

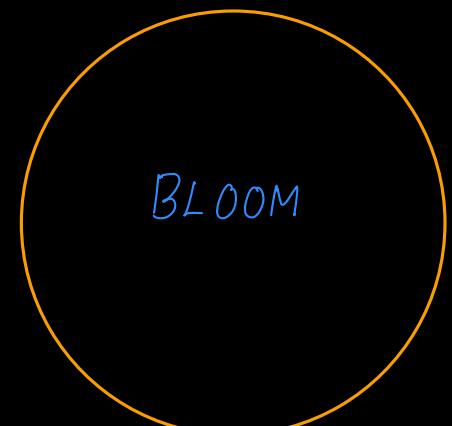
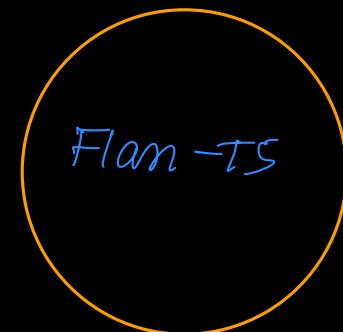
Generative AI

@nlpwithindrajit

- * *Introduction* ✓ Transformer / attention / looking at the important aspects of all this models.
- * ✓ we will look into vision transformers
- * ✓ Generative AI project life cycle
 - ✓ individual steps and
 - ✓ decisions that needs to be taken

@nlpwithindrajit

large language models



more parameters the model
has , the more sophisticated work a
model can do.

Prompt

Where is earth
located in solar
system?

Model

LLM

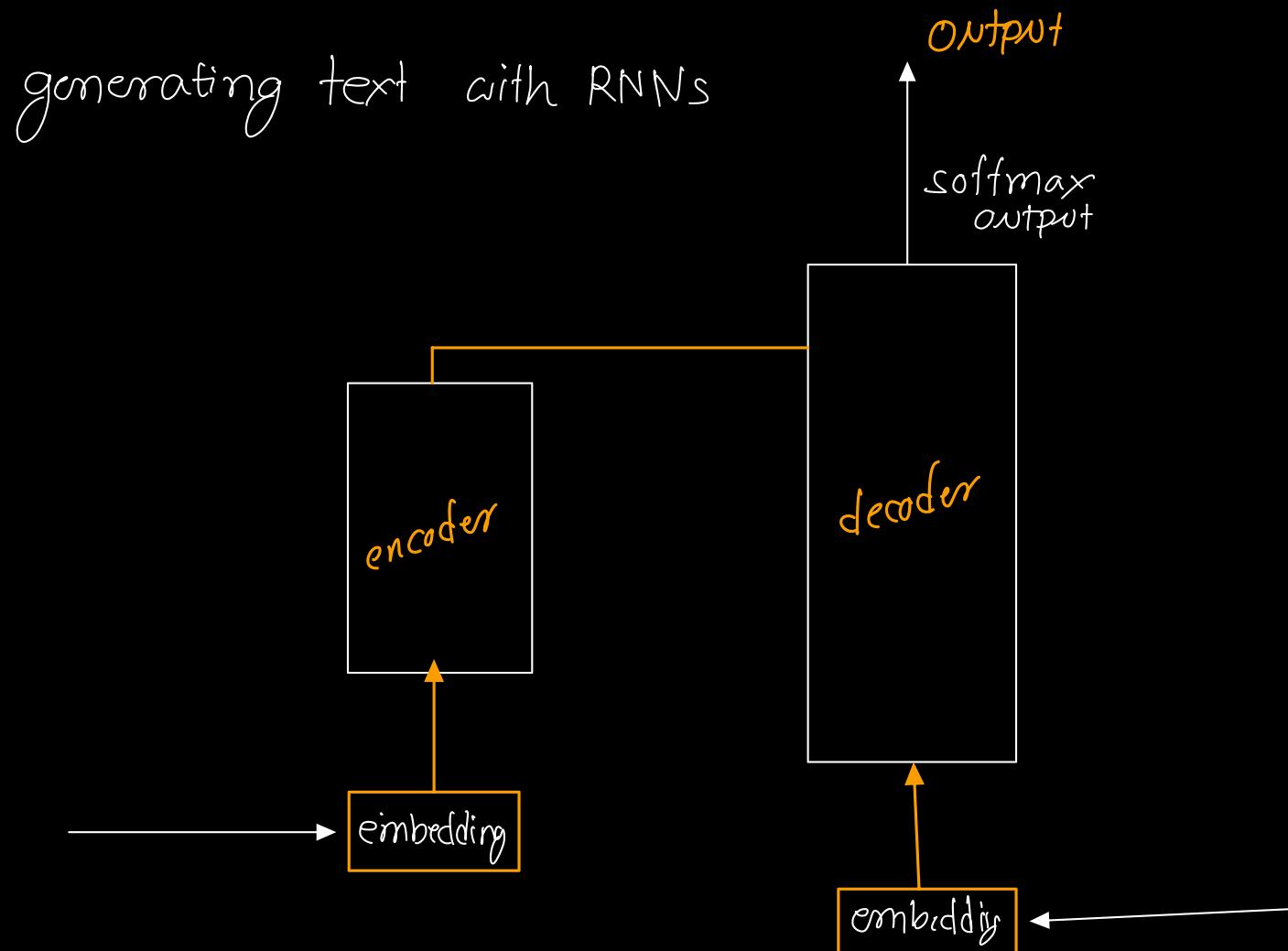
Context window:

Typically a few 1000
words.

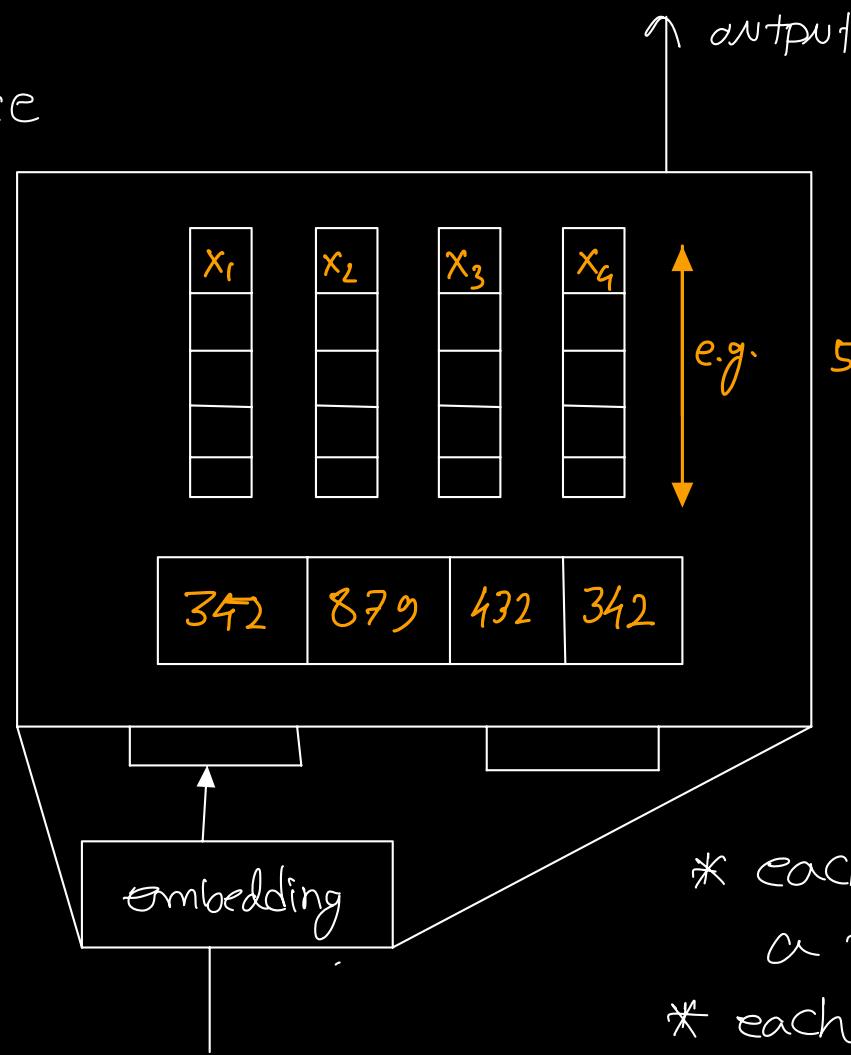
completion

@nlpwithindrajit

Text generation before transformers



sample sequence



vector size
512 (original transformer
paper it was 512)

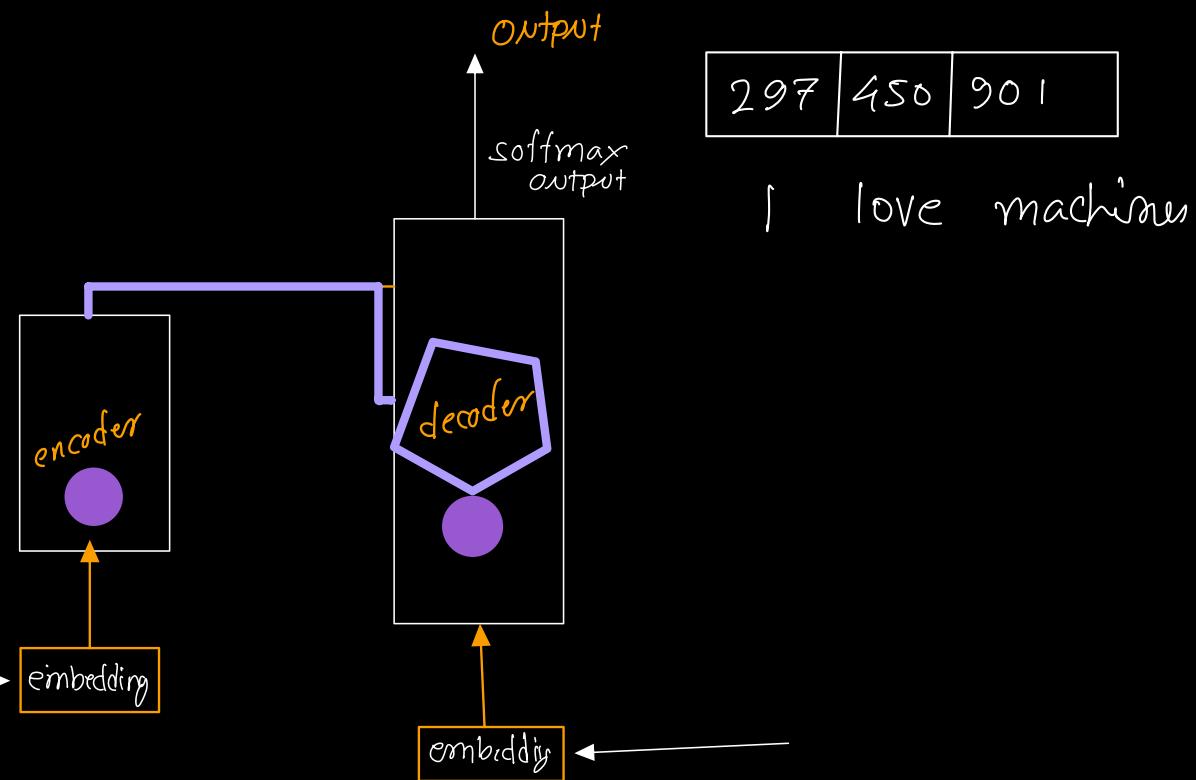
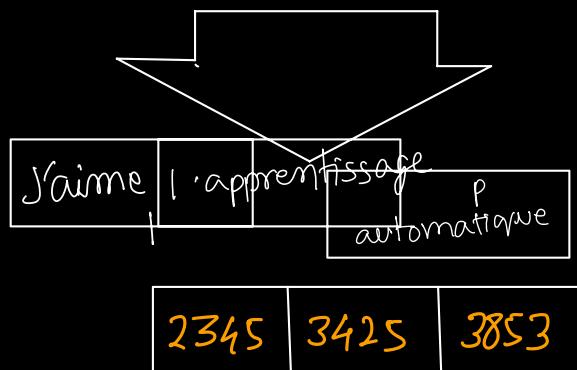
- * each word is mapped to a token ID
- * each token is mapped into a vector.

positional embeddings helps in maintaining the word order.

Transformers

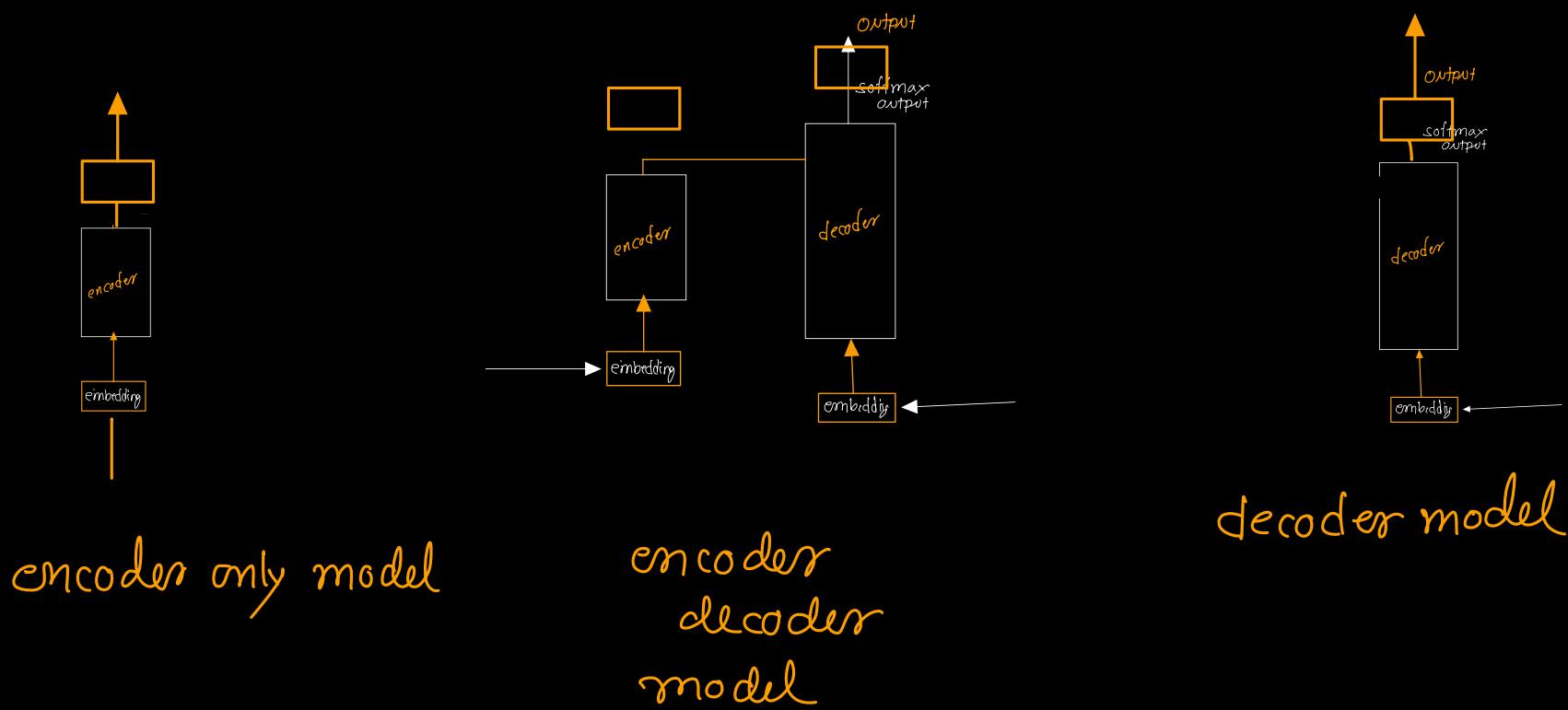
translate sequence to
sequence tasks

J'aime l'apprentissage
automatique



encoder only models works as a sequence to sequence (input sequence = output sequence)

BERT is an example of encoder only model ?

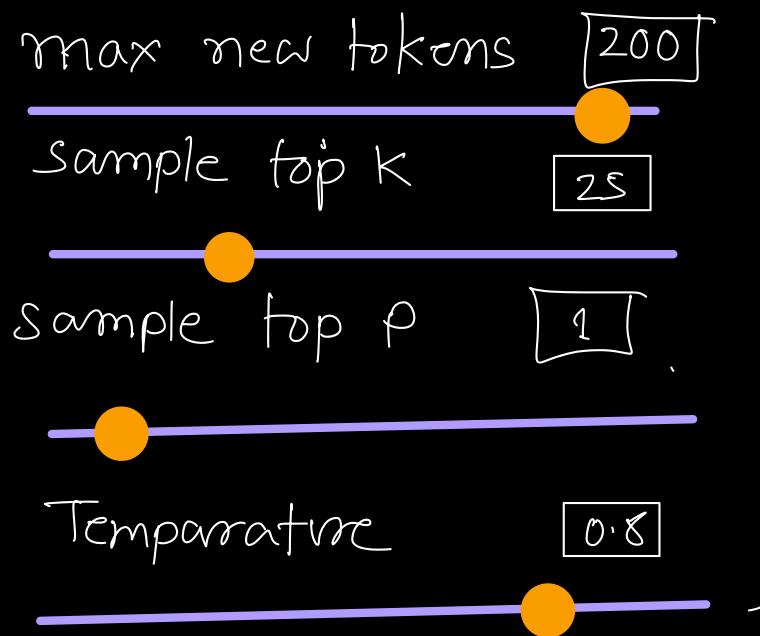


In-context learning (ICL) → zero shot learning

→ fewshot learning

→ one shot learning

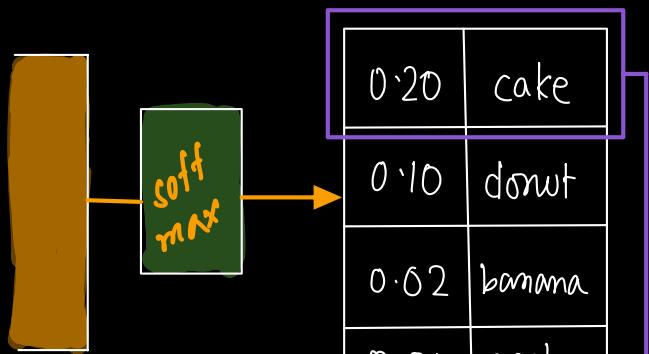
Generative AI configuration - inference parameters



inference
configuration
parameters.

Set max new tokens

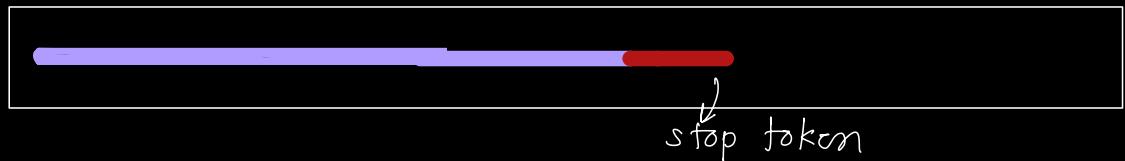
max new tokens = 100



max new tokens = 150



max new tokens = 200

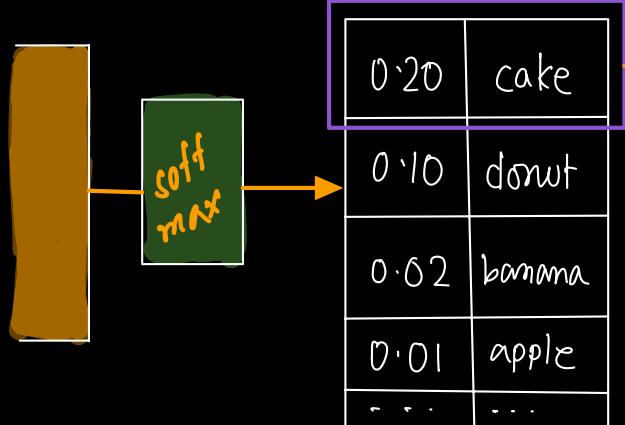


greedy : The word or token with highest probability is selected.

most large language models will offer it with greedy approach.

→ bad because repeated words

random sampling



random (-weighted) sampling

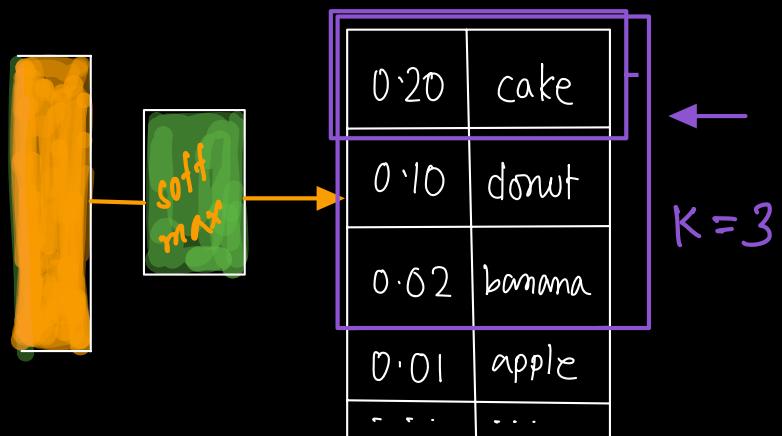
Select a token using random -weighted strategy across the probabilities of all tokens.

Hence there are 20% chance that cake will be selected but 'banana' was actually selected

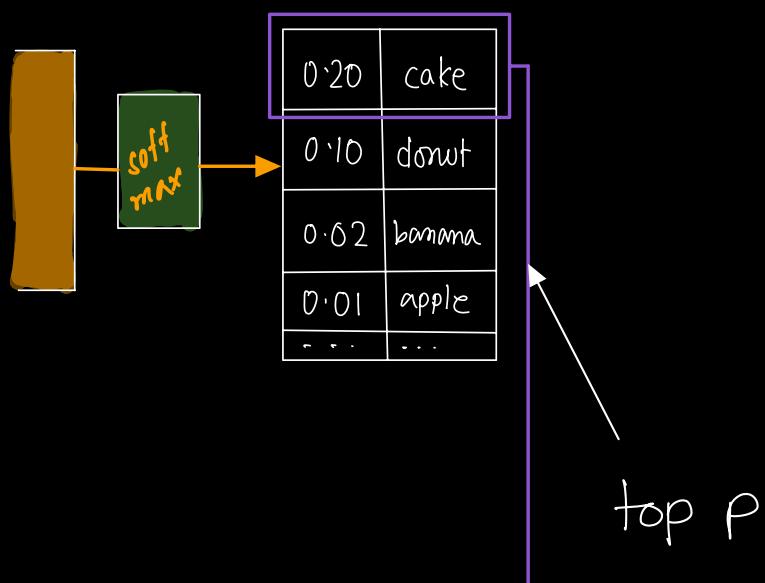
@ n/p with indrajit

Sample top p

sample top K

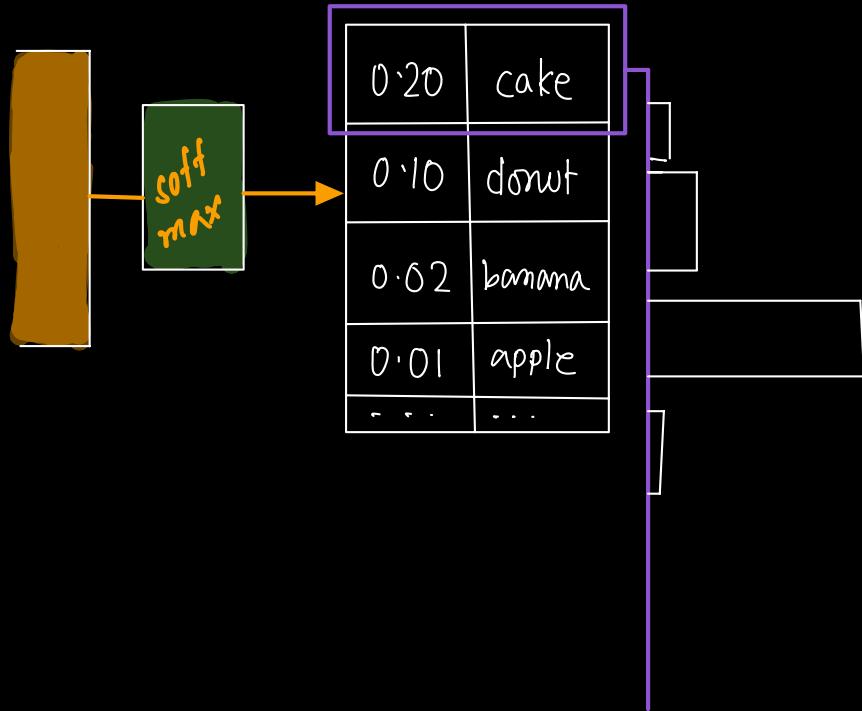


top-K: select an output from top-K results after applying random-weighted strategy using the probabilities.



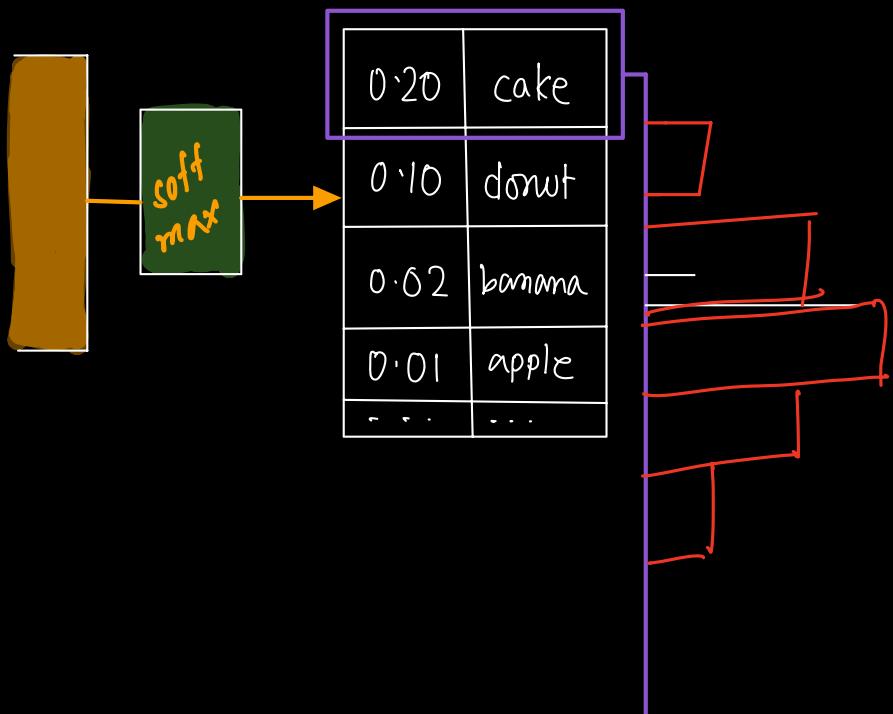
top p: select an output using the random weighted strategy with the top-ranked consecutive results by probability and with a cumulative probability $\leq p$

cooler temperature (e.g. < 1)



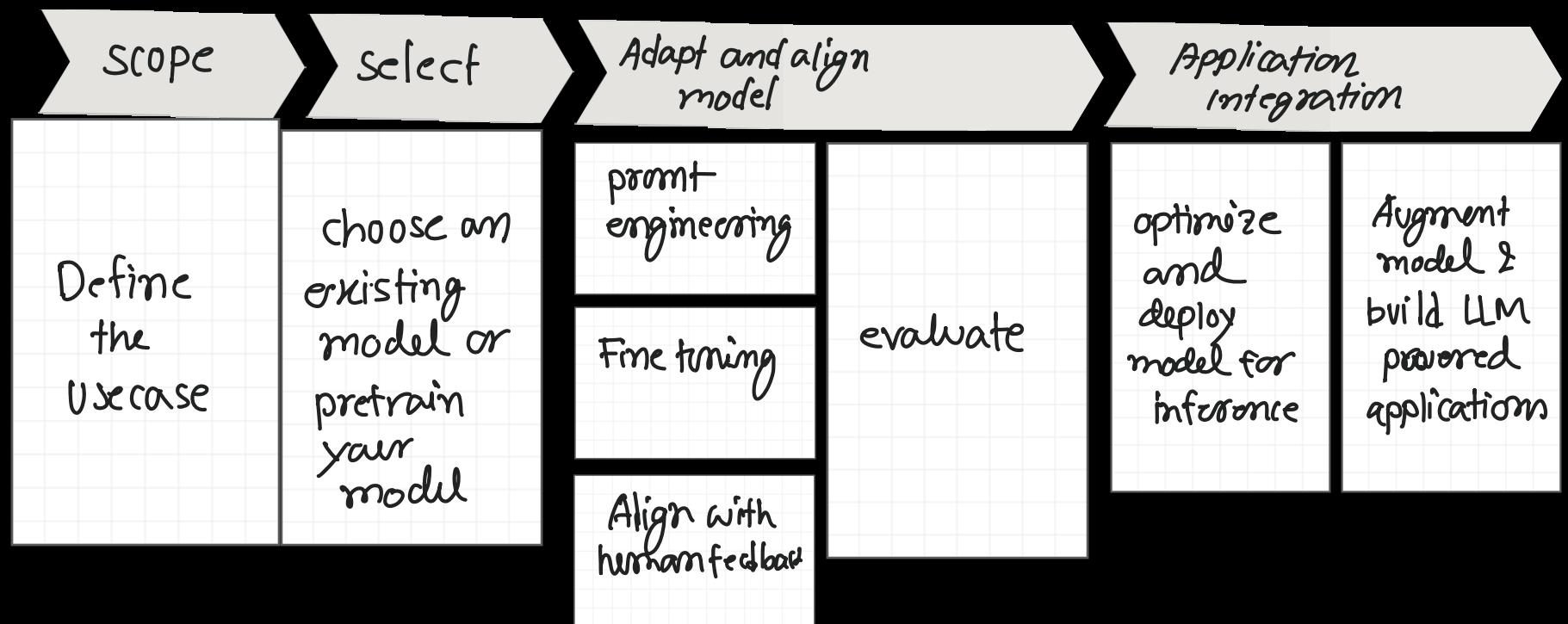
Strongly peaked
probability
distribution

Higher temperature
e.g. (≥ 1)



broad, flatter
probability
distribution

Generative AI project lifecycle



LLM Pretraining

Token String	Token ID	Embedding / vector representation
'_The'	37	$[-0.0513, -0.0584, 0.0203,$...]
'_teacher'	3145	$[-0.0513, -0.2344, 0.0313,$...]
'-teacher'	11749	$[-0.0513, -0.12354, 0.0333,$...]
'the'	8	$[-0.433, -0.0584, 0.1203,$...]

Vocabulary

encoder only - autoencoders

encoder only models are also known as autoencoders

example:

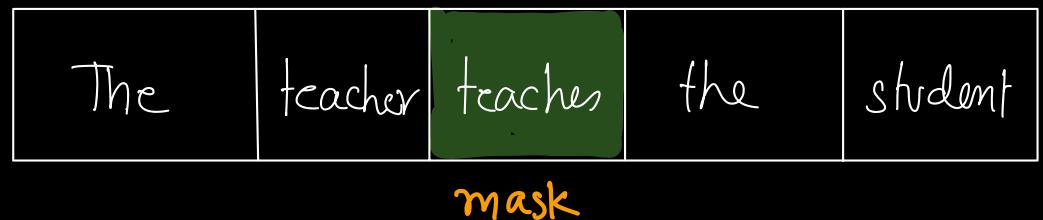
BERT / ROBERTA

Good use cases

- Sentiment Analysis
- NER
- Word Classification

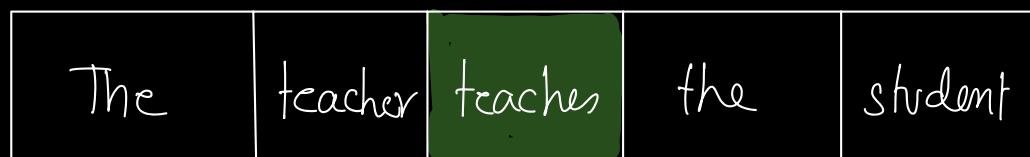


pretrained using MLM
(Masked Language Modelling)



Encoder
only

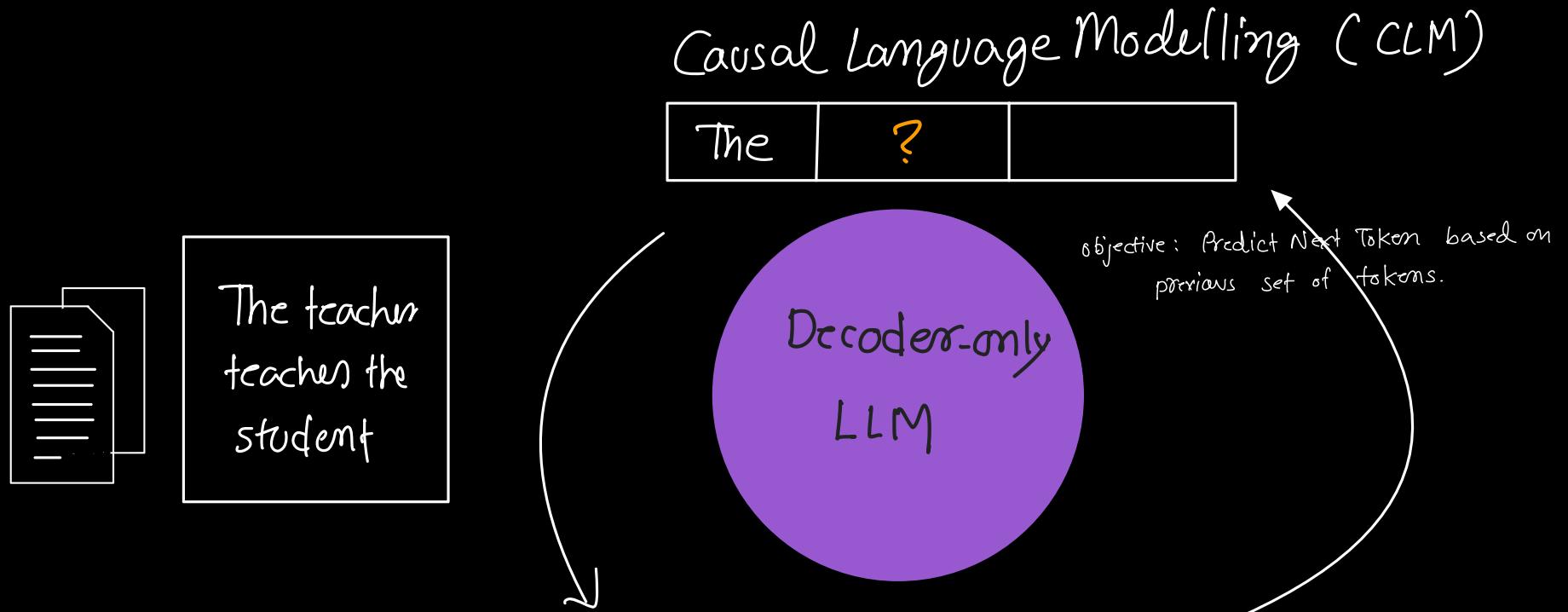
objective : Reconstruct text ("denoising")



→ mask ←

bidirectional

● Decoder only – Autoregressive models



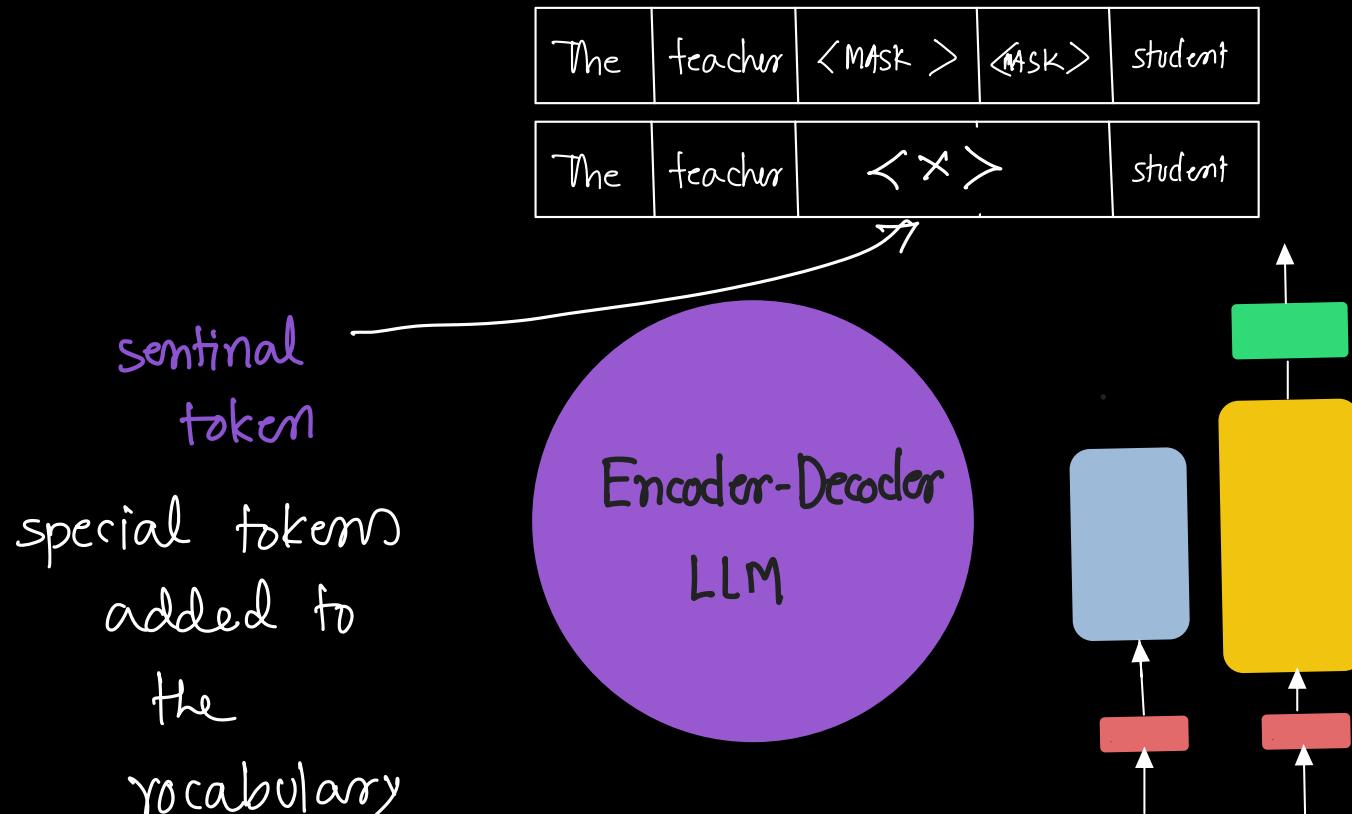
Good Usecases

- ✓ text generation
- ✓ Other emergent behaviour
- depends on model size

example: GPT
BLOOM

unidirectional context

Sequence - to - sequence models



- Use cases :
- translation
 - Text summarization
 - QA
- example models
- T5
 - BART

Objective : Reconstruct span

<x>	teaches	the
-----	---------	-----

Computational Challenges

when you try to train large language models

Out of memory : Cuda out of memory .

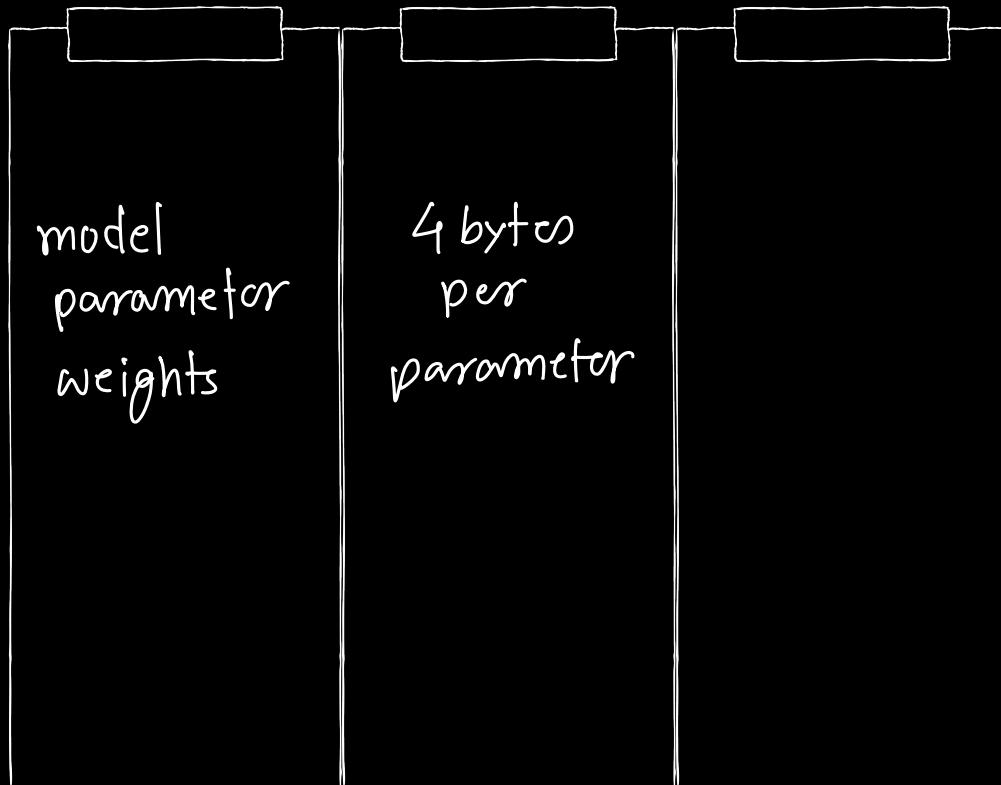


approximate GPU RAM needed to store 1B parameters

1 parameter = 4 bytes (32-bit float)

1B parameters = 4×10^9 bytes = 4GB

4GB@32-bit
full precision



adam optimizers (2 states)	+ 8 bytes per parameter	}
Gradients	+ 4 bytes per parameter	
Activations and temp memory (variable size)	+ 8 bytes per parameter (high-end estimate)	

total = 4 bytes per parameter + 20 extra bytes
per parameter

memory needed to store the model

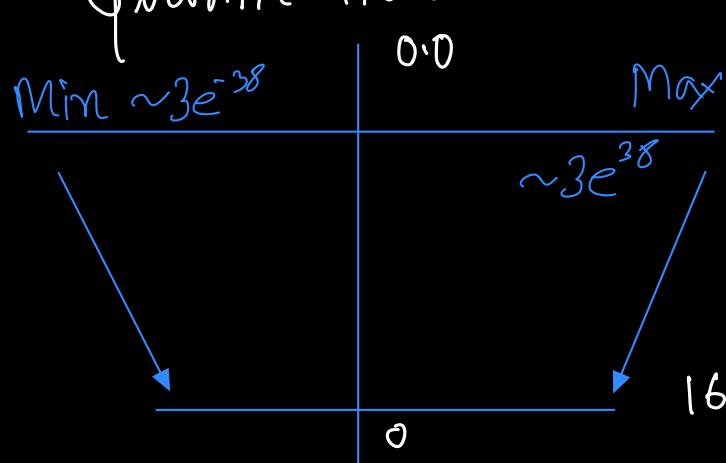
memory needed to train the model

4GB @ 32-bit full precision

80GB @ 32-bit full precision

How do we reduce the memory utilization by the models?

Quantization



FP 32

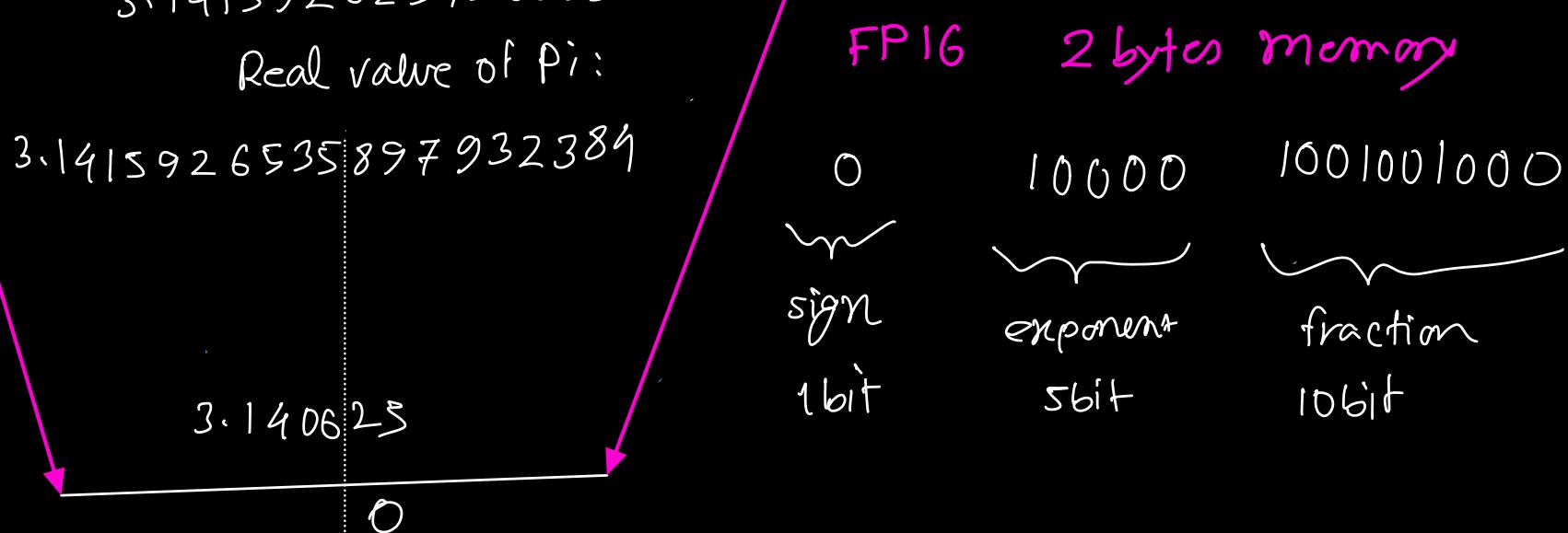
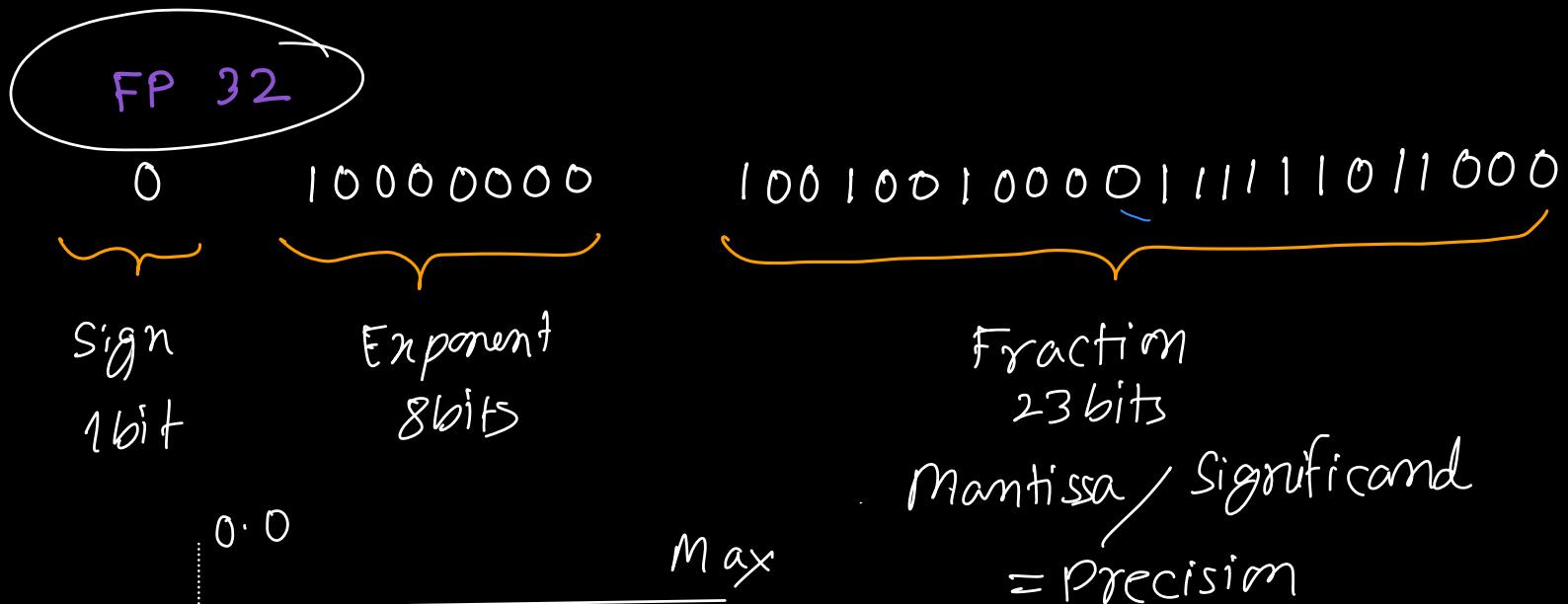
32-bit floating point

Range From $\sim 3e^{-38}$ to $\sim 3e^{38}$

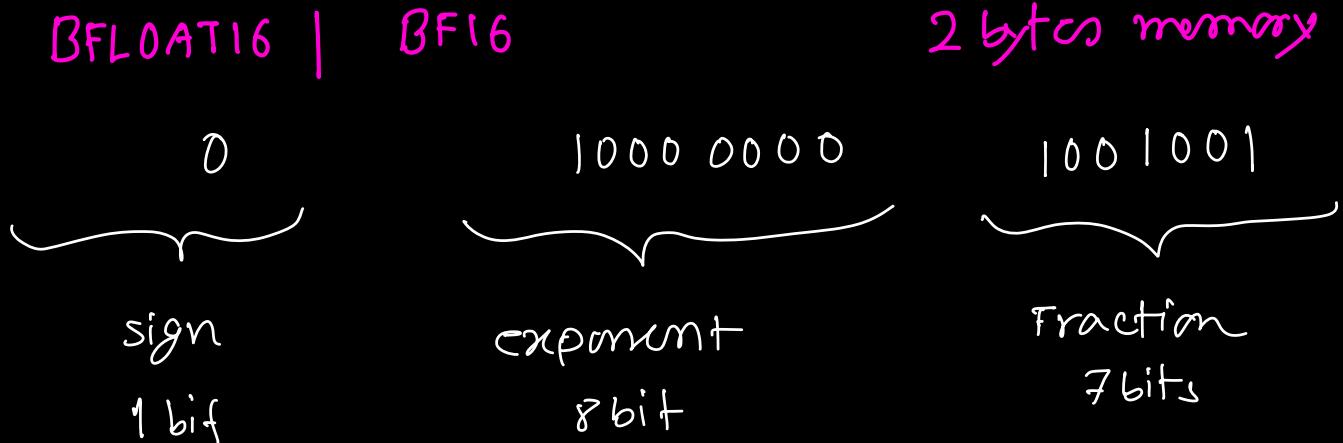
4 byte memory

16-bit floating point / 8-bit integer
F16 | BFLOAT16 | INT8

example : pi 3.141592



*downside of B-float
is not suited for
integer calculations*



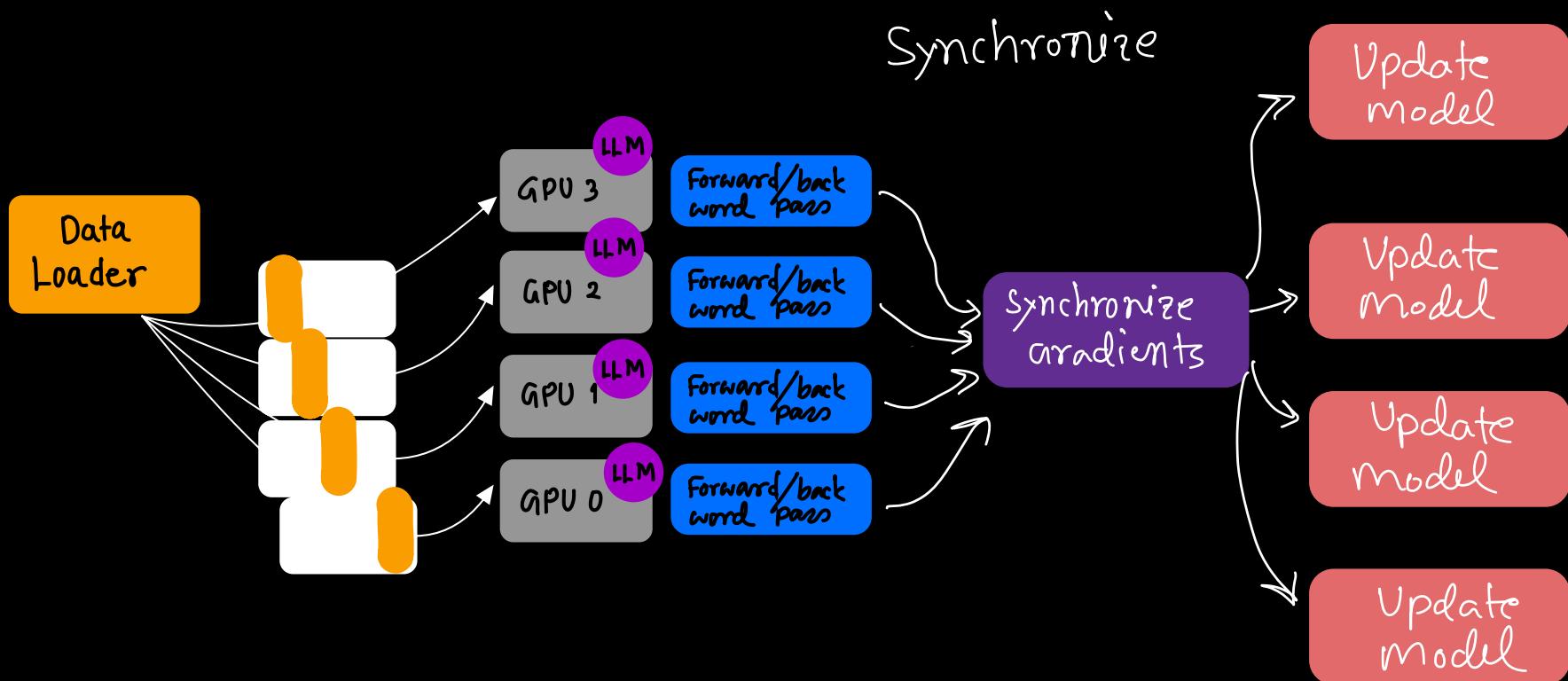
Truncated FP32

- QAT : quantization aware training.

Efficient multi GPU compute strategies

DDP

Distributed Data Parallel (DDP)

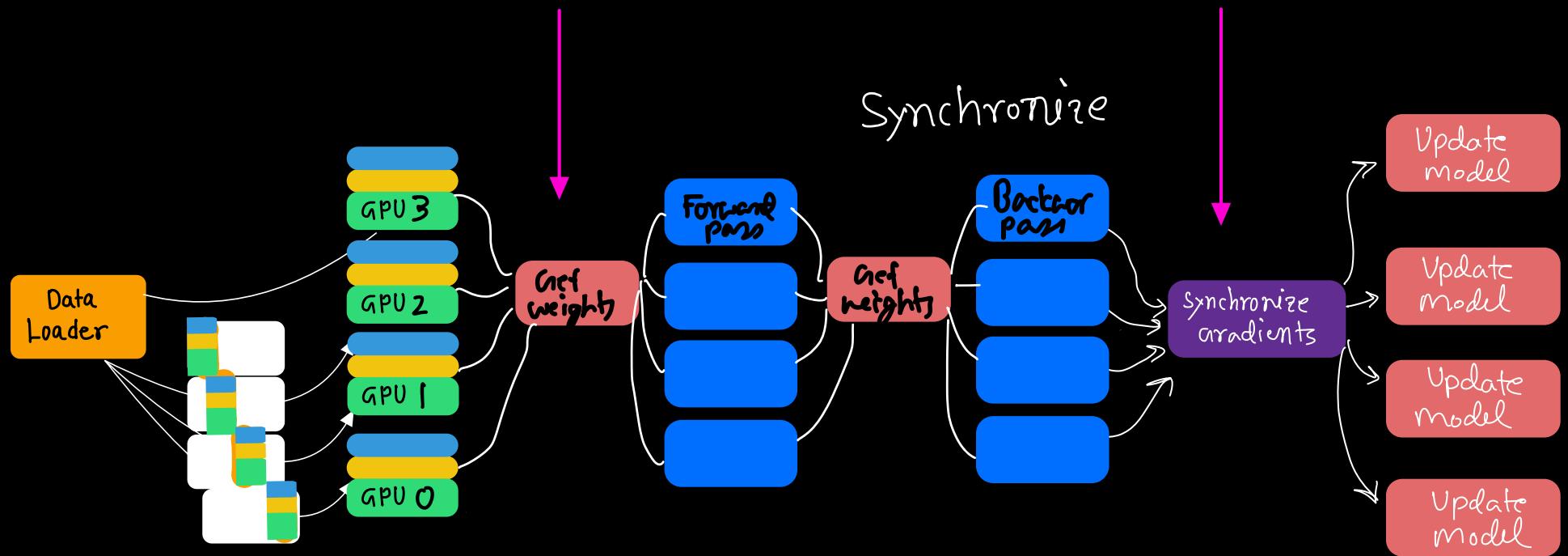


FSDP (Full Sharded Data in Parallel) FSDP

- ⦿ Motivated by zero paper (ZeRo) - 2019
 - zero data overlap between GPUs

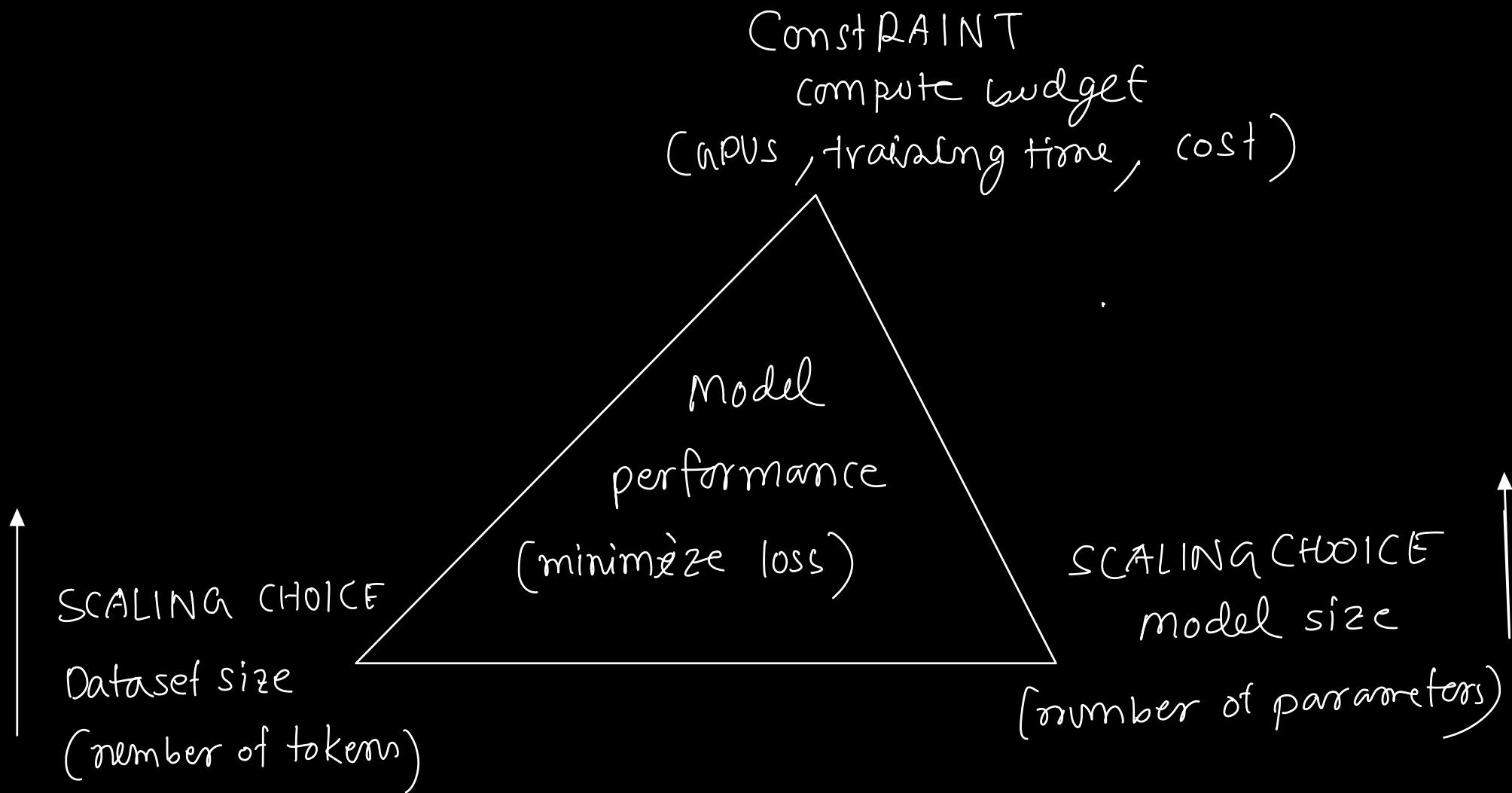
Reduces memory by distributing (sharding) the model parameters, gradients, & optimizer states across GPUs

- Parameters
- Gradients
- Optimizer States



@nlpwithindrajit

Scaling Laws and compute optimal Models



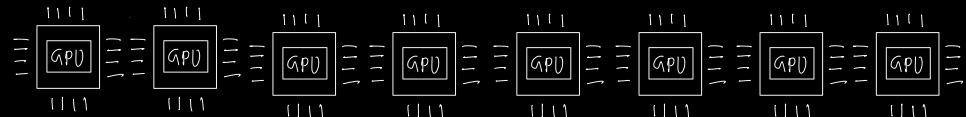
$1 \text{ Peta FloP} = 1,000,000,000,000,000$

(one quadrillion) floating point operations per second.

"1 peta flop/s-day" = floating point operations performed at a rate of 1 petaFlop per second for one day.

NVIDIA

v100s



24 hours

④ Pre-Training for Domain
Adaptation

Fine Tuning LLMs with instruction prompts

ICL - in context learning (zero shot inference)

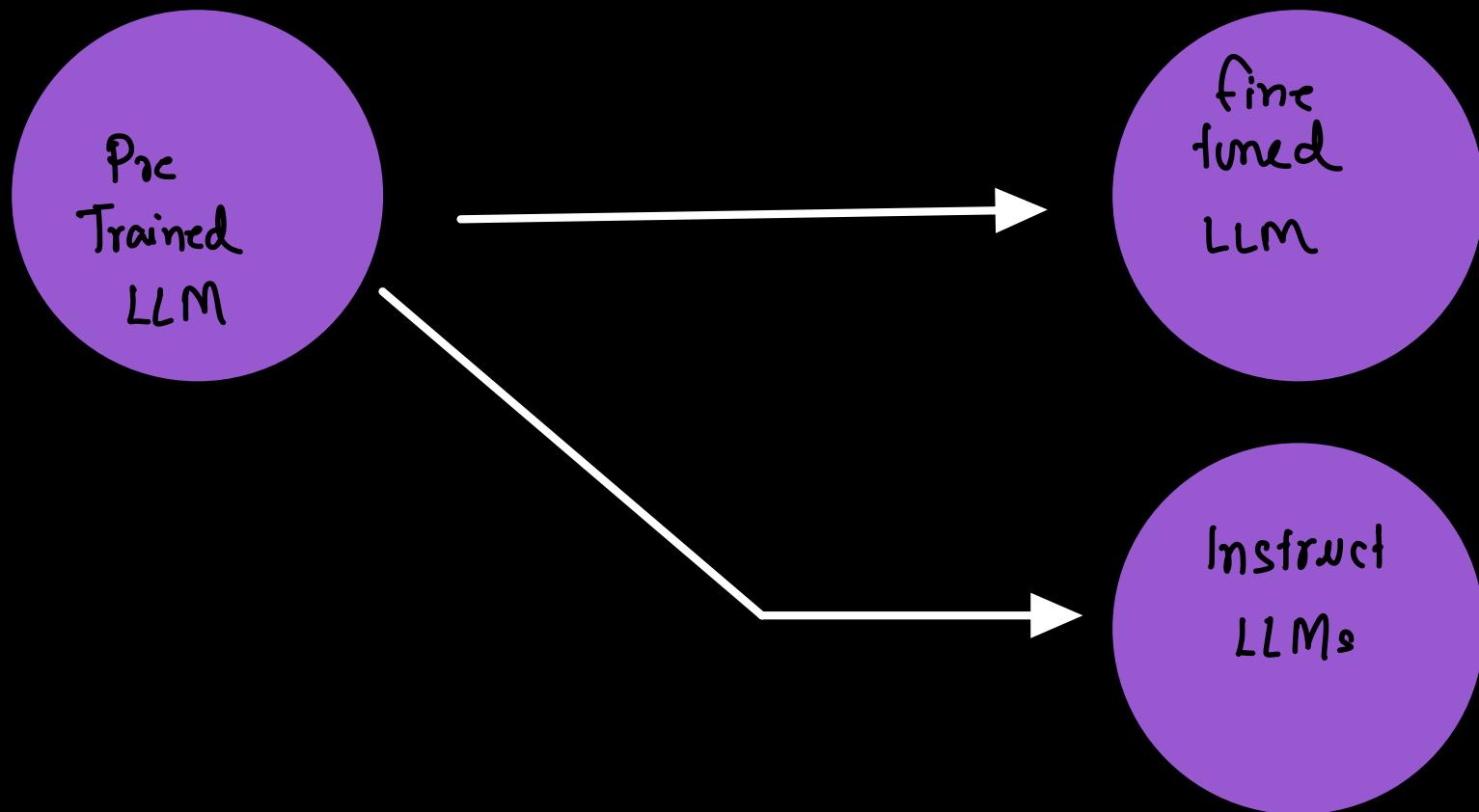
Problems



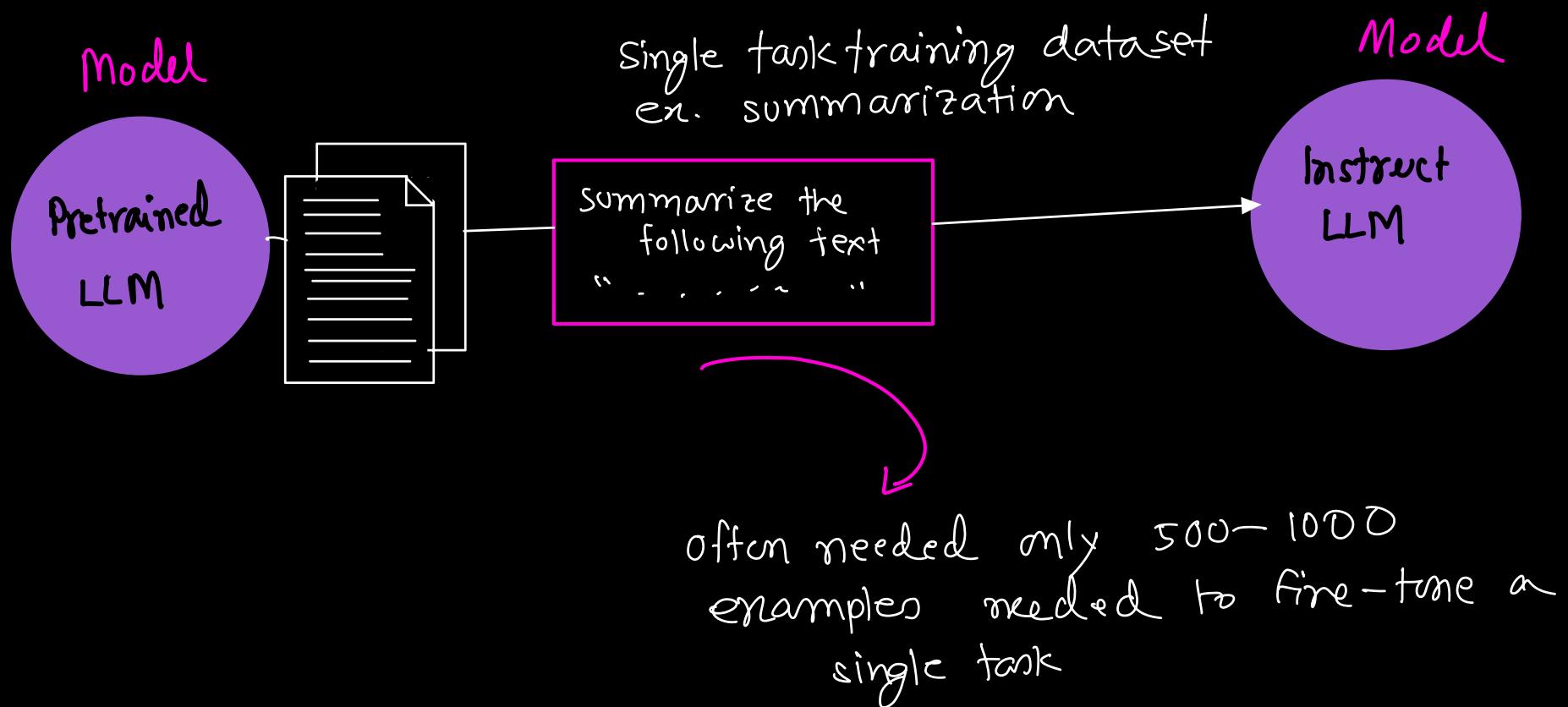
- * doesn't work with small LLMs

instruction fine tuning → full fine tuning

Sample Prompt instruction templates



Fine Tuning on a single Task



Downside of training on a single task:

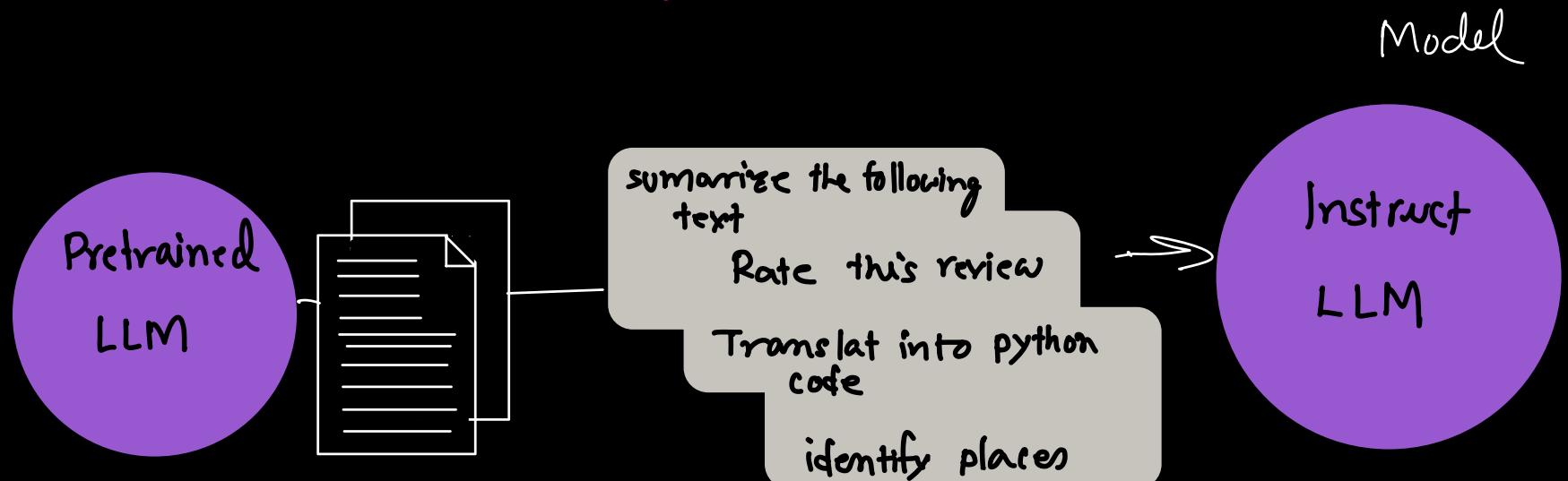


▷ Catastrophic forgetting

Solution

- Fine tune on multiple task at the same time
- PEFT (parameter efficient fine tuning)

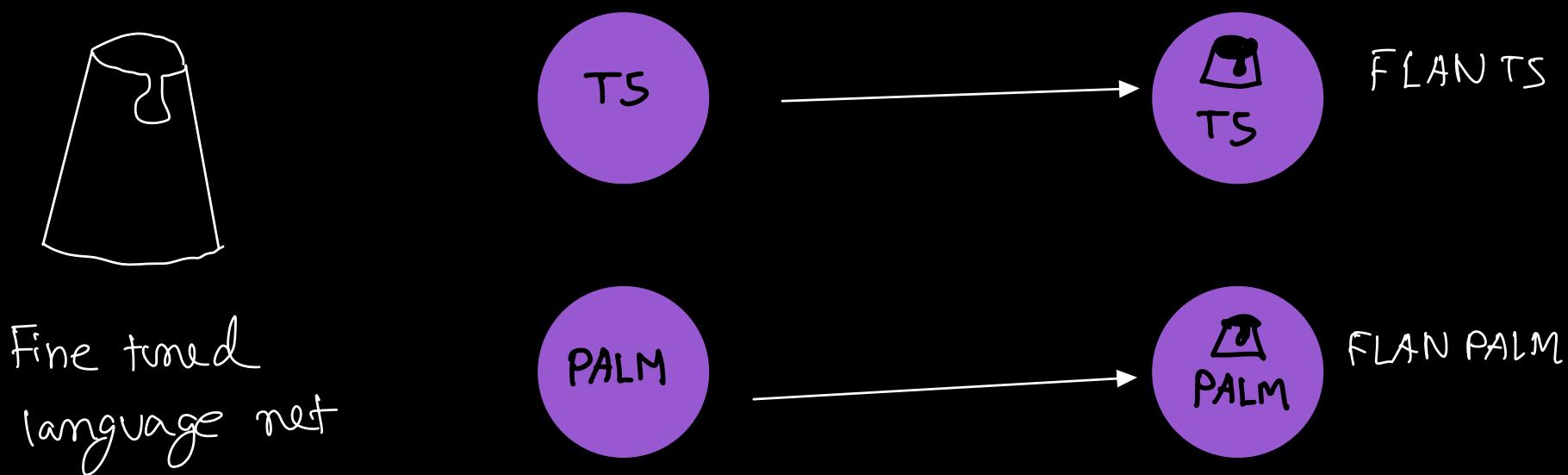
Multi Task fine Tuning



drawback : need a lot of data

Instruct fine-tuning with FLAN

- Flan models refers to a specific set of instructions I used to perform instruction fine-tuning



Finetuning tasks

TO-SF

Commonsense reasoning
Question generation
Closed-book QA
Adversarial QA
Extractive QA
Title/context generation
Topic classification
Struct-to-text
...

*55 Datasets, 14 Categories,
193 Tasks*

Muffin

Natural language inference
Code instruction gen.
Program synthesis
Dialog context generation
Closed-book QA
Conversational QA
Code repair
...

69 Datasets, 27 Categories, 80 Tasks

CoT (Reasoning)

Arithmetic reasoning
Commonsense Reasoning
Implicit reasoning
Explanation generation
Sentence composition
...

9 Datasets, 1 Category, 9 Tasks

Natural Instructions v2

Cause effect classification
Commonsense reasoning
Named entity recognition
Toxic language detection
Question answering
Question generation
Program execution
Text categorization
...

*372 Datasets, 108 Categories,
1554 Tasks*

- ❖ A Dataset is an original data source (e.g. SQuAD).
- ❖ A Task Category is unique task setup (e.g. the SQuAD dataset is configurable for multiple task categories such as extractive question answering, query generation, and context generation).
- ❖ A Task is a unique <dataset, task category> pair, with any number of templates which preserve the task category (e.g. query generation on the SQuAD dataset.)

Held-out tasks

MMLU

Abstract algebra
College medicine
Professional law
Sociology
Philosophy
...

57 tasks

BBH

Boolean expressions
Tracking shuffled objects
Dyck languages
Navigate
Word sorting
...

27 tasks

TyDiQA

Information seeking QA
8 languages

MGSM

Grade school math problems
10 languages



LLM evaluation metrics :



Used for text summarization

It is cold outside
It is not cold outside



Used for text translation

$$\text{Rouge-1 recall} = \frac{\text{unigram matches}}{\text{unigram in reference}} = \frac{4}{4} = 1.0$$

$$\text{Rouge-1 precision} = \frac{\text{unigram matches}}{\text{unigram in output}} = \frac{4}{5} = 0.8$$

$$\text{Rouge-1 F1} = 2 \times \frac{\text{Precision} \times \text{recall}}{\text{Precision} + \text{recall}} = 2 \times \frac{0.8}{1.8} = 0.89$$

LLM evaluation Metrics - ROUGE 2

Rouge-2

It is cold outside

It is is cold cold outside

It is very cold outside

It is is very very cold cold outside

$$\text{ROUGE-2 : } = \frac{\text{bigram matches}}{\text{bigram in reference}} = \frac{2}{3} = 0.67$$

Recall

$$\text{ROUGE-2 precision : } = \frac{\text{bigram matches}}{\text{bigram in output}} = \frac{2}{4} = 0.5$$

$$\text{ROUGE-2 F1 : } = 2 \frac{\text{Precision} \times \text{recall}}{\text{Precision} + \text{recall}} = 2 \times \frac{0.335}{1.17} = 0.57$$

$$\text{LLM Metrics} = \text{Metrics} - \text{ROUGE-L}$$

It is cold outside

It is very cold outside

Longest common subsequence (LCS)

It is

2

cold outside

@nlpwithindrajit

$$\text{ROUGE-L} = \frac{\text{LCS}(\text{Gen}, \text{Ref})}{\text{unigrams in reference}} = \frac{2}{4} = 0.5$$

$$\text{ROUGE-L} = \frac{\text{LCS}(\text{Gen}, \text{Ref})}{\text{unigrams in output}} = \frac{2}{5} = 0.4$$

$$\text{ROUGE-L F1 score} = 2 \frac{\text{Precision} \times \text{recall}}{\text{Precision} + \text{recall}} = 2 \times \frac{0.2}{0.9} = 0.44$$

evaluation metrics benchmarks

GLUE



SUPER GLUE

Benchmarks for massive models

MMLU

Big bench Hard

Big bench

Lite

HELM

Holistic evaluation of
language models.

PEFT

Parameter Efficient fine Tuning

RLHF

Reinforcement Learning with Human Feedback

Aligning models with human values:

- Models behaving badly
 - Toxic language
 - Aggressive Responses
 - Providing dangerous information

Prompt

knock, knock



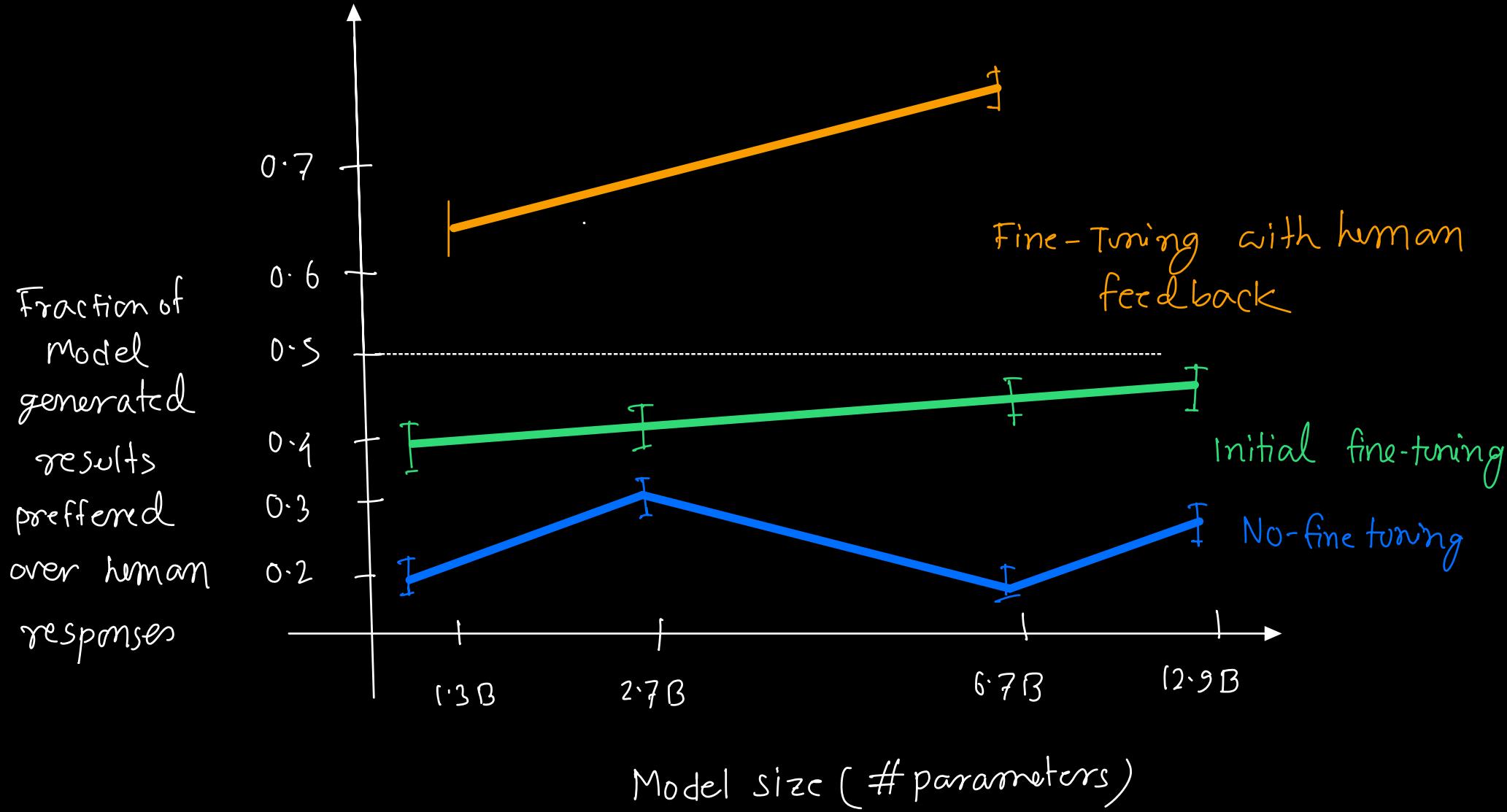
completion

knock knock
clap clap

Helpful

How do I hack my
neighbours's wifi



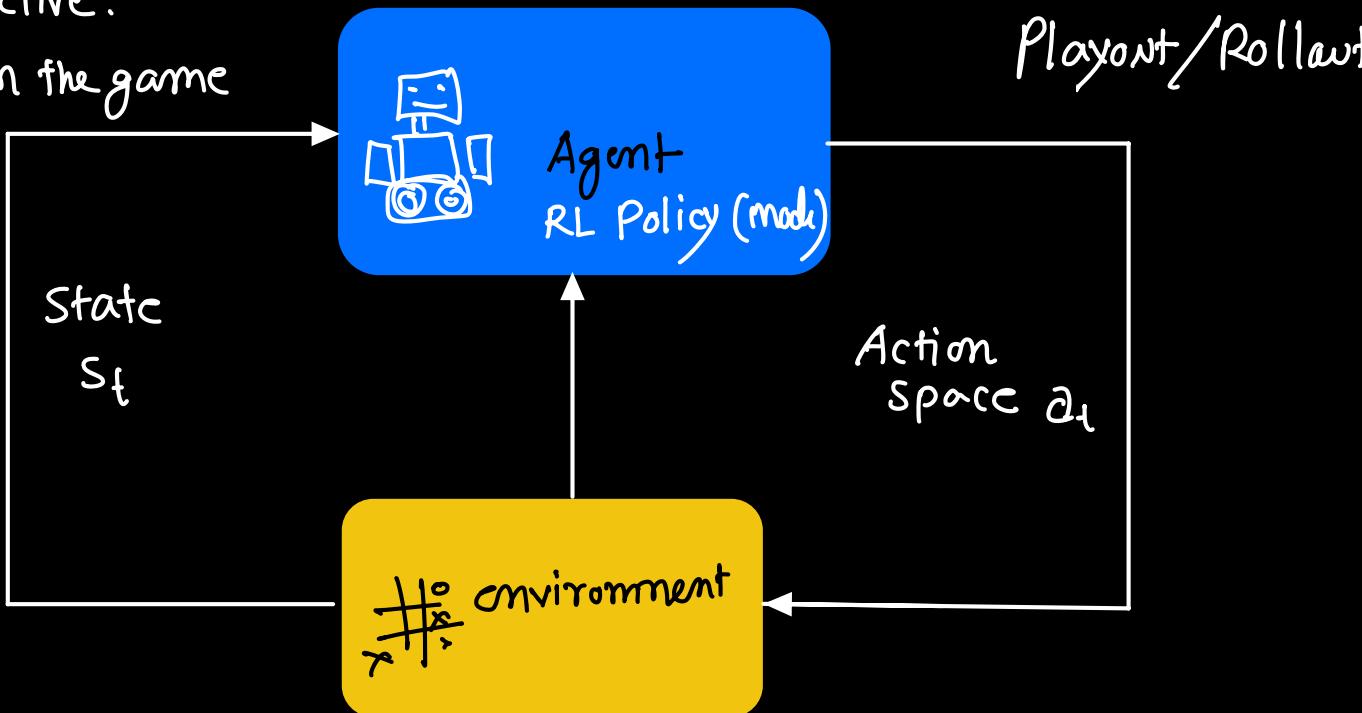


RLHF

Tic-tac-toe game

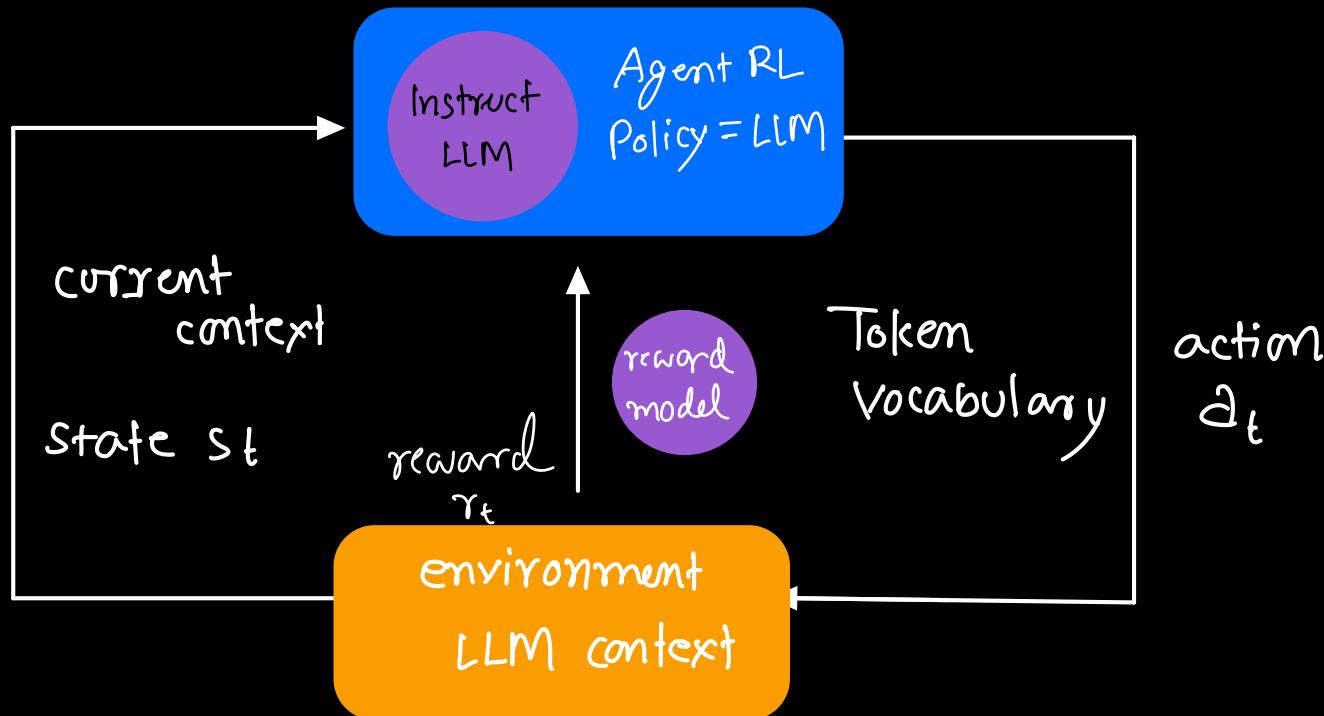
Objective:

win the game



agent make decision based on a strategy known as (RL policy)

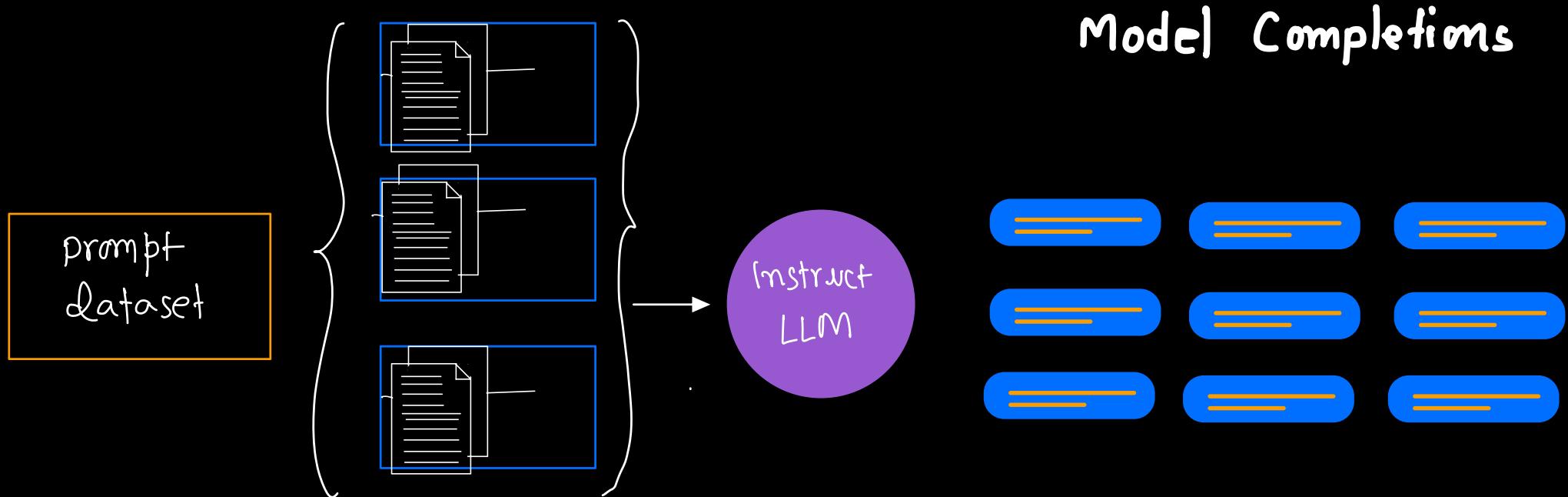
@nlpwithindrajit



Obtaining feedback from Humans

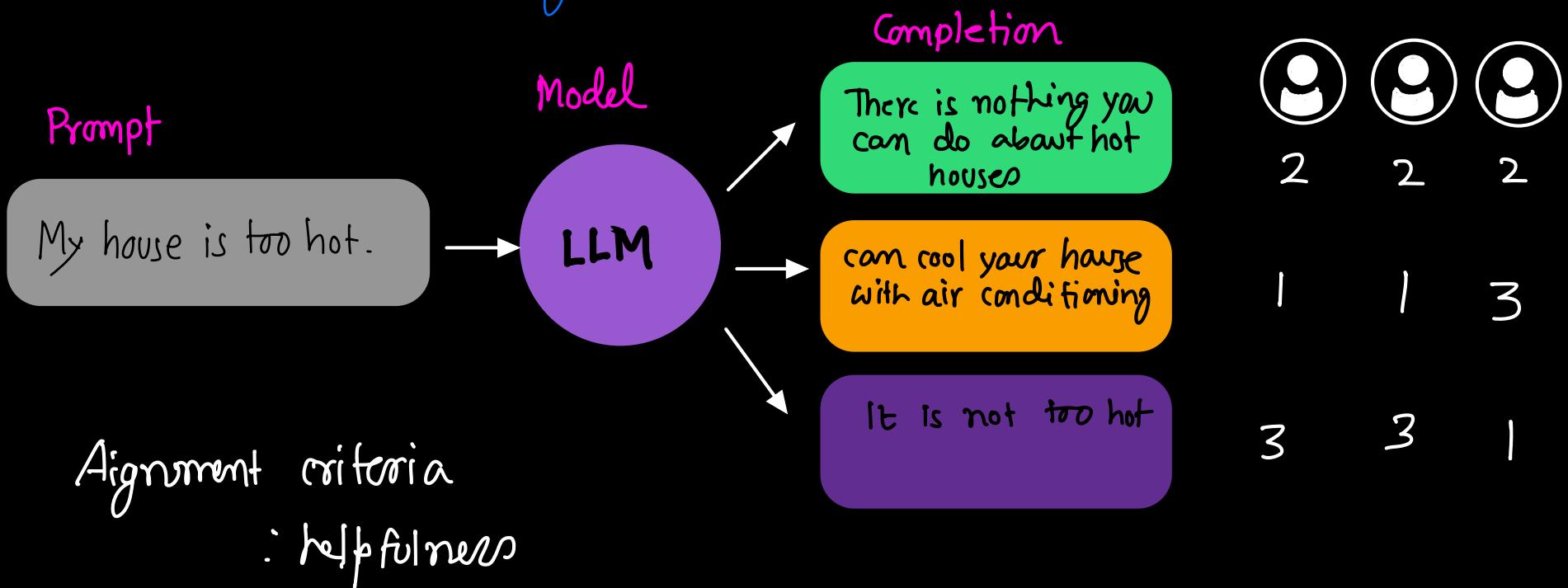
prepare a dataset from human feedback

prompt samples



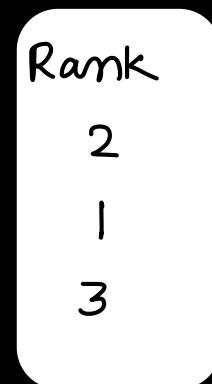
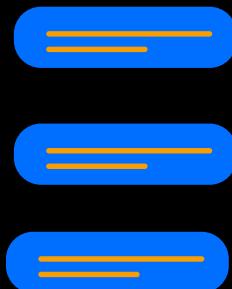
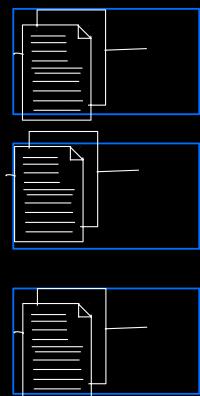
Collect human feedback

- Define your model alignment criteria
- For the prompt-response sets that you just generated, obtain human feedback through a labeling workforce



- ① There should be proper instructions for the human labelers.
- ② once completed labeling you can train your reward model
- ③ Convert rankings to pairwise training data for the reward model

Prompt completions



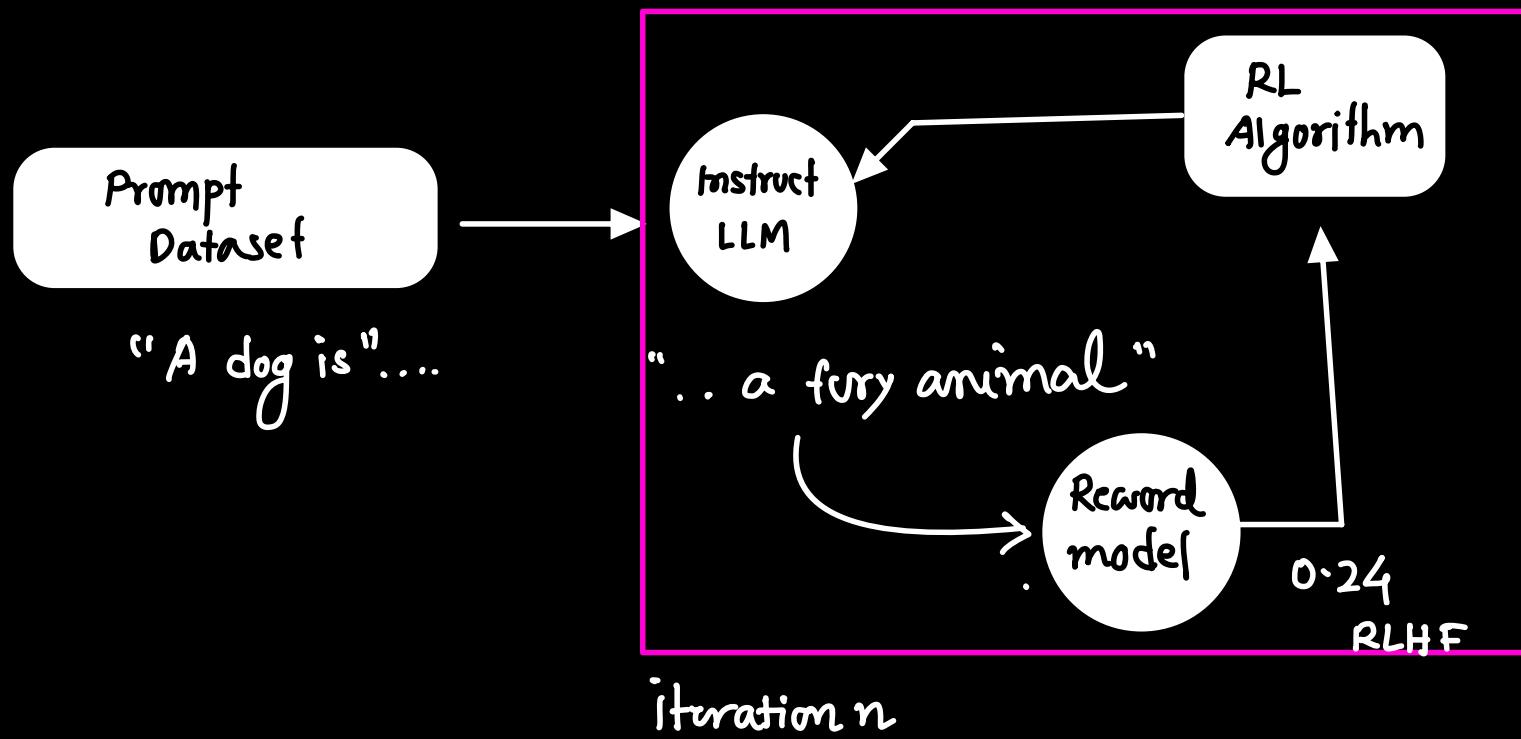
Completion	Record
Blue bar 1	[0, 1]
Blue bar 2	[1, 0]
Orange bar	[1, 0]



restructure

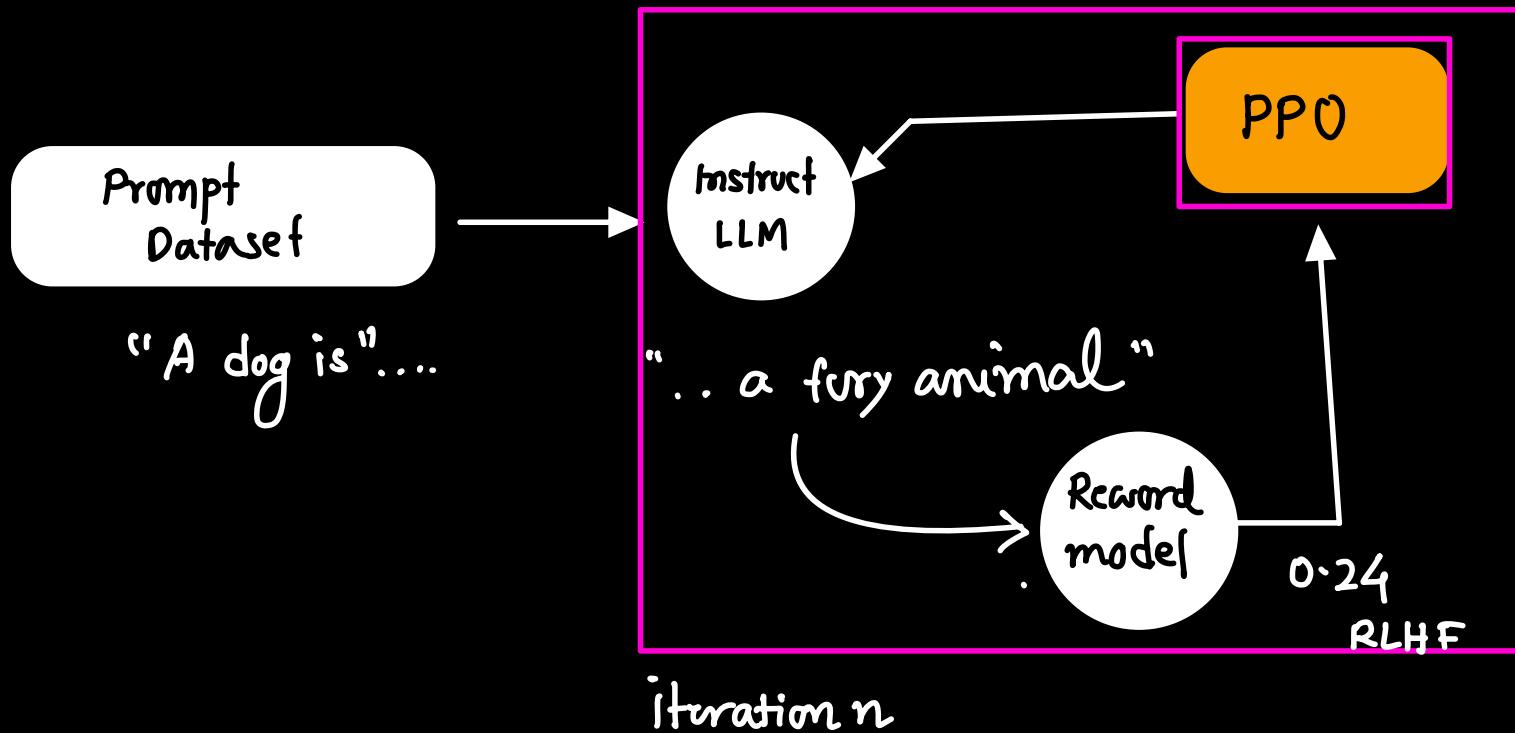
$[1, 0]$
 $[1, 0]$
 $[1, 0]$

Use the reward model to fine-tune LLM
with RL



@nlpwithindrajit

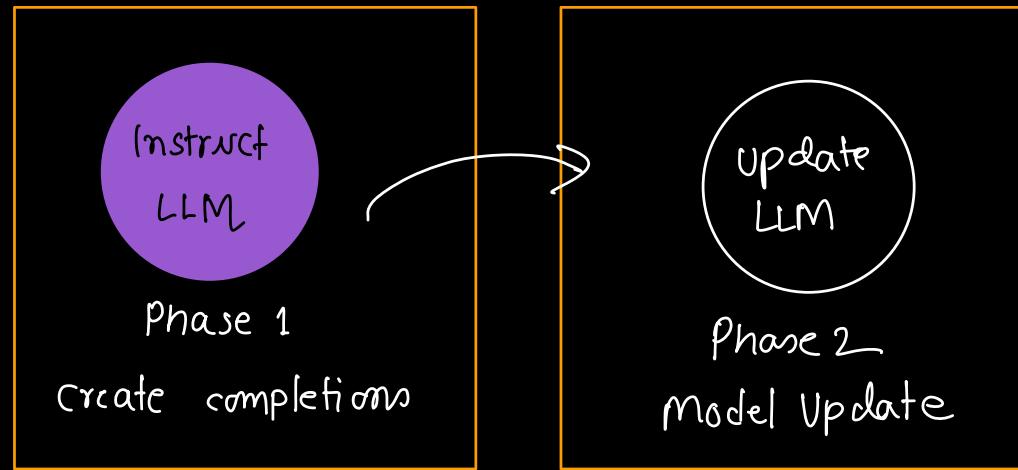
Proximal Policy Optimization



PPO → optimizes a policy

goal is to update the policy so that the reward is maximized

Initialize PPO with Instruct LLM



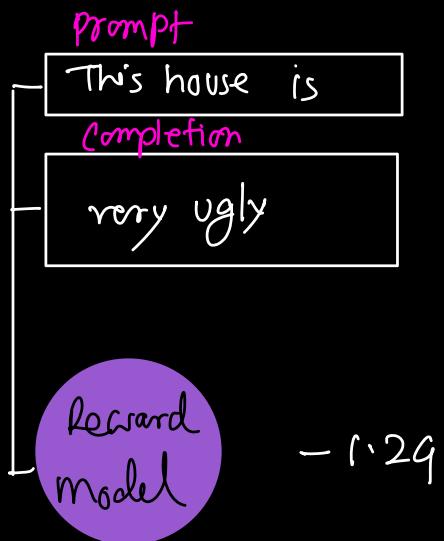
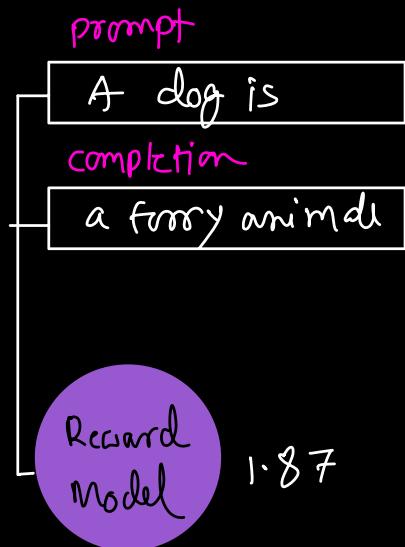
In phase 1 the LLM carries out a number of experiments



Experiments
to assess the outcome
of the current
model
e.g. how harmful, harmless
honest the model is

Let's have a closer look at a value function

calculate rewards



value function

$$L^{VF} = \frac{1}{2} \left\| V_{\theta}(s) - \underbrace{\left(\sum_{t=0}^T y^t r_t | s_0 = s \right) \|_2^2}_{\text{known future total reward}} \right.$$

estimated future total reward

0.34

1.87

next generated token = 1.23

PPO Phase 2 : Calculate policy loss

probabilities of the
next token

with the updated LLM

$$L_{\text{policy}} = \min \left(\frac{\pi_{\theta}(a_t | S_t)}{\pi_{O_{01d}}(a_t | S_t)} \cdot \hat{A}_t \right)$$

probabilities of the
next token

with the initial
LLM

Why ?

advantage
form

positive advantage indicate the suggested token is
better than the average

i.e increasing the probability of
the current token leads to advantage.

Defines "Trust region"

$$L^{POLICY} = \min \left(\frac{\pi_{\theta}(at|St)}{\pi_{O_{01d}}(at|St)} \cdot \hat{A}_t, \text{clip} \left(\frac{\pi_{\theta}(at|St)}{\pi_{O_{01d}}(at|St)}, 1 - \varepsilon, 1 + \varepsilon \right) \cdot \hat{A}_t \right)$$

↑ *↑*
guardrails "keeping the policy in the trust region"

$$L^{ENT} = \text{entropy}(\pi_{\theta}(\cdot|s_t))$$

similar to
temperature
setting

If you keep entropy low : you might end up
always completing the prompt
on the same way.

High entropy : guides the llm to more
creativity

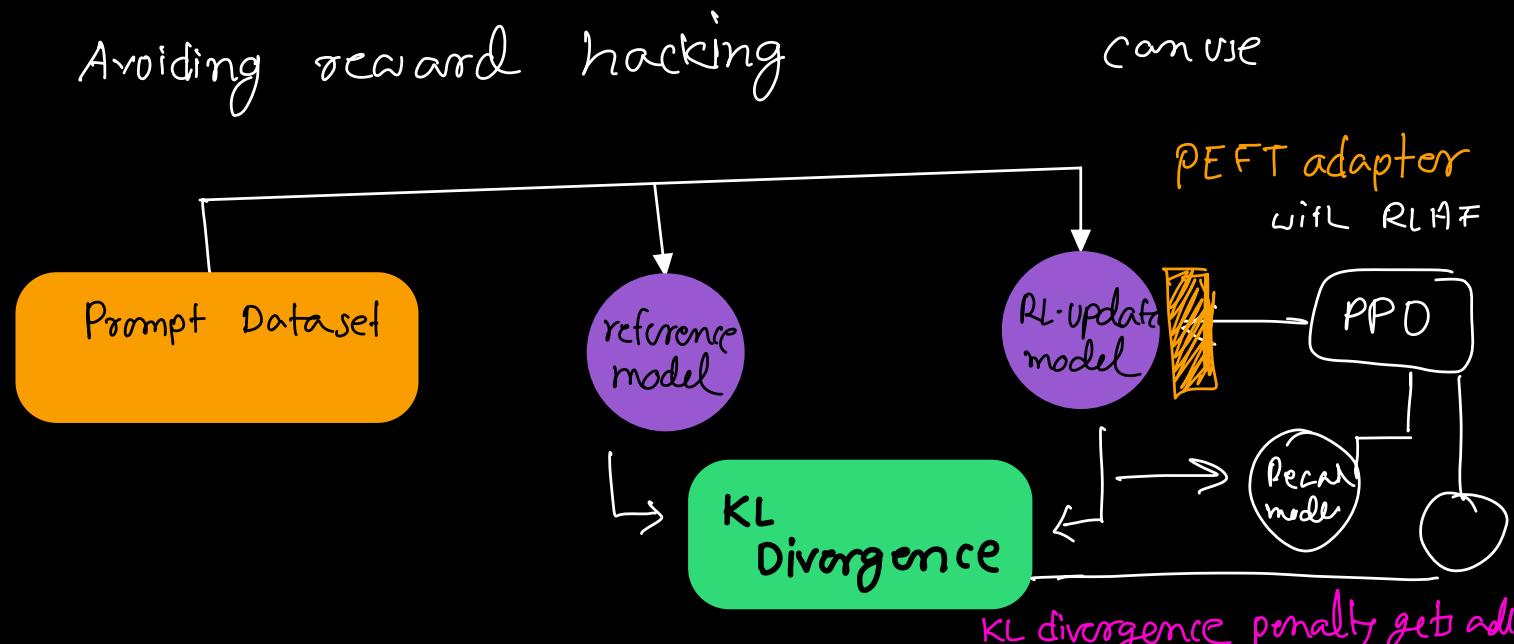
PPO : objective function

Hyperparameters

$$L^{PPO} = L^{POLICY} + C_1 L^{VF} + c_2 L^{ENT}$$

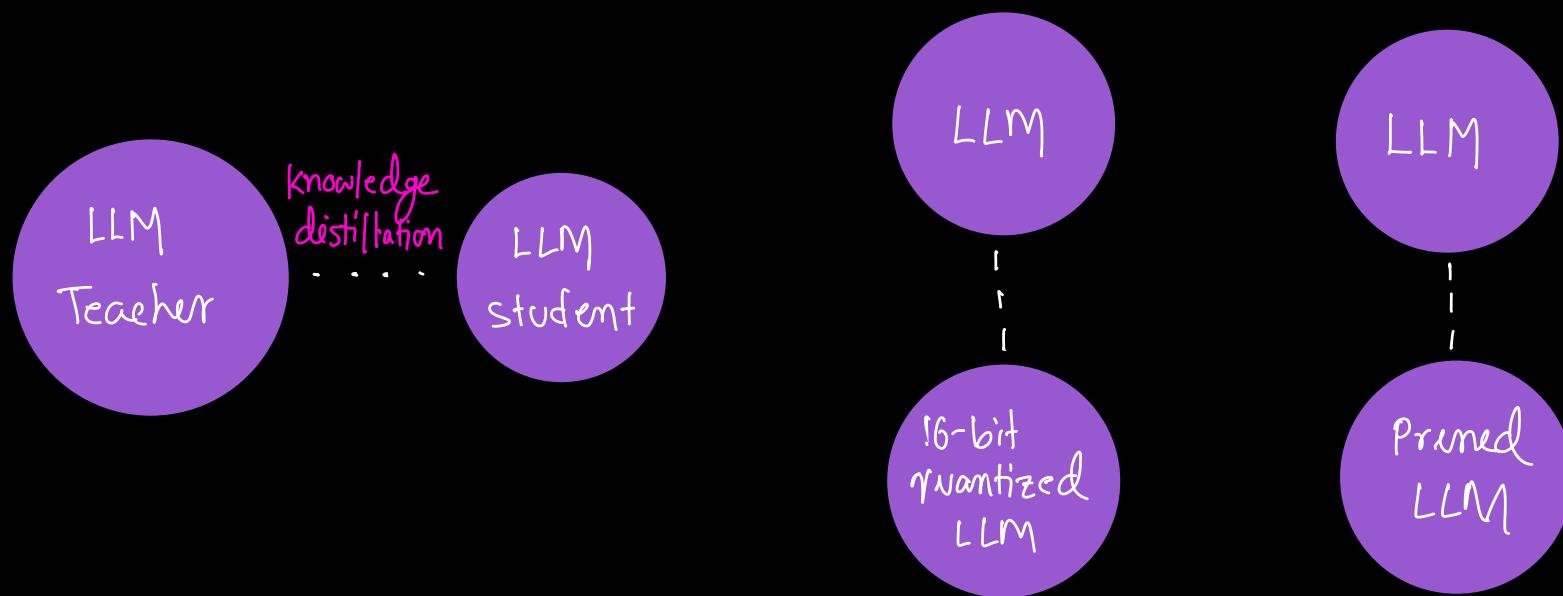
{ policy loss } { value loss } { Entropy loss }

once it completes one iteration PPO , it starts another iteration of PPO and so on.



LLM Powered Applications

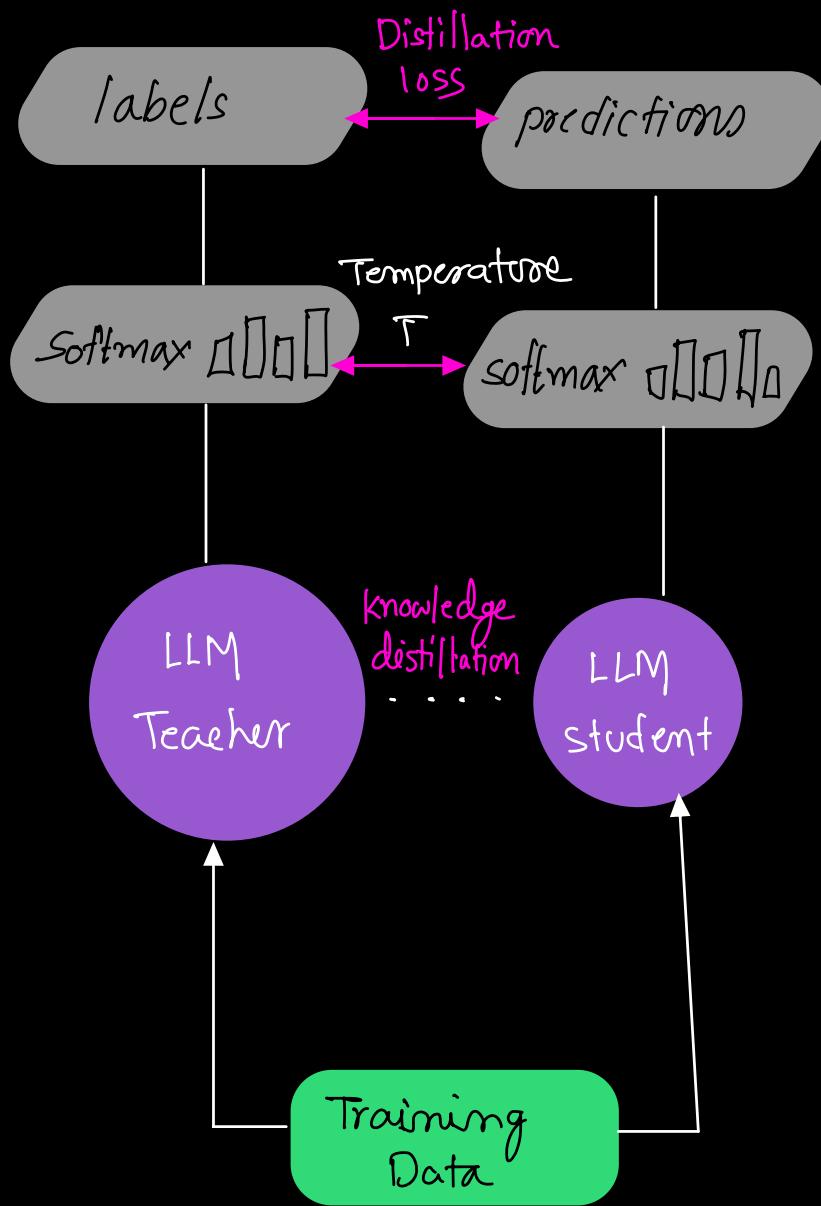
Distillation



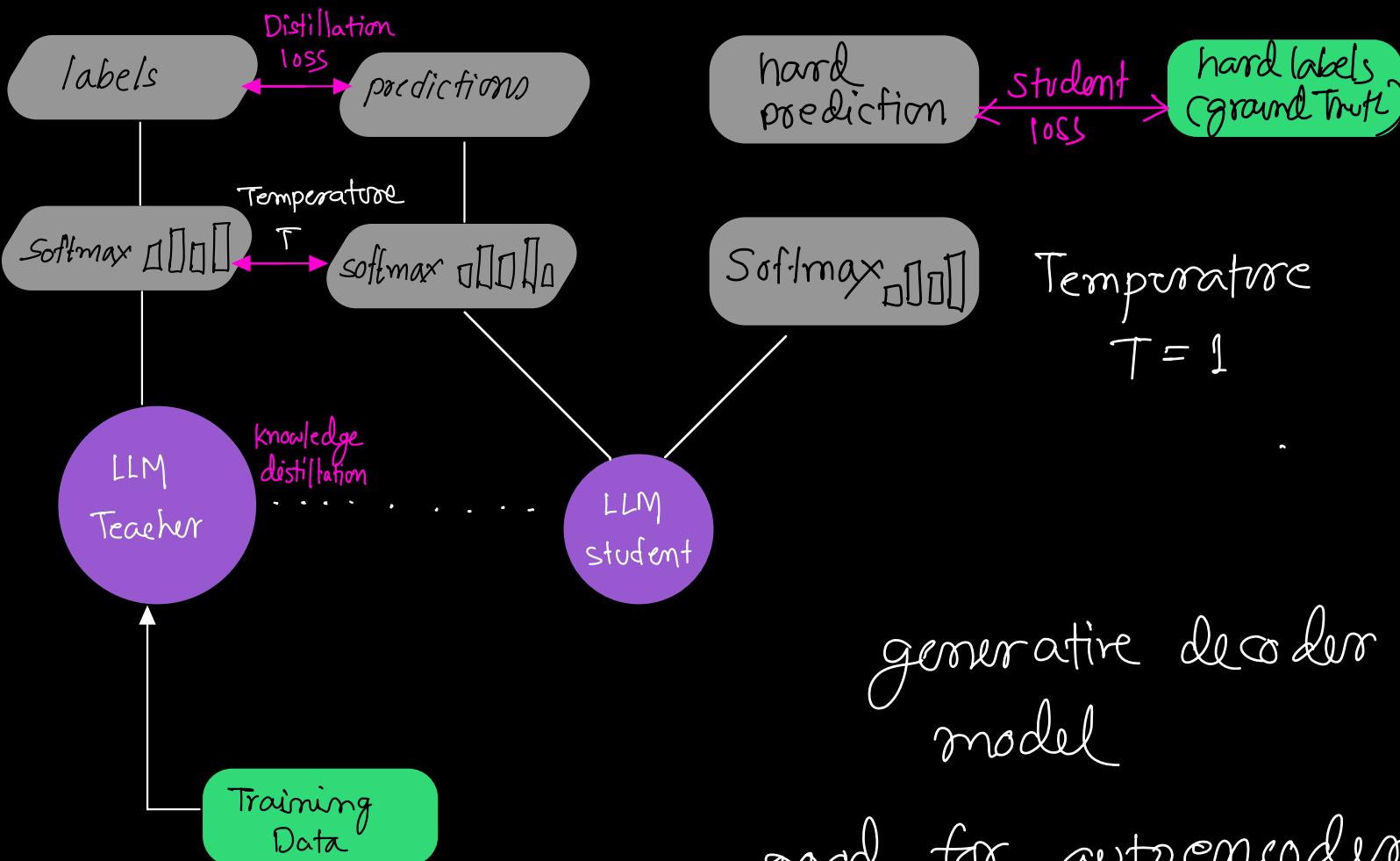
LLM optimization Technique

Quantization

Pruning



Train a smaller
student model
from a larger
teacher model



Temperature
 $T = 1$

generative decoder
model
good for autoencoders
like BERT

Post Training Quantization (PTQ)

Pruning

Remove model weights with values
close to equal to zero

ReLU

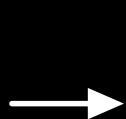
Models having difficulty

prompt

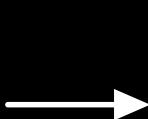
Model

Completion

Who is prime minister of UK ?



LLM



Boris Johnson
out of date

What is $40366/439$?



LLM

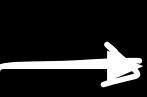


92.49 wrong

What is a Martian Dune? ?



LLM

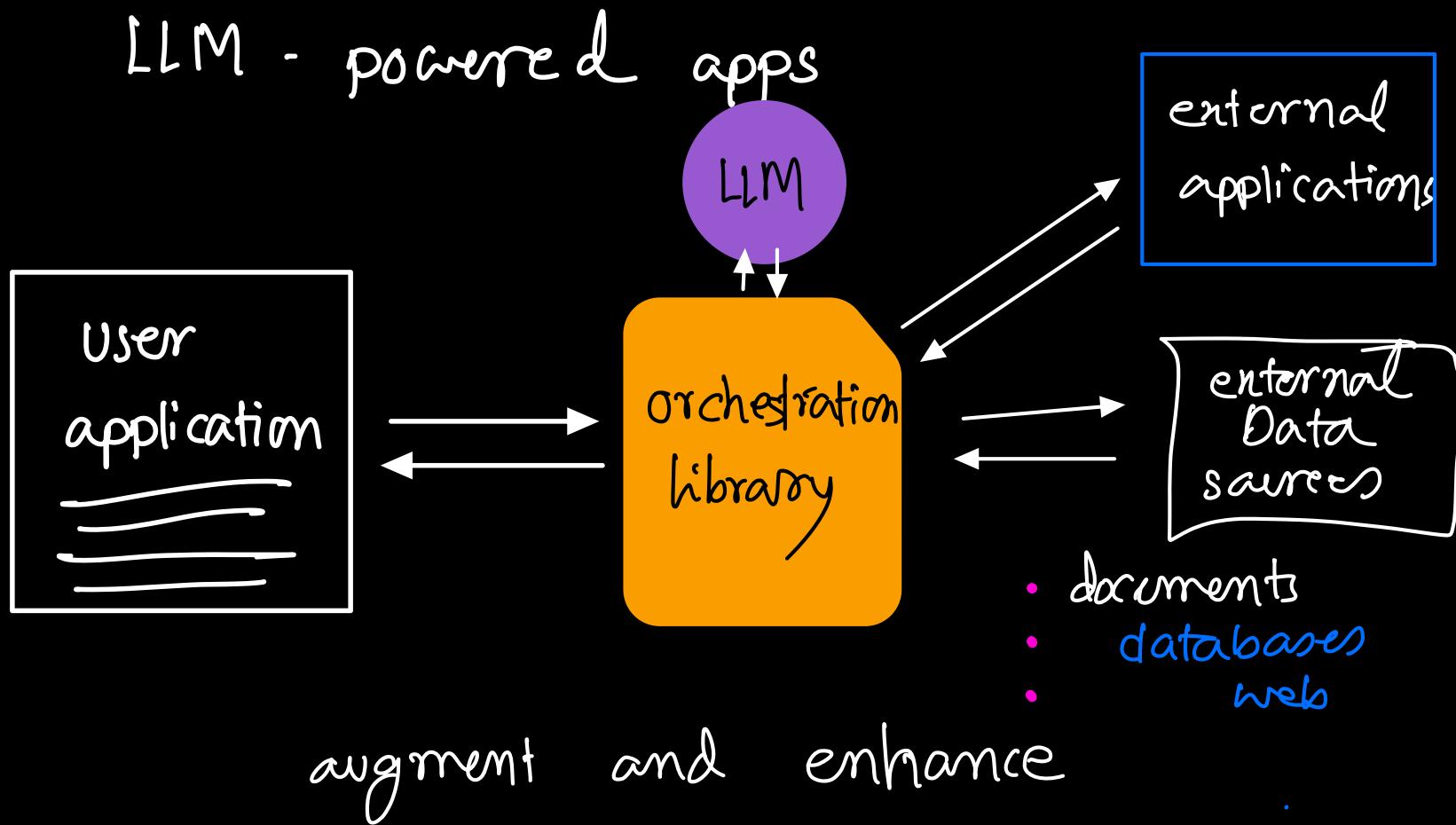


?

hallucination

@nlpwithindrajit

Solution

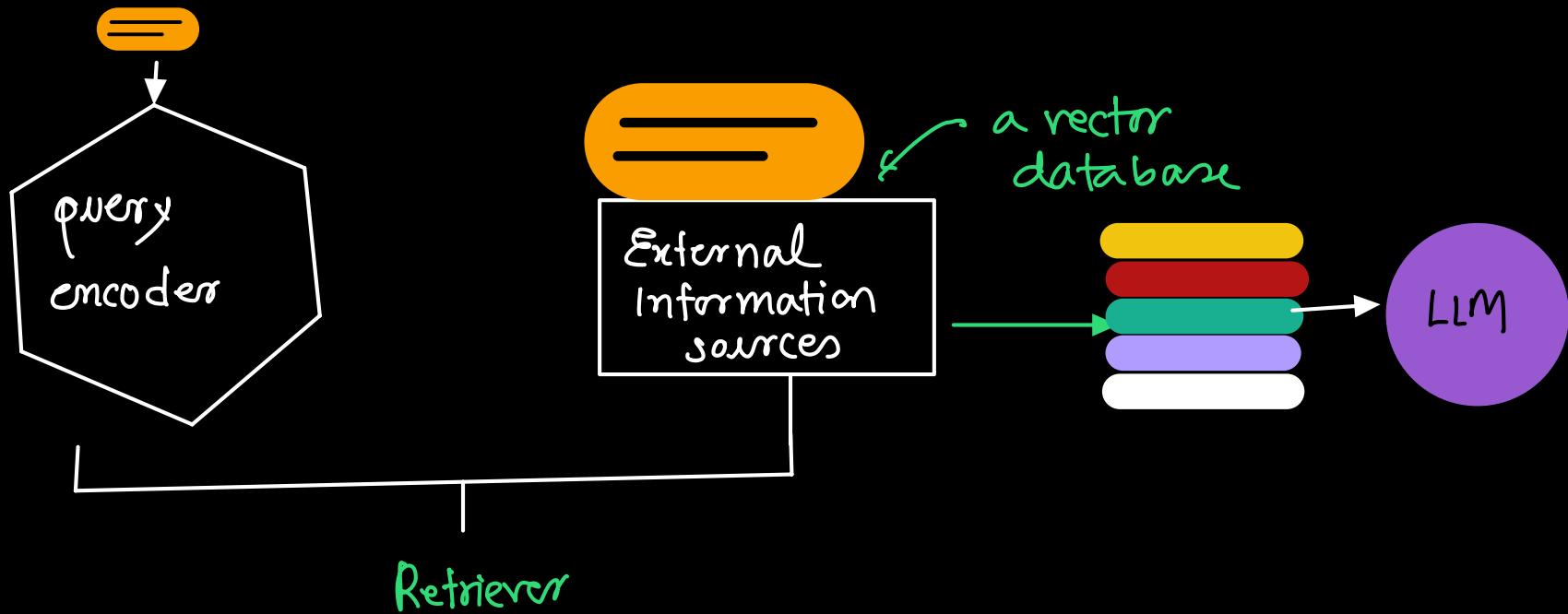


Retrieval augmented generation (RAG)

- a framework to build LLM powered systems that make use of external data systems and applications.
- provide access to data that's not seen during training
- one approach could be to give external / additional data during the inference time.
-

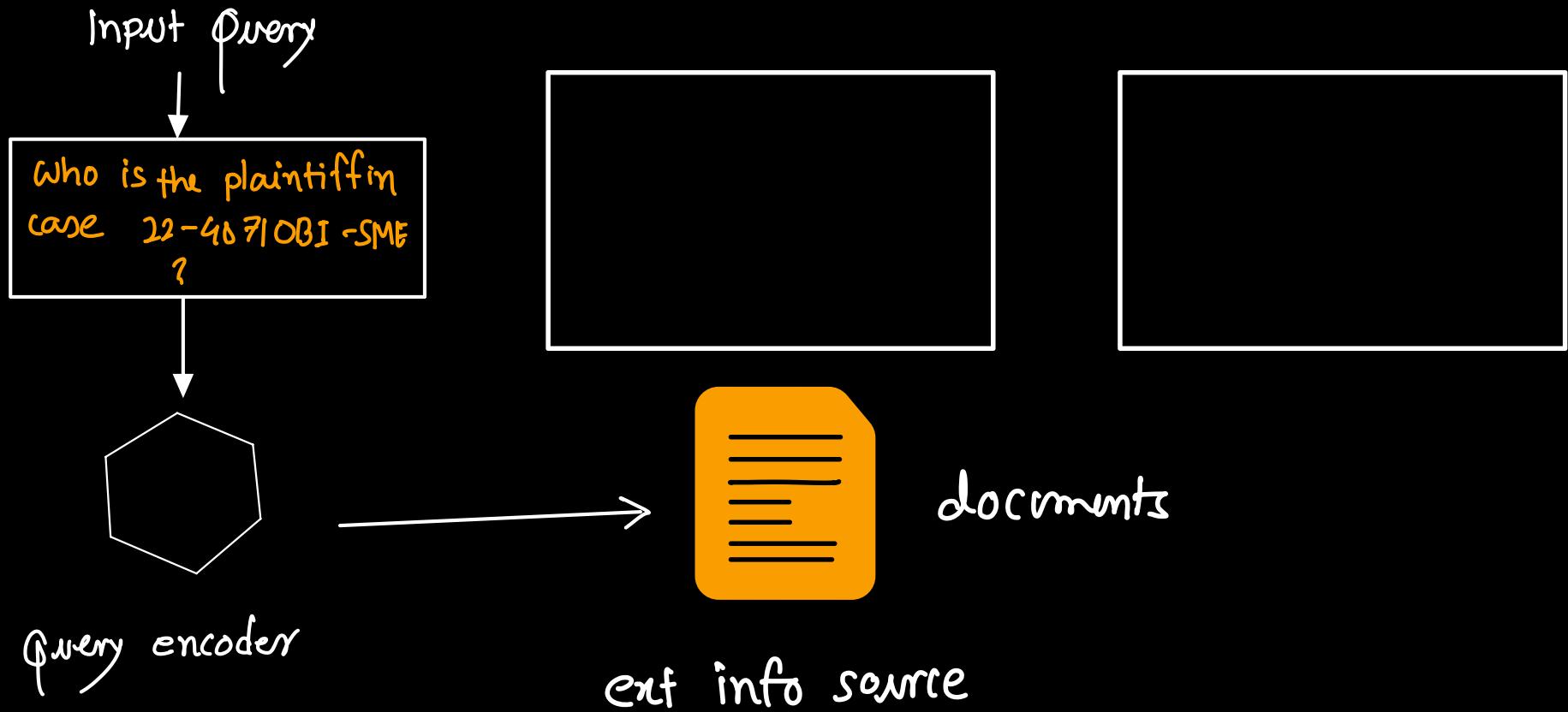
@nlpwithindrajit

RAG



@nlpwithindrajit

Searching Legal documents

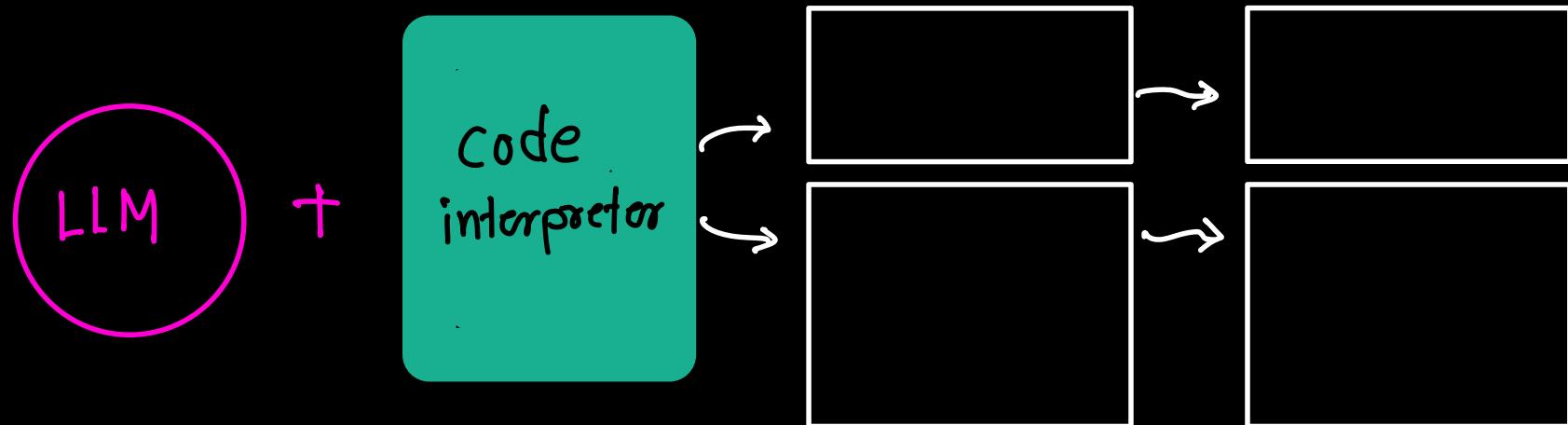


RAG Helps in reducing hallucination problems.

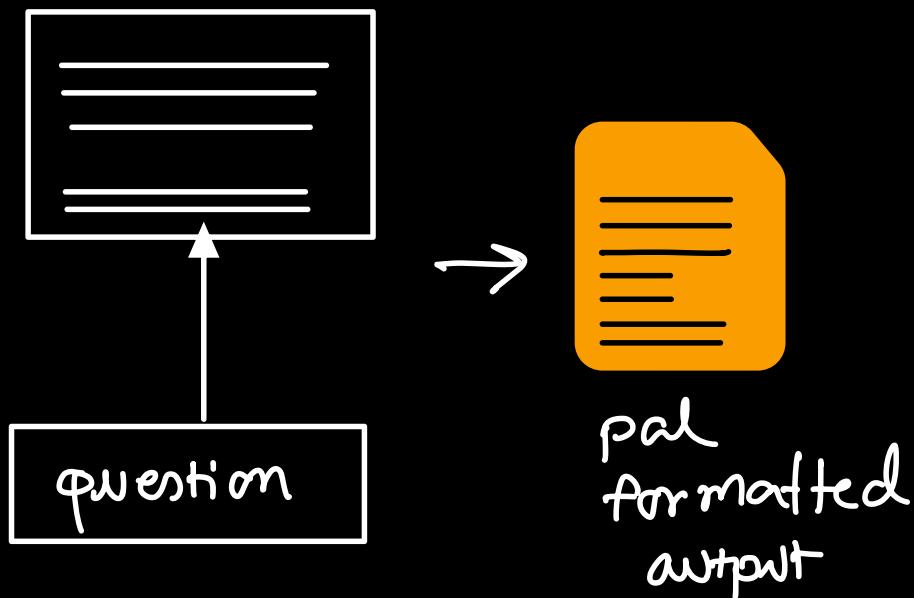
RAG takes small chunks of external data and store them & process them through LLMs to be stored on vector databases.

COT - chain of thought

PAL - program - aided language models



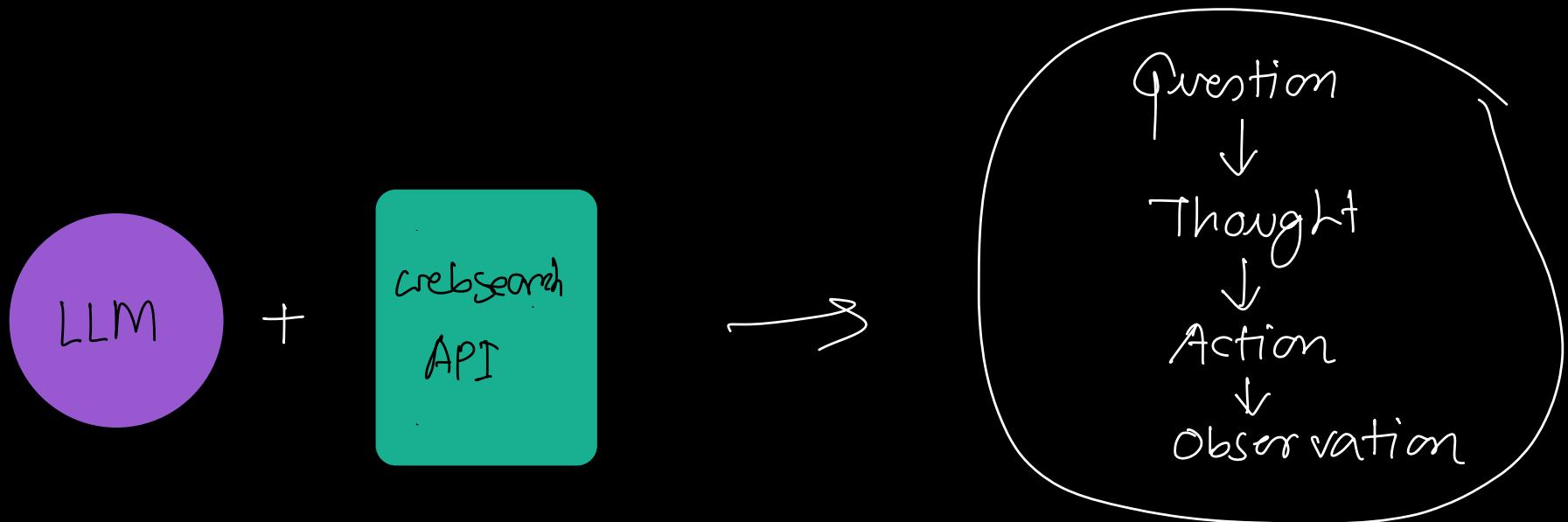
Pal prompt
template



@nlpwithindrajit

ReAct : Combining reasoning and action

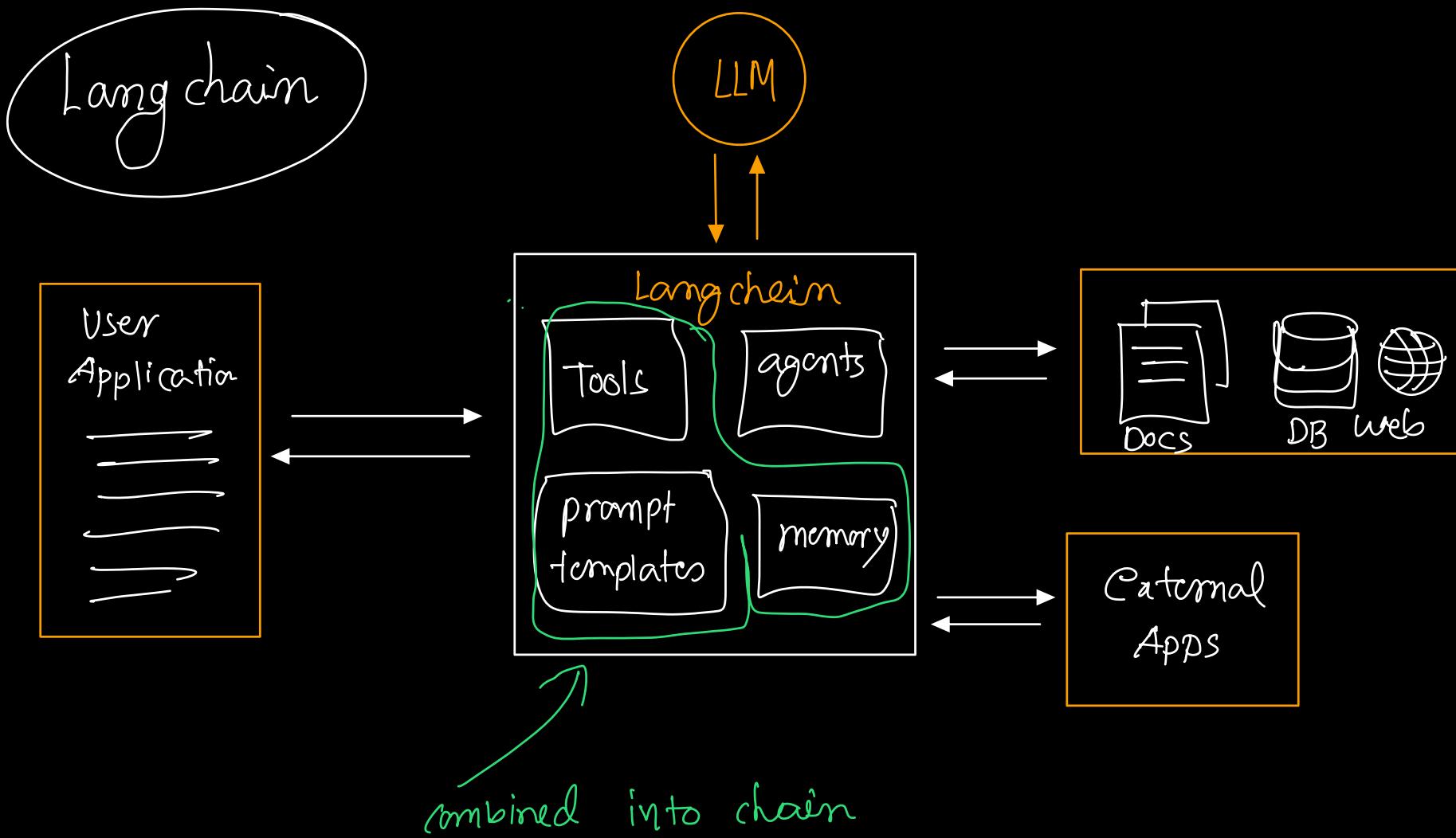
ReAct : Synergizing Reasoning and Action in LLMs



HOpPot QA : multi step question answering

Fever :

@nlpwithindrajit



@nlpwithindrajit

@nlpwithindrajit

Thank You

inspired by Andrew Ng's course

Generative Ai with Large Language
models. (coursera)