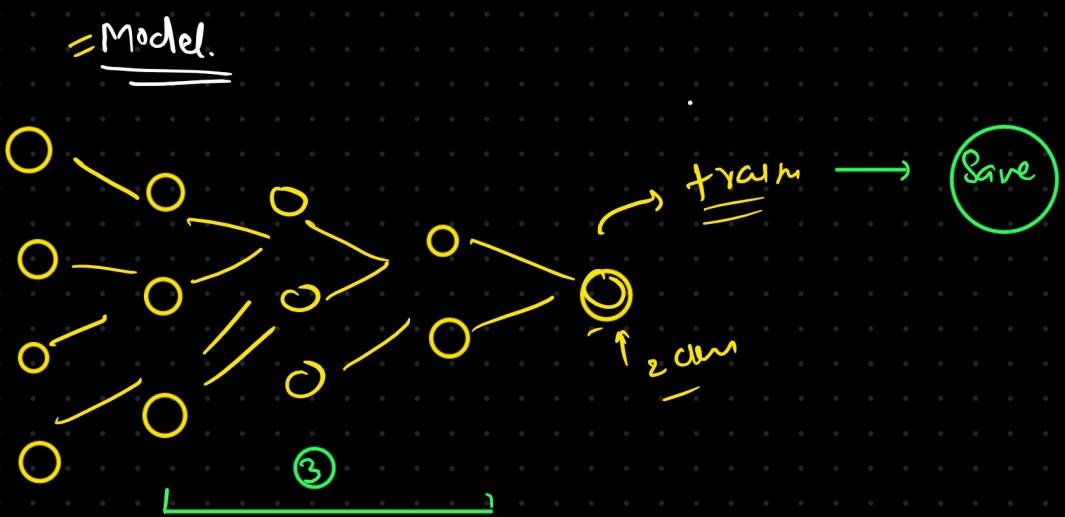


Quantization  $\Rightarrow$  Parameter ( $w, b$ )

Architecture

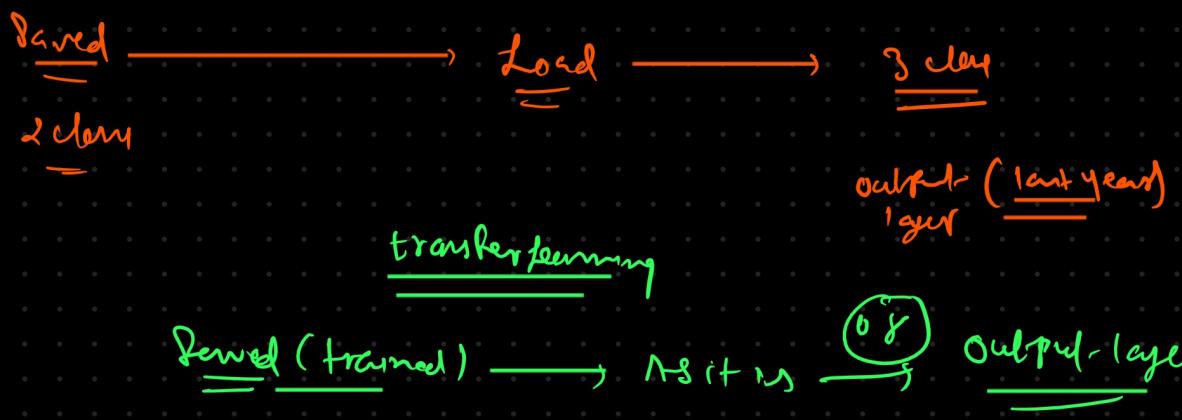
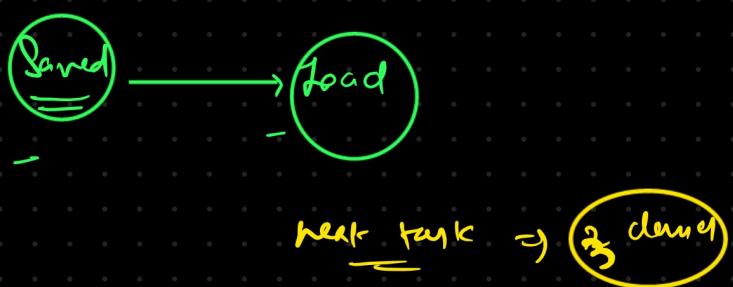
$w = \text{float-32}$   $\xrightarrow{=}$   
 $w = 23.72$   $\xrightarrow{(\text{float-32})}$   $= 12.8 \Rightarrow 1 \text{ byte}$   
 $\downarrow$   
 $1 \text{ byte} \approx 32 \text{ bits}$

float-32  
 int-16,32  
 $\downarrow$   
 memory  $\rightarrow$  RAM



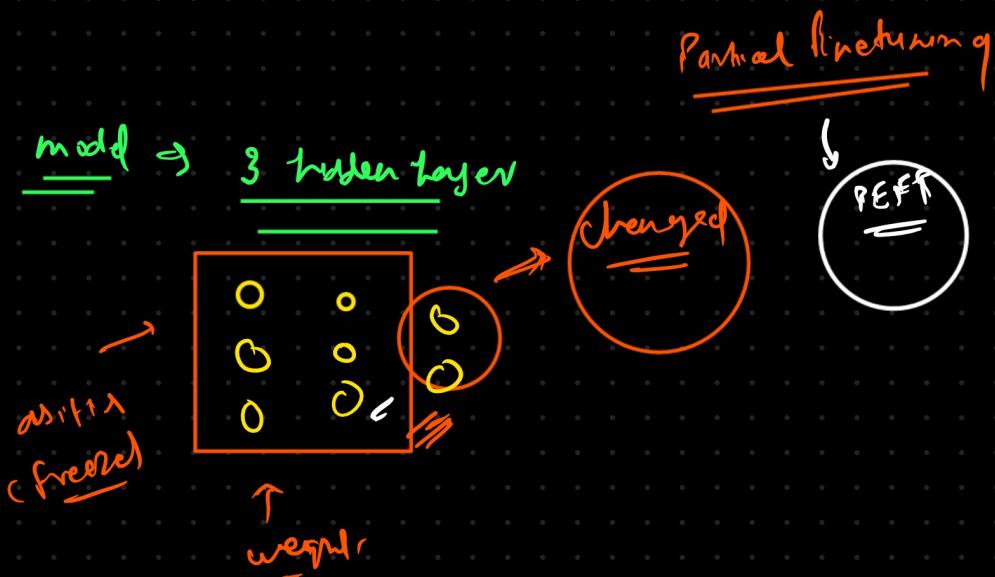
$$(4 \times 3 + 3) + (3 \times 3 + 3) + (3 \times 2 + 2) + (2 \times 1 + 1)$$

$$12 + 9 + 6 + 2 + 2 = 38$$

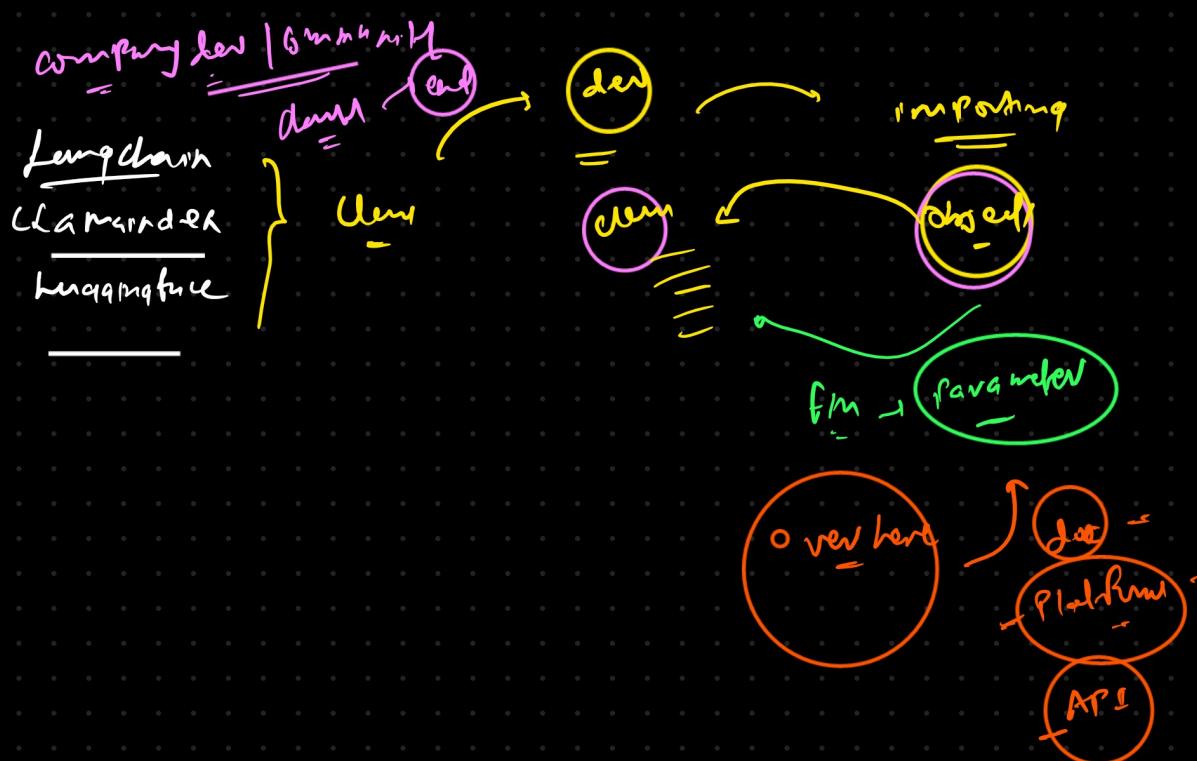


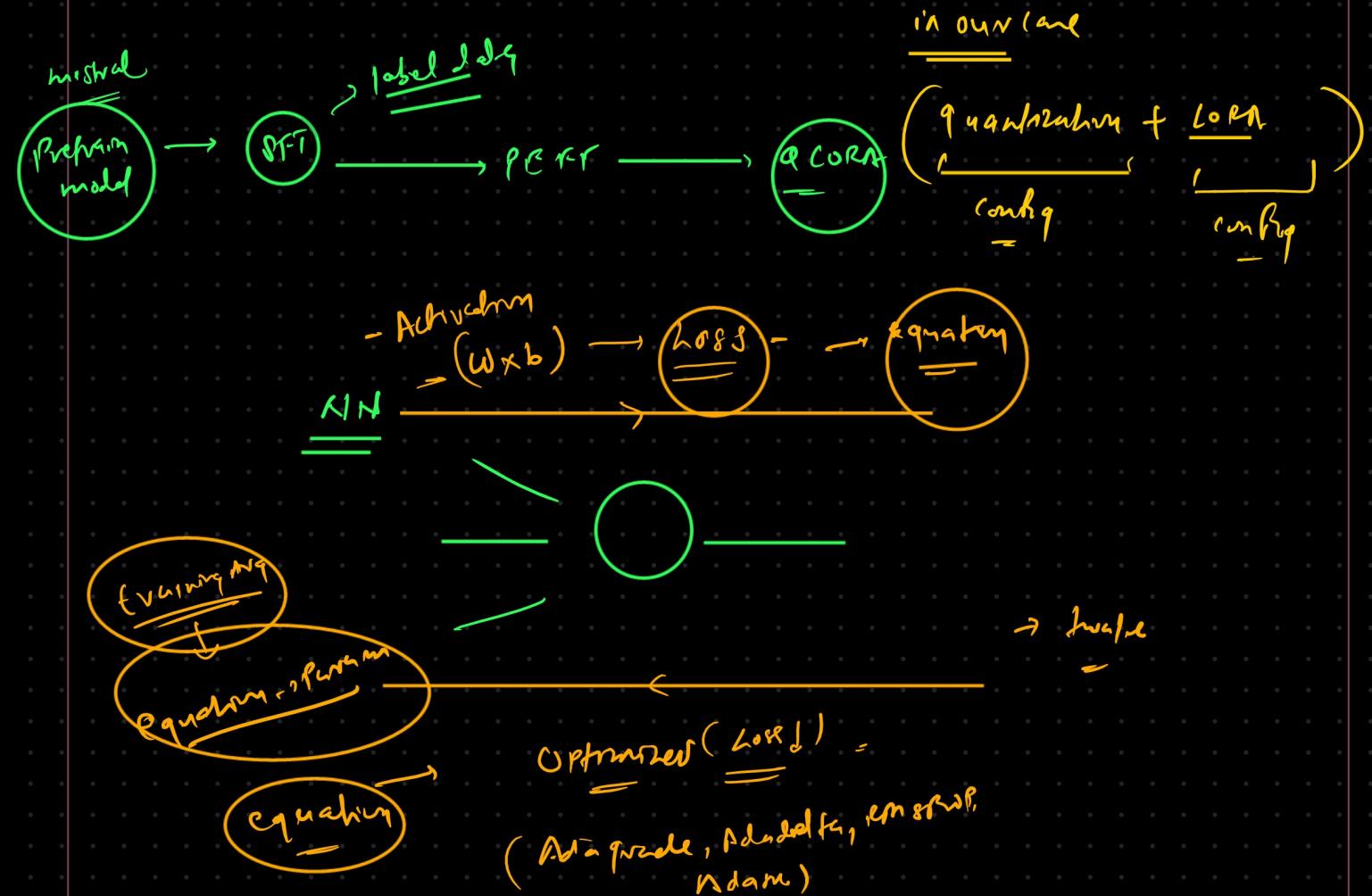
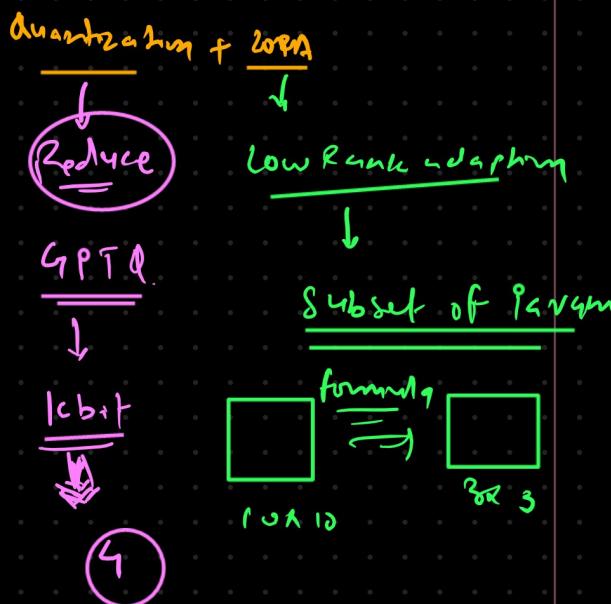
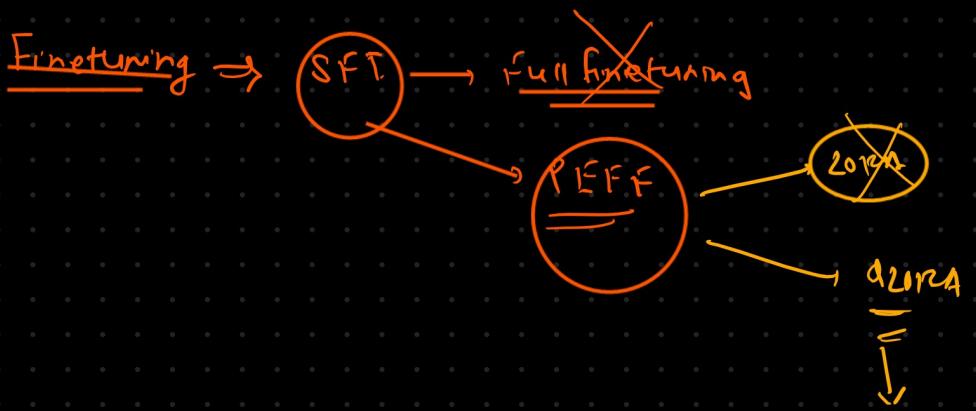
Rest of layer keep as it is (freezing)

## Transfer learning



3 layers → full finetuning





e Poch, batch size

- 1 Implement Lora in Bert finetuning
- 2 Perform  $\text{PEFT} \rightarrow \text{GLORA}$  over the Zephyr finetuning

optional Poison

{ create your own verb  
over the business file}

Lora 3.5 ratio  
ratio

token & load  $\Rightarrow$  inference  
data specific  
working fine or not

## Safegun

fine-tune  $\rightarrow$  generate info

only

Sir, Can we prevent pretrained model knowledge from finetuned model (for specific chatbot model)

(Prompting)