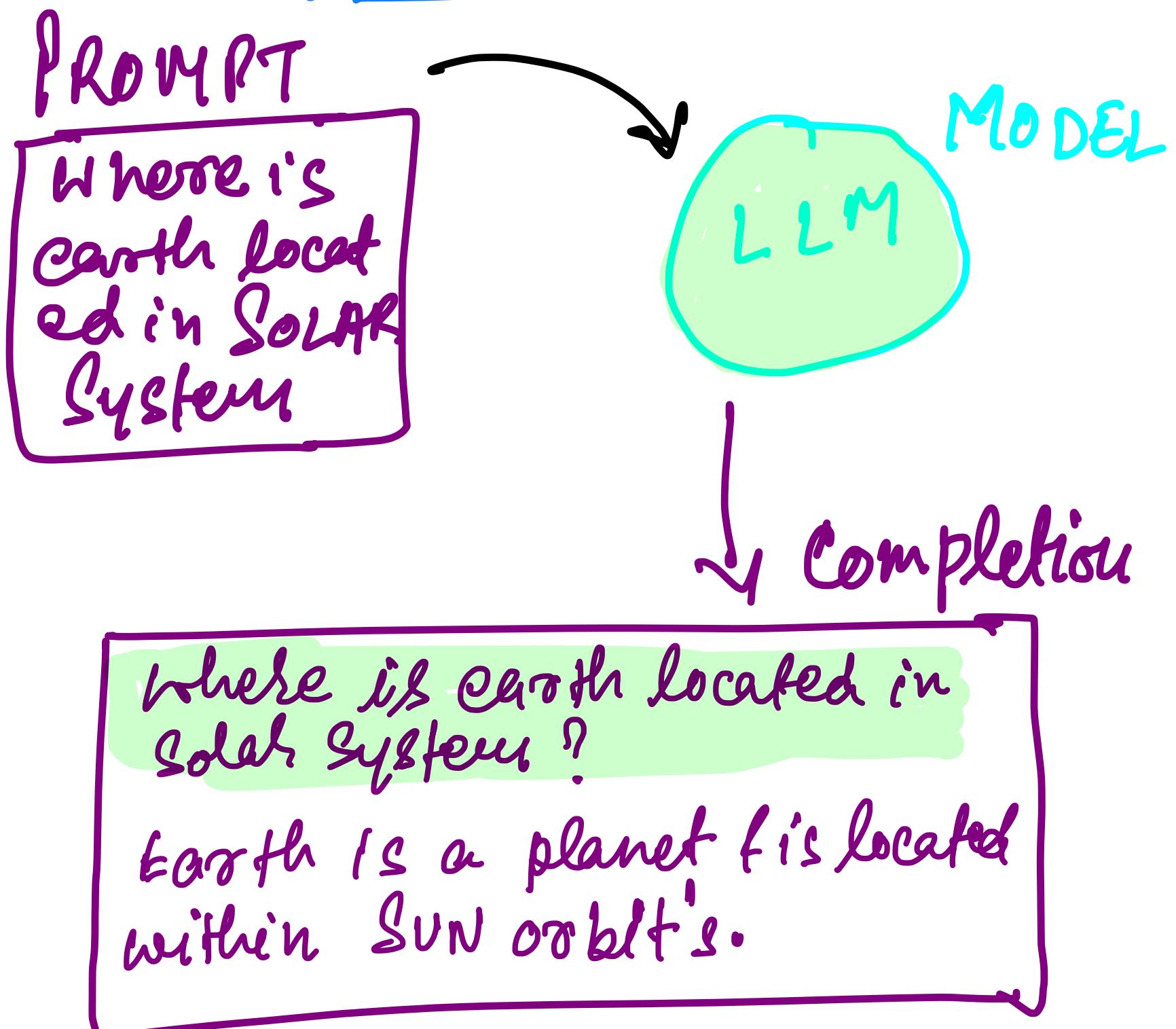


WEEK-01

Generative AI & LLM :-



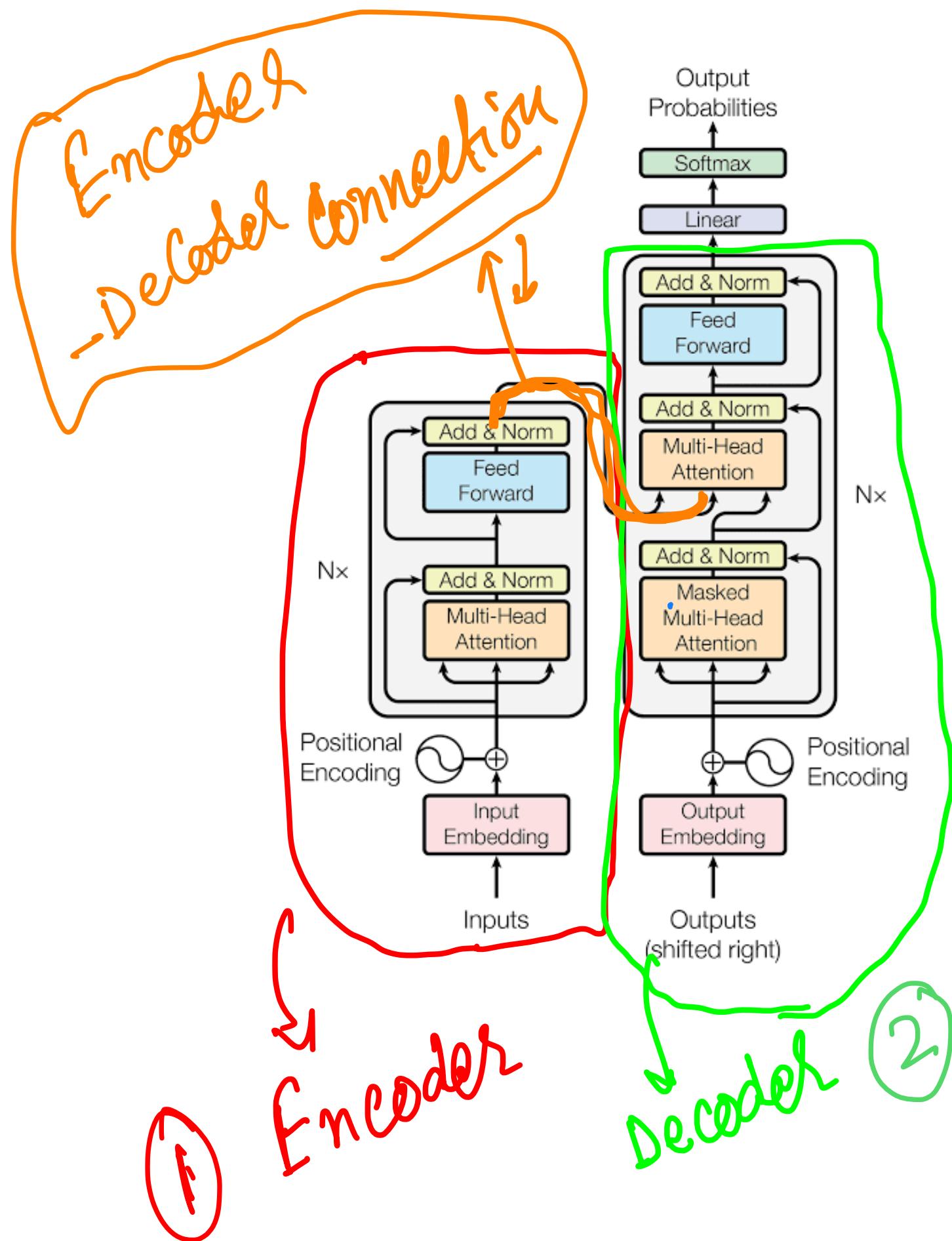
LLM USE CASES :-

1. Summarize 2. Translation

SCALE OF LLM :-

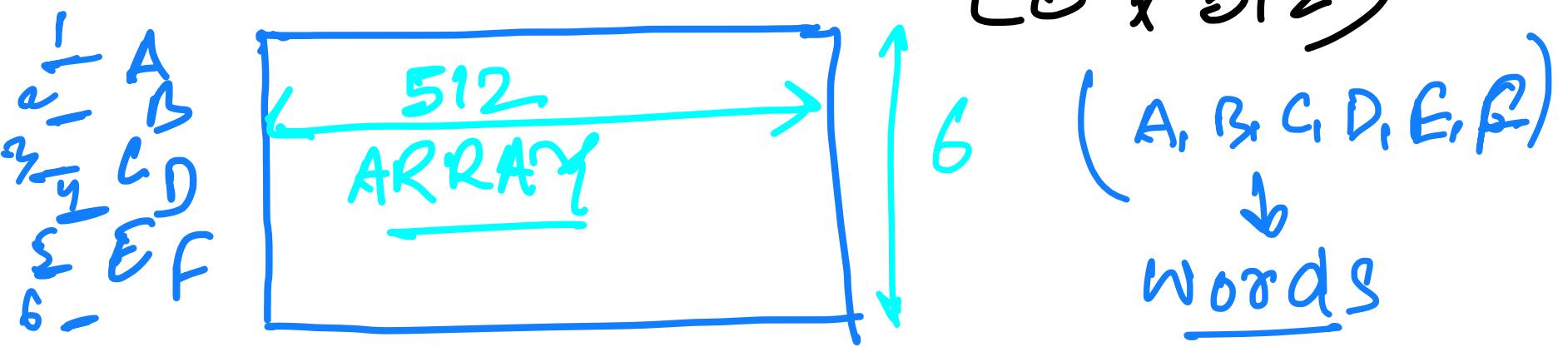


TRANSFORMER ARCHITECTURE



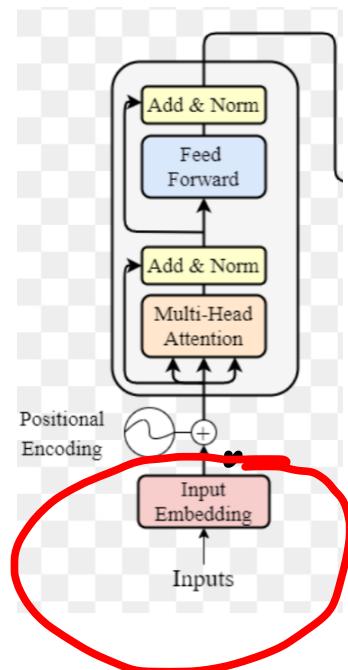
NOTATIONS :-

Input Matrix (sequence, d_{model})
 (6×512)



Encoder:

INPUT EMBEDDINGS ?



Words are mapped to tokens.

You
↓
105 ↴ map

Cat
↓
6587

is
↓
5475

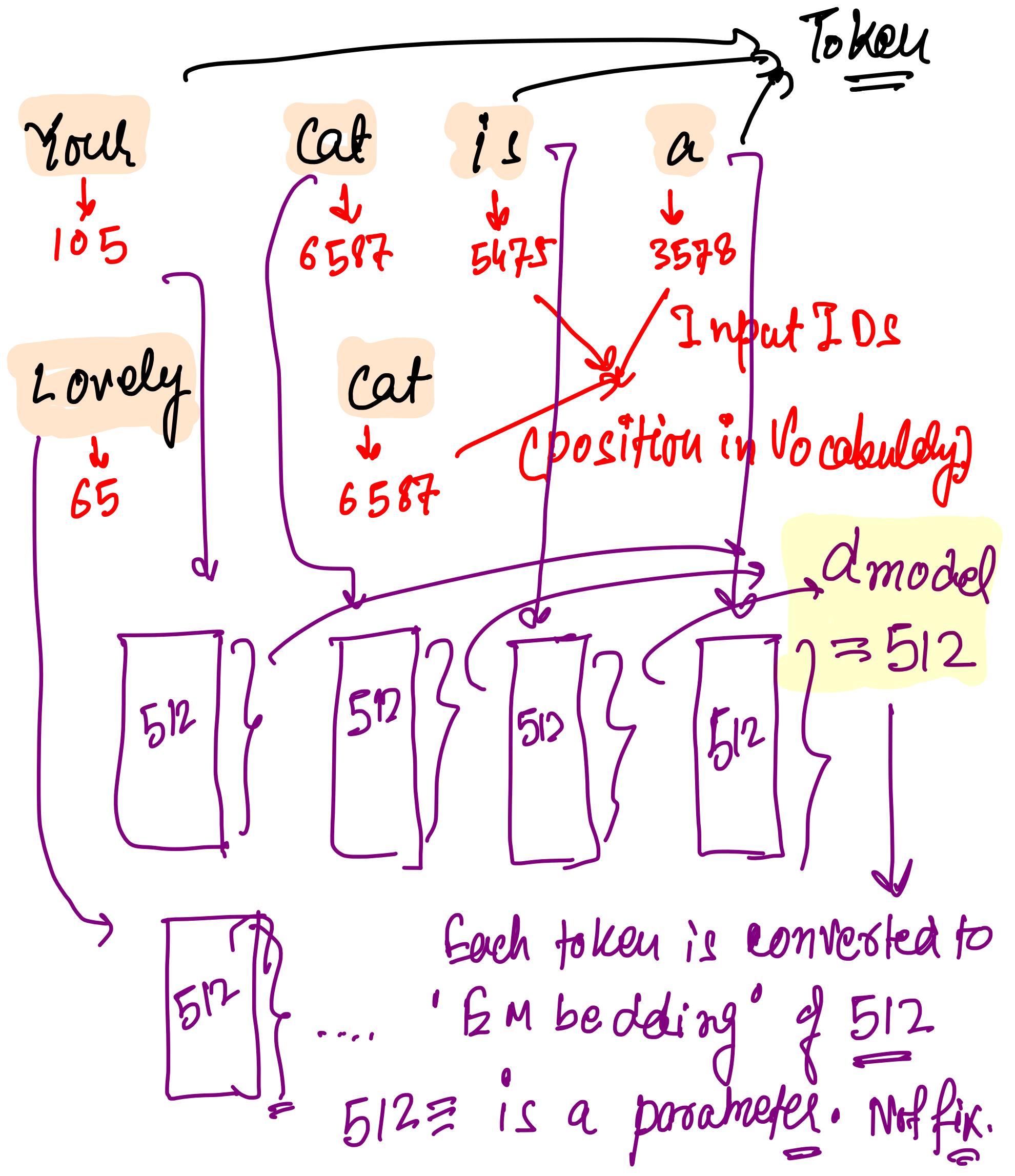
a
↓
3578 ↴ map

Lovely
↓
65 ↴ map

Cat
↓
6587

↳ to position in Vocab

Input IDs
(position in Vocabulary)



Positional Encoding:

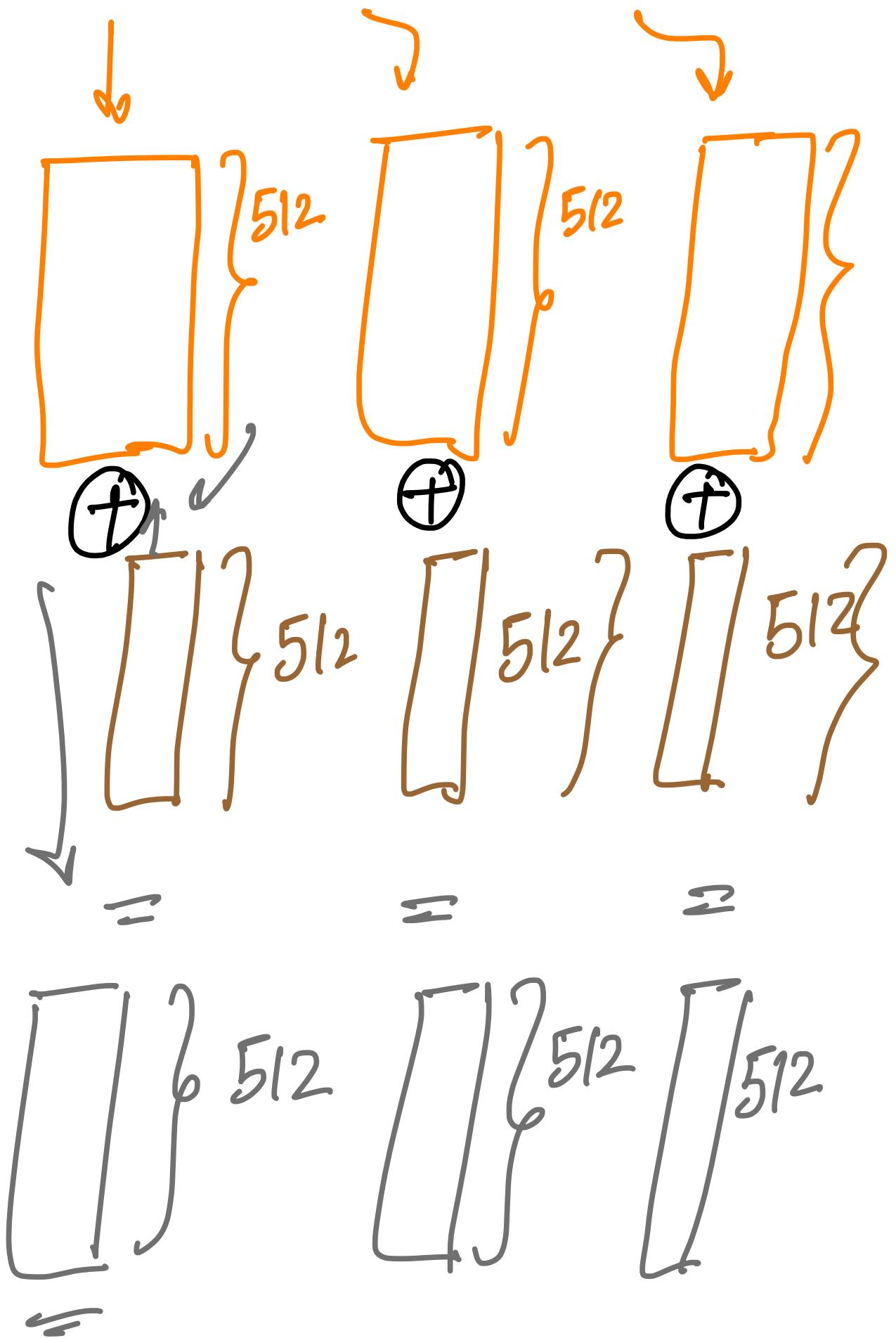
- ④ We want each word to carry some information about its position in the sentence.
- ④ We want model to treat words that appear close to each other as 'CLOSE'. & words that are distant as 'DISTANT'.

Original Seq = Your Cat is ...

Embeddings
Vector of 512 =

Positional Embedding, Vector
512

Encoder
Input, Vector
512



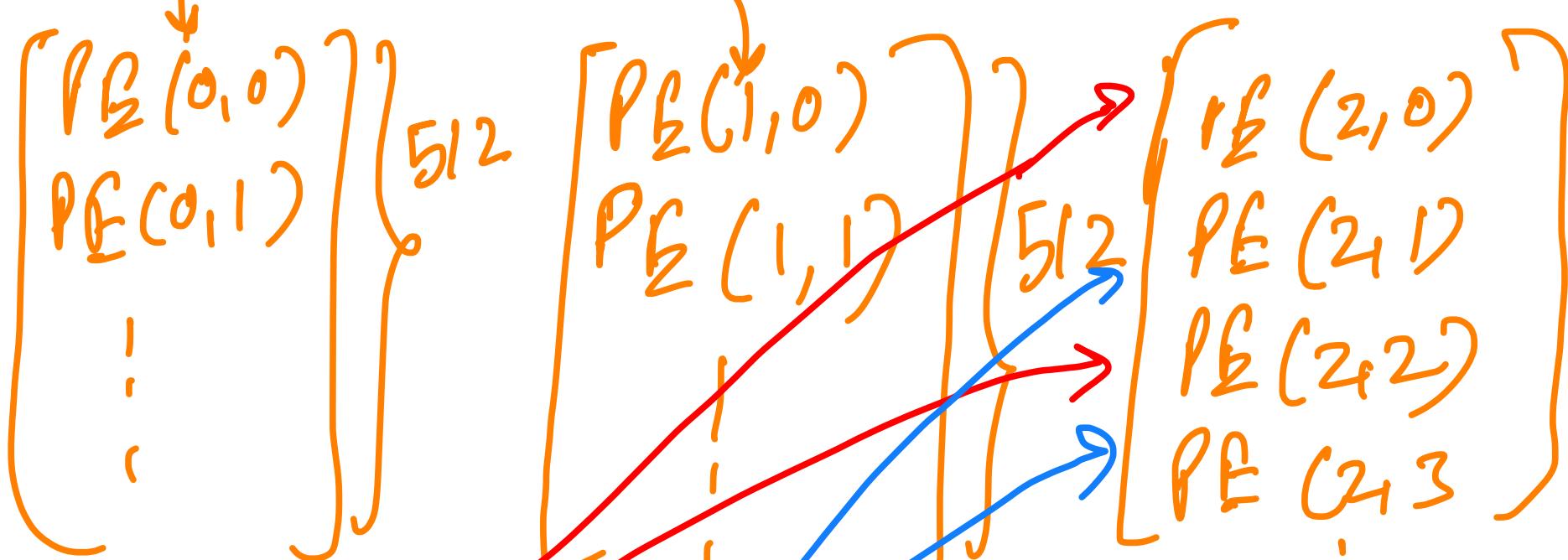
How Positional Encoding (PE) Calculated?

Your Cat is

pos=0

pos=1

pos=2



$$PE_2(pos, di) = \sin\left(\frac{pos}{10000^{di/d_{\text{model}}}}\right)$$

even rows

$$PE_2(pos, di+1) = \cos\left(\frac{pos}{10000^{di/d_{\text{model}}}}\right)$$

odd rows

Encoder: Multihead Attention

Let's first understand 'Single Head Attention'.

Self-Attention:- Self-attention allows the model to relate words to each other.

$$\text{Attention}(Q, K, V) = \text{softmax} \left(\frac{QK^T}{\sqrt{dk}} \right)V$$

Let's take simple example 6
word size = 512.



$$seq = 6$$

$$d_{\text{model}} = 512$$

This can be re-presented by matrix Q, K, V , these are

same matrix re-presenting input of 'SEX' words with dimension of '512' vector length. ($d_{model} = 512$)

$$Q = K = V =$$

Matrix (5×512)

Softmax

$$\underbrace{Q \times K^T}_{\sqrt{d_{model}} = \sqrt{512}} = \text{matrix } =$$

$Q \quad (6, 512)$

$K^T \quad (512, 6)$

$\sqrt{d_{model}} = \sqrt{512}$

	Your	Cat	is	a	lonely	Cat	Σ
Your	.263	.119	-	-	-	-	1
Cat	-	-	-	-	-	-	1
is	-	-	-	-	-	-	1
a	-	-	-	-	-	-	1
lonely	-	-	-	-	-	-	1
Cat	-	-	-	-	-	-	1

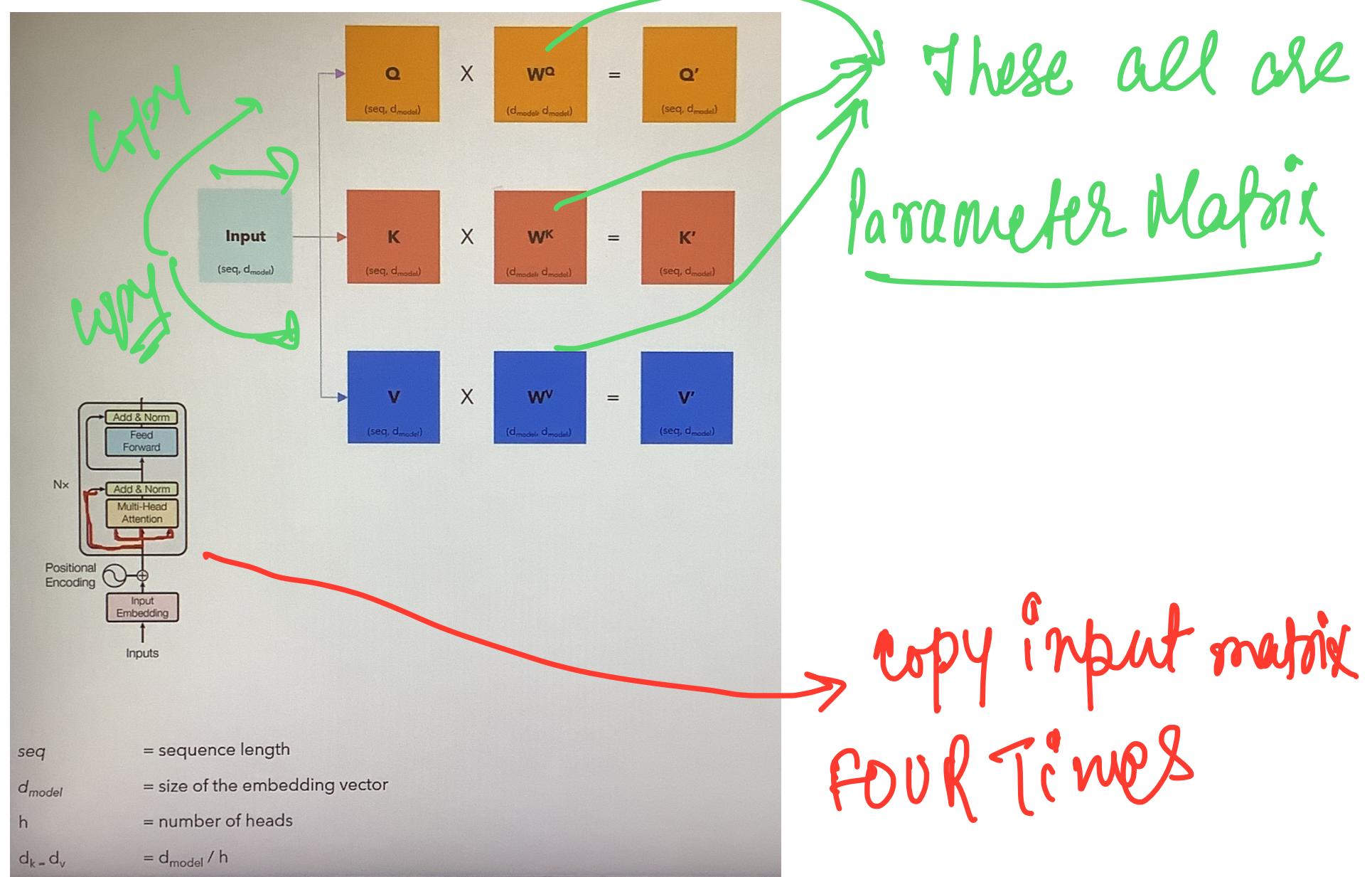
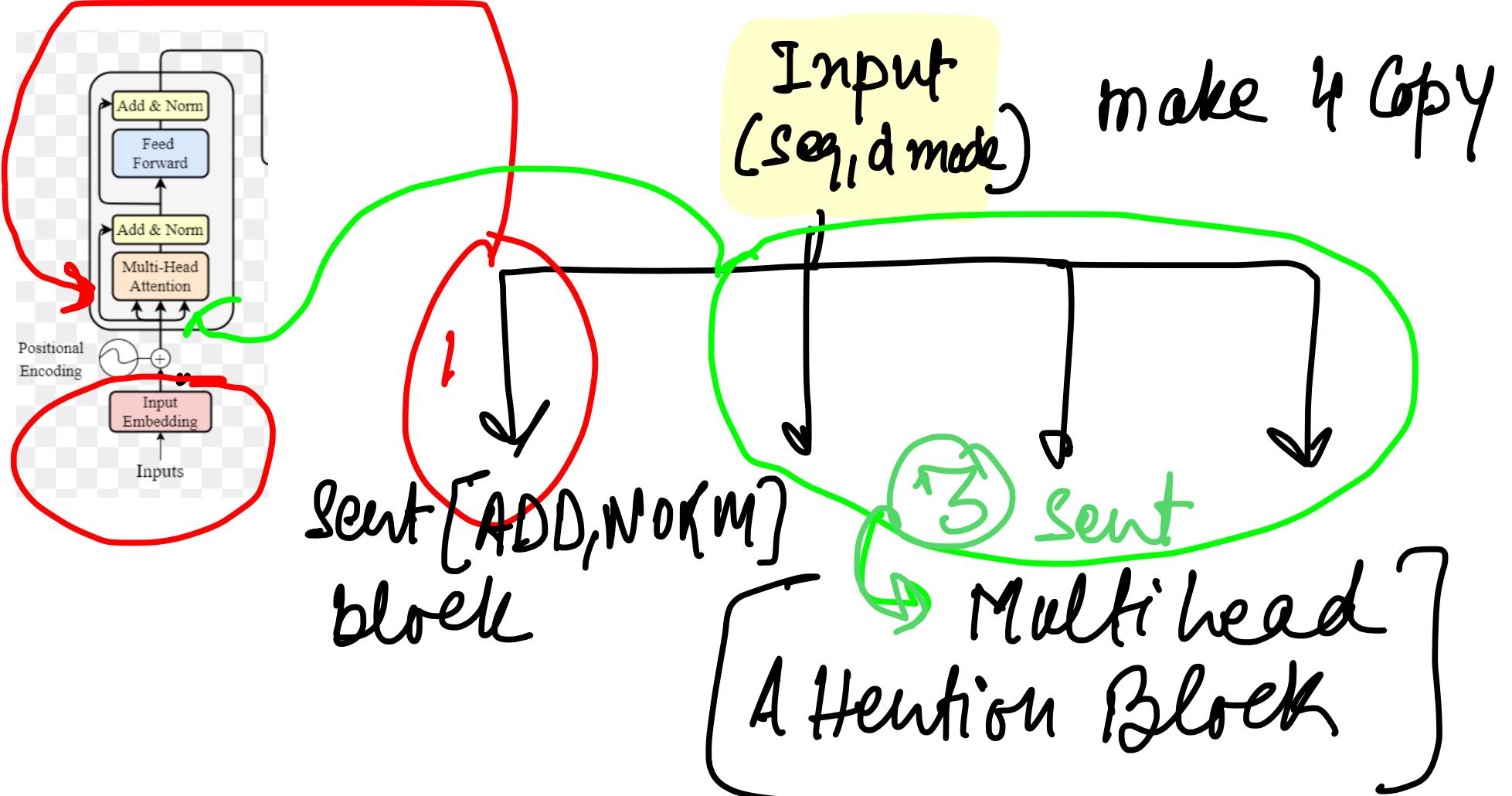
This is intermediate matrix of Attention which. Each cell is DOT Product of Embeddings. Hence each cell represent a 'Score' How intense is a relationship bet ween those words.

(6 × 6)

× ^V
(6 × 512) =

Attention
(6 × 512)

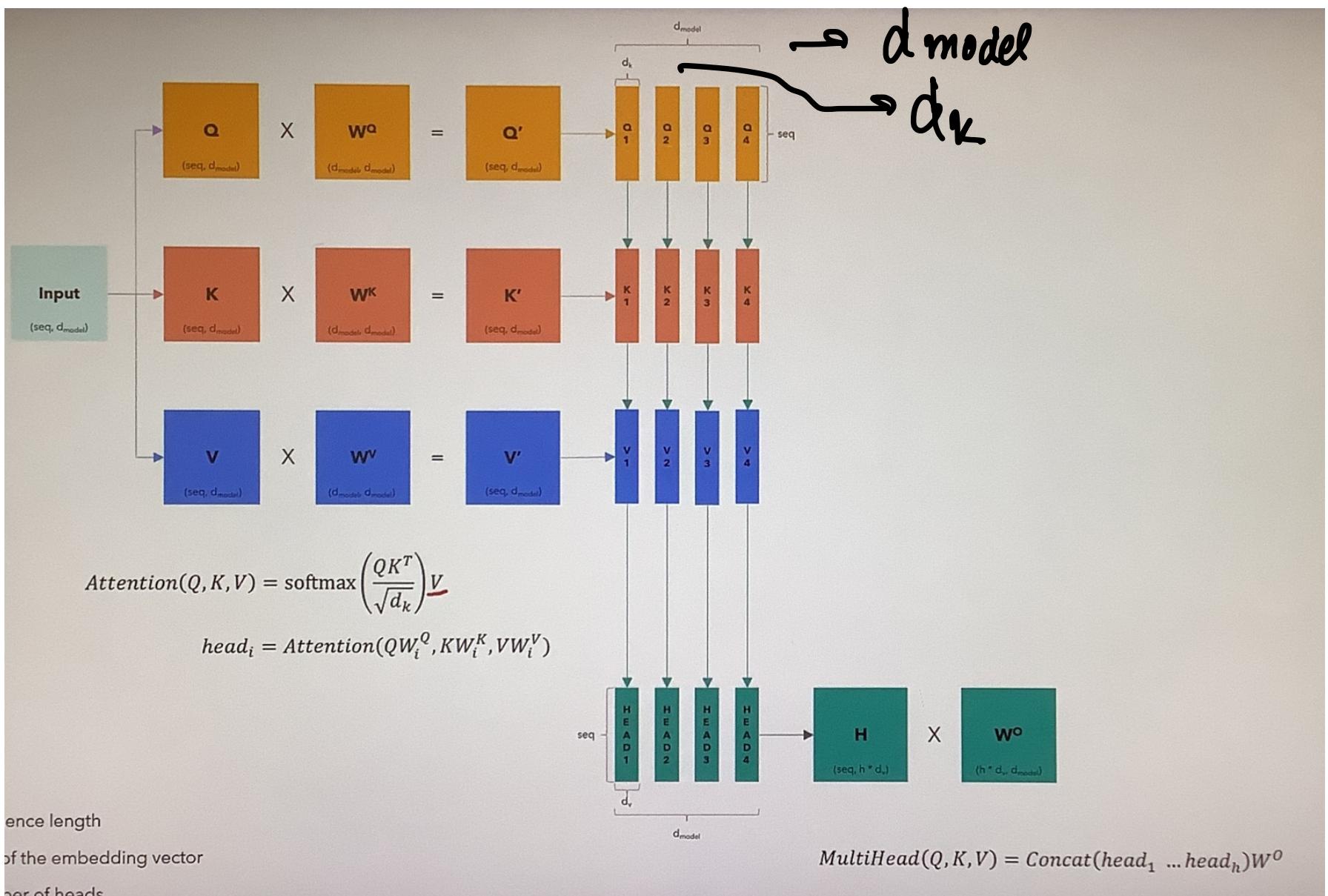
- ④ Each row of Attention matrix captures not only the meaning (given by embedding) or the position in sentence (represented by the positional encodings) but also each words INTERACTIONS with other words.



W^Q, W^K, W^V = Parameter Matrix.
 h ($h = \text{no of heads}$)

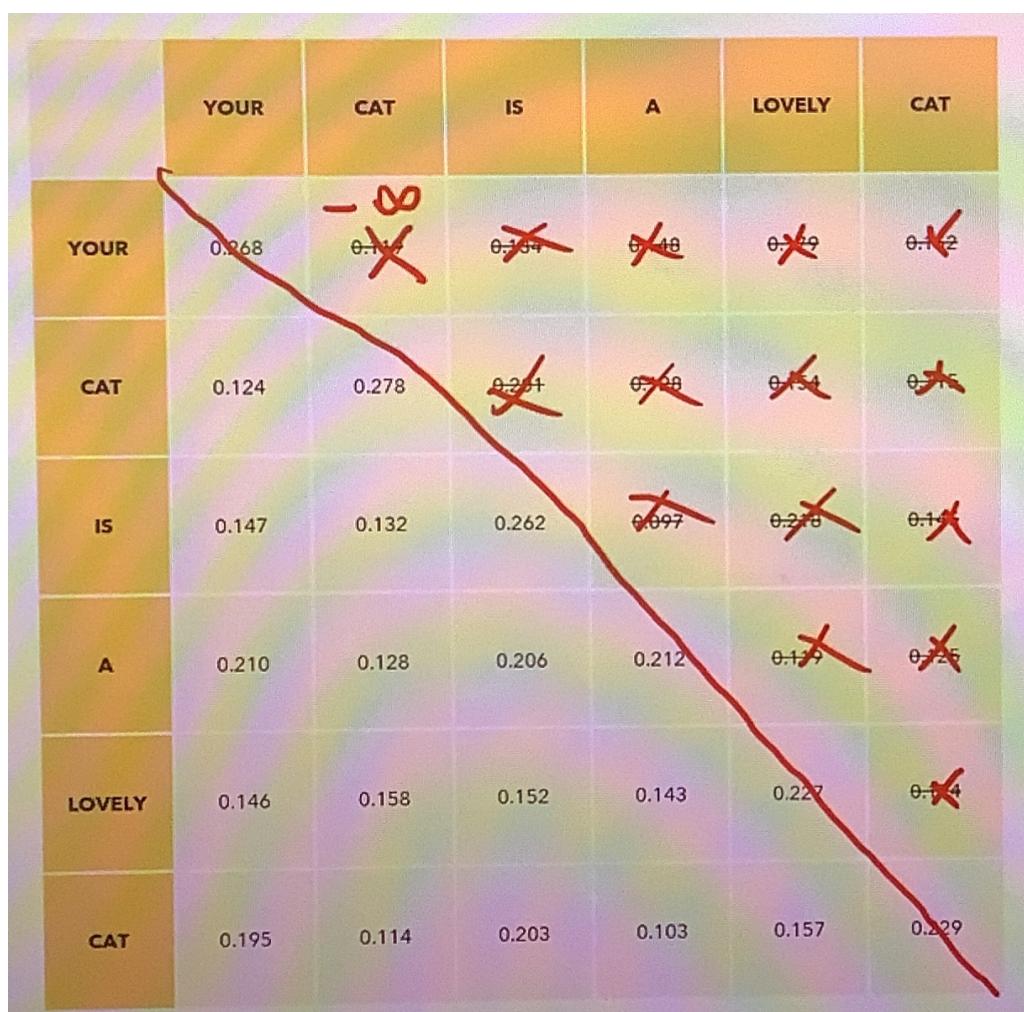
④ Each head watching full sentence but part of embedding.

$(d_K = 512/4)$
 $d_{\text{model}} = \text{embedding size.}$



Masked Multi-Head Attention :-

Our goal is to make model causal
! It means output at certain position can only depend on the words on the previous positions.
The model must not be able to see the Future Words.



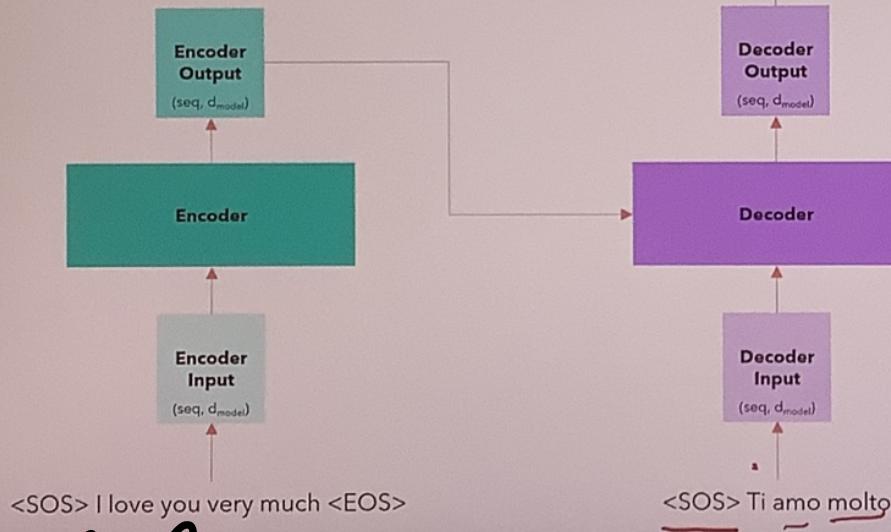
We achieve it by masking = $-\infty$ the next word & then lot of max turn that into ZERO.

Training

Time Step = 1

It all happens in one time step!

The encoder outputs, for each word a vector that not only captures its meaning (the embedding) or the position, but also its interaction with other words by means of the multi-head attention.

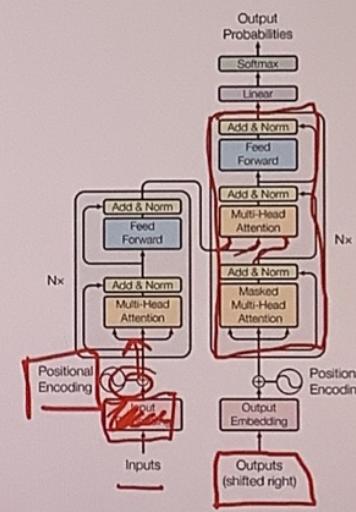


Label

Ti amo molto <EOS>

* This is called the "label" or the "target"

Cross Entropy Loss



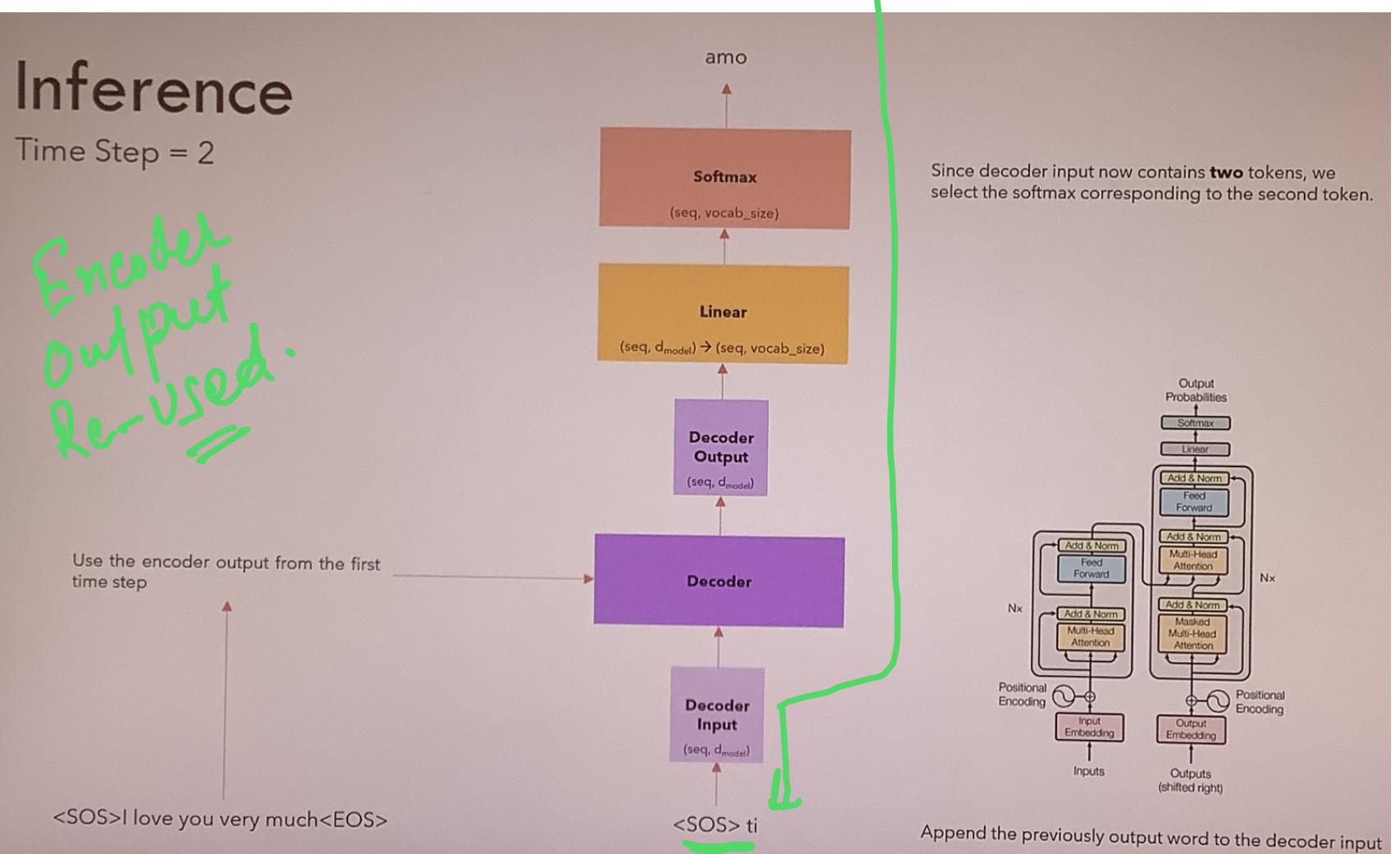
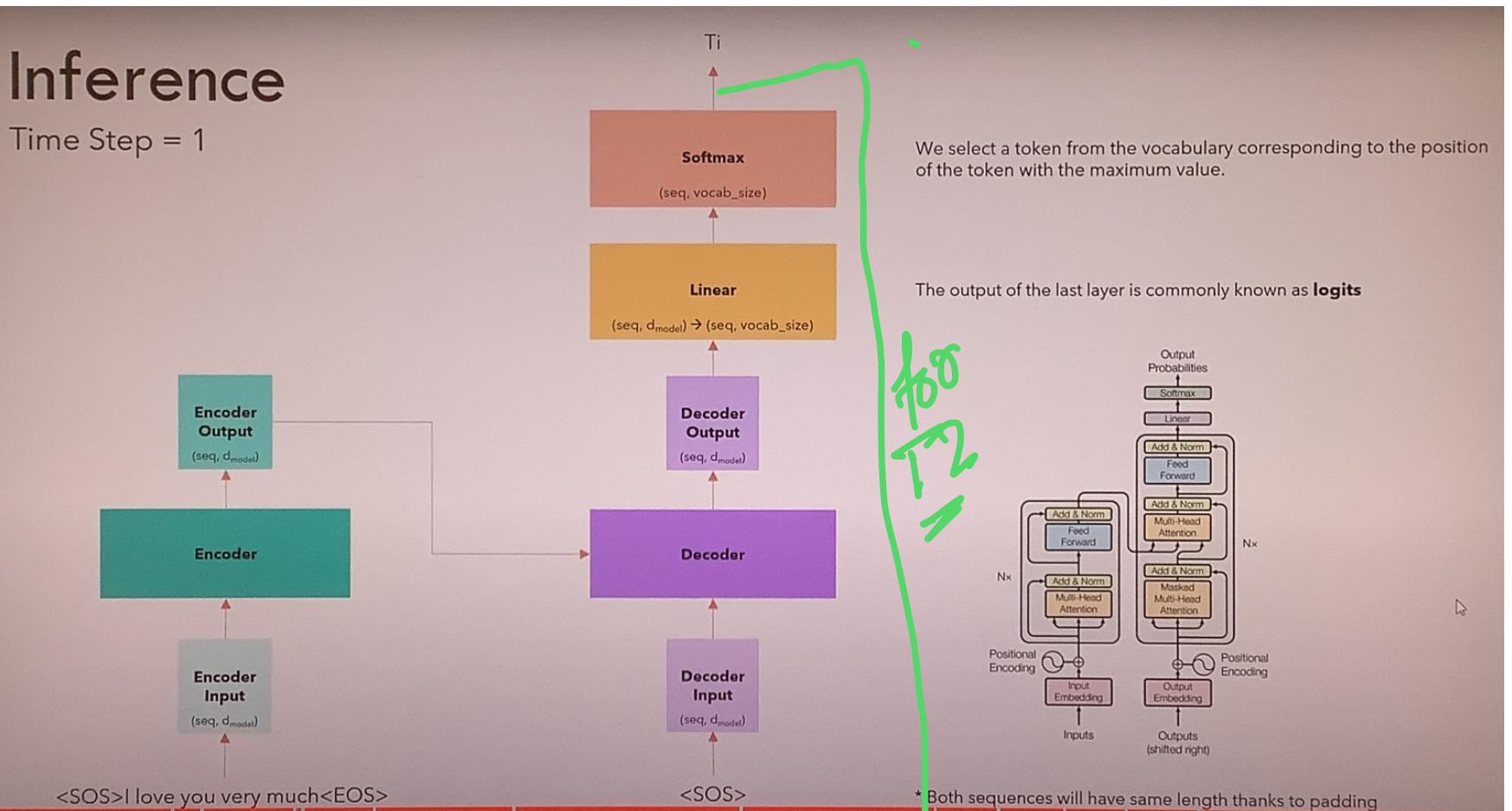
We prepend the <SOS> token at the beginning. That's why the paper says that the decoder input is shifted right.

Umar Jamil - <https://github.com/hkproj/transformer-from-scratch-notes>

↑ Input

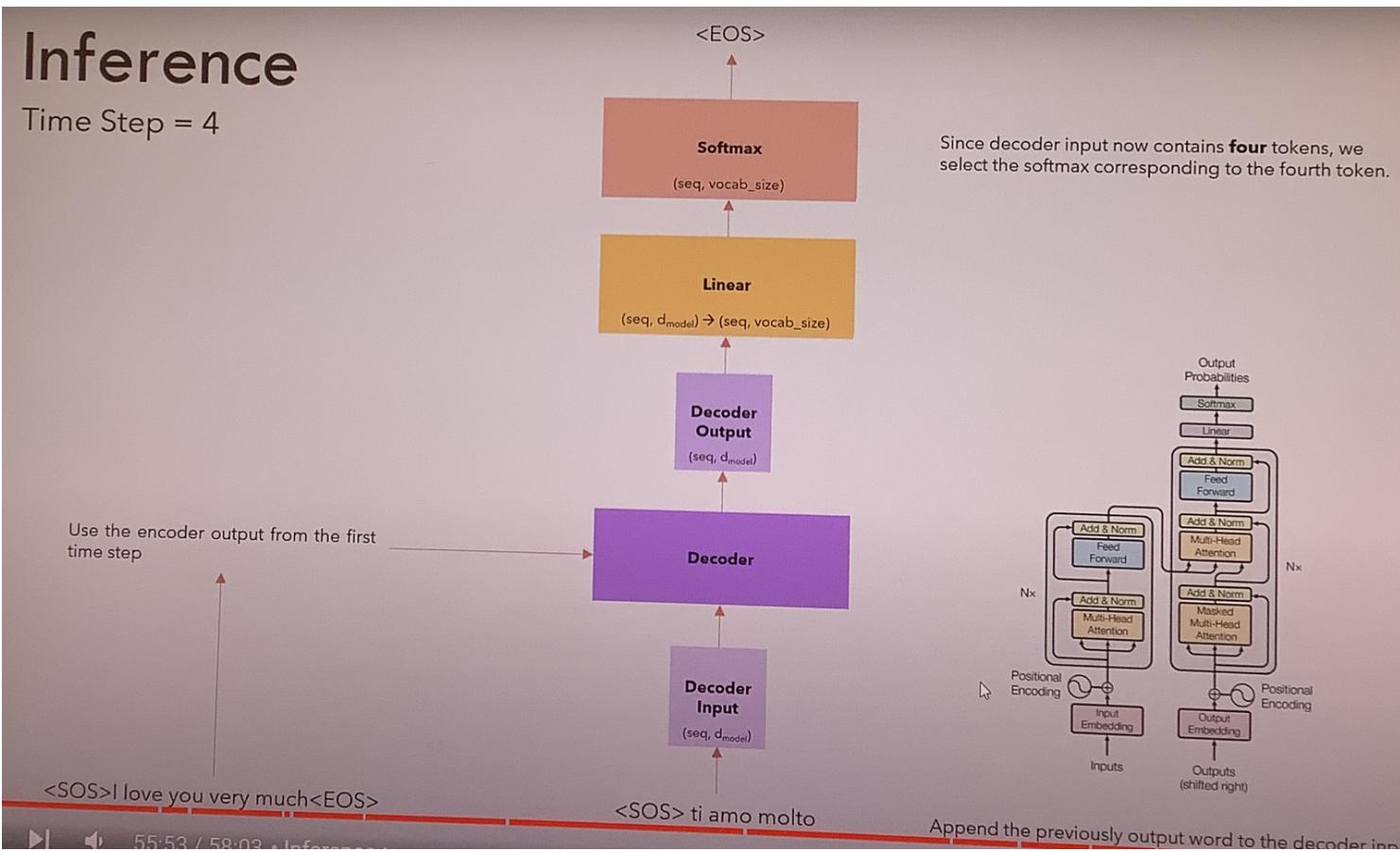
④ One long sequence can be trained parallelly faster.

Inference :- Inference is done Token-by-Token.



Inference

Time Step = 4



Top
B

We keep on
iterating

until

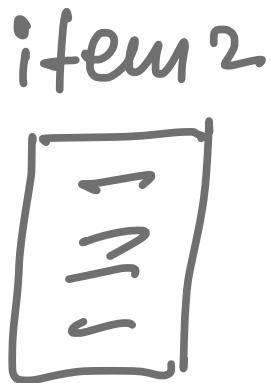
decoder part produce `<EOS>` End If
`Sentence>` token. Here we have
used GREEDY Strategy. Beam search
is better at each step we select TopB
words & evaluate all the possible
next words for each of them at
each step keeping top B most probable
sequences. This is the "BEAM
SEARCH"

Add f Norm :- At this stage we perform 'Layer Normalization' not 'Batch Normalization'

Ex: A batch with 3 items having 512 vector size.



$$\mu_1 \\ \sigma_1^2$$



$$\mu_2 \\ \sigma_2^2$$



$$\mu_3 \\ \sigma_3^2$$

$$\hat{x}_j = \frac{x_j - \mu_j}{\sqrt{\sigma_j^2 + \epsilon}} * \gamma + \beta$$

↳ learning Params.

Prompts & Completions :-

PROMPT

Where is Ganymede located in the Solar System?

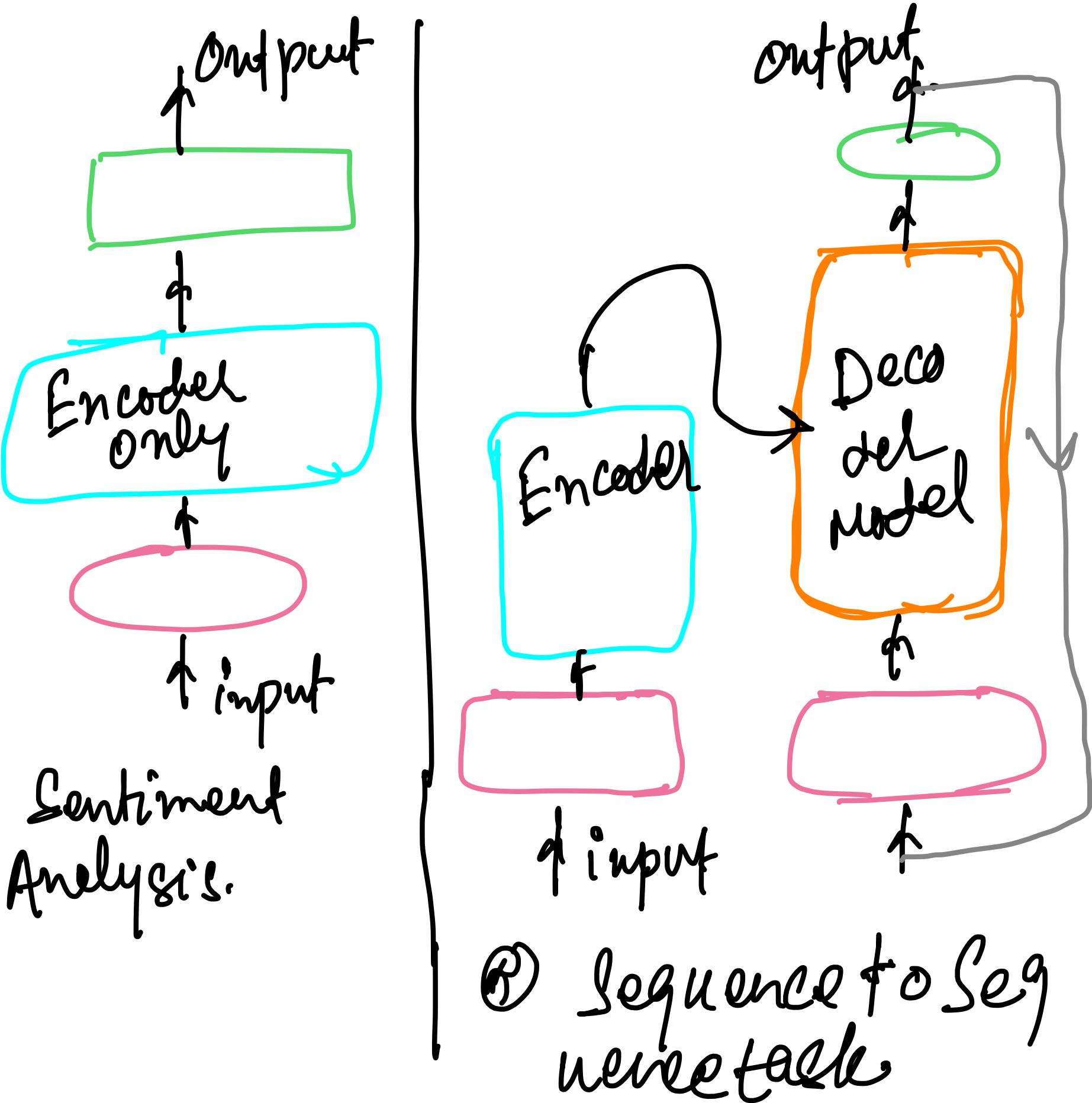
Model

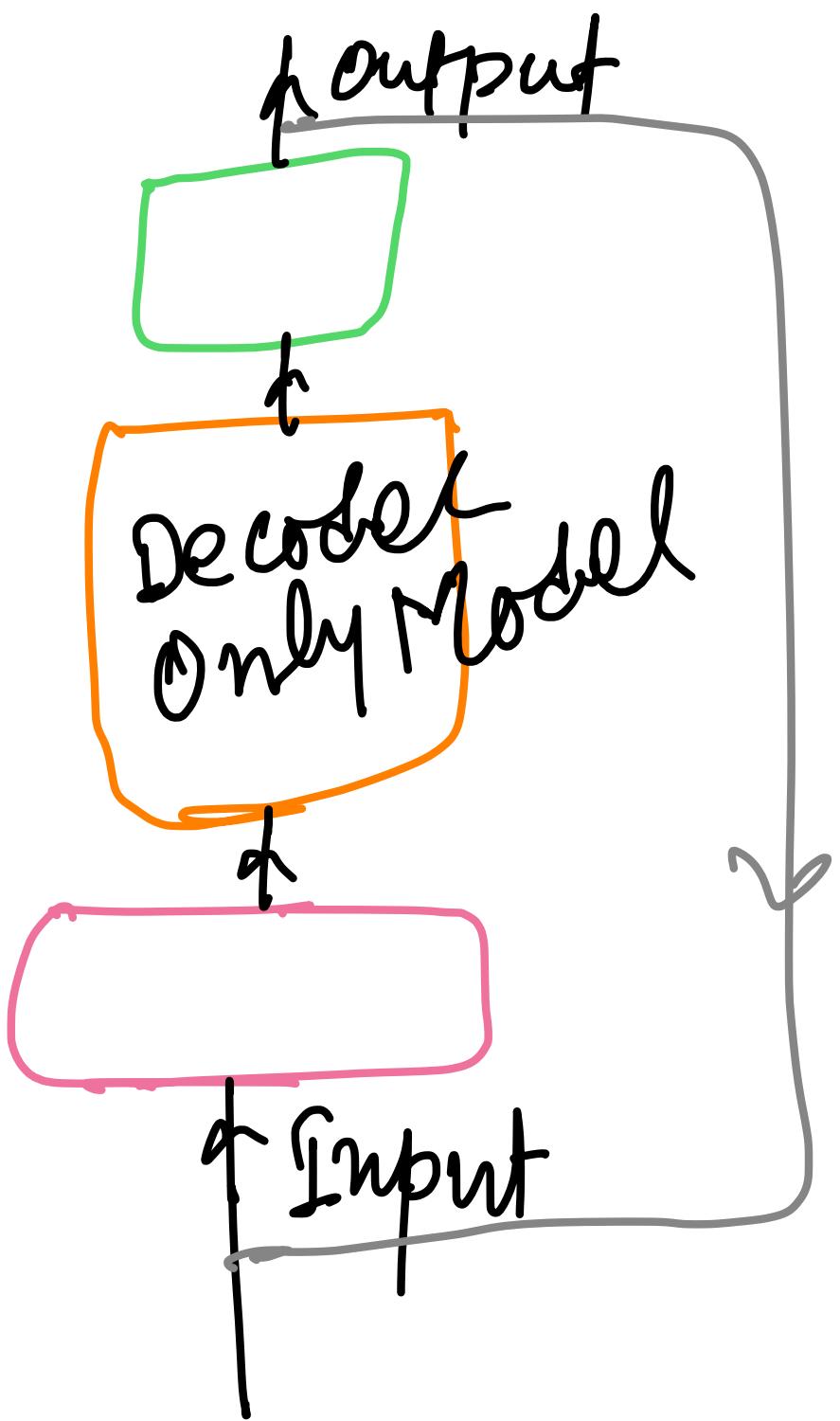


- * Context window
(typically a few 1000 words)

Completion

Where is Ganymede located in the Solar System?
Ganymede is a moon of Jupiter ...





② GPT family of
models like
Bloom, LAMA, etc..

Generative AI Project Life Cycle

Scope :- Define the use case.

Select :- Choose existing model or pre-train yours.

Adapt and Aligned model :-

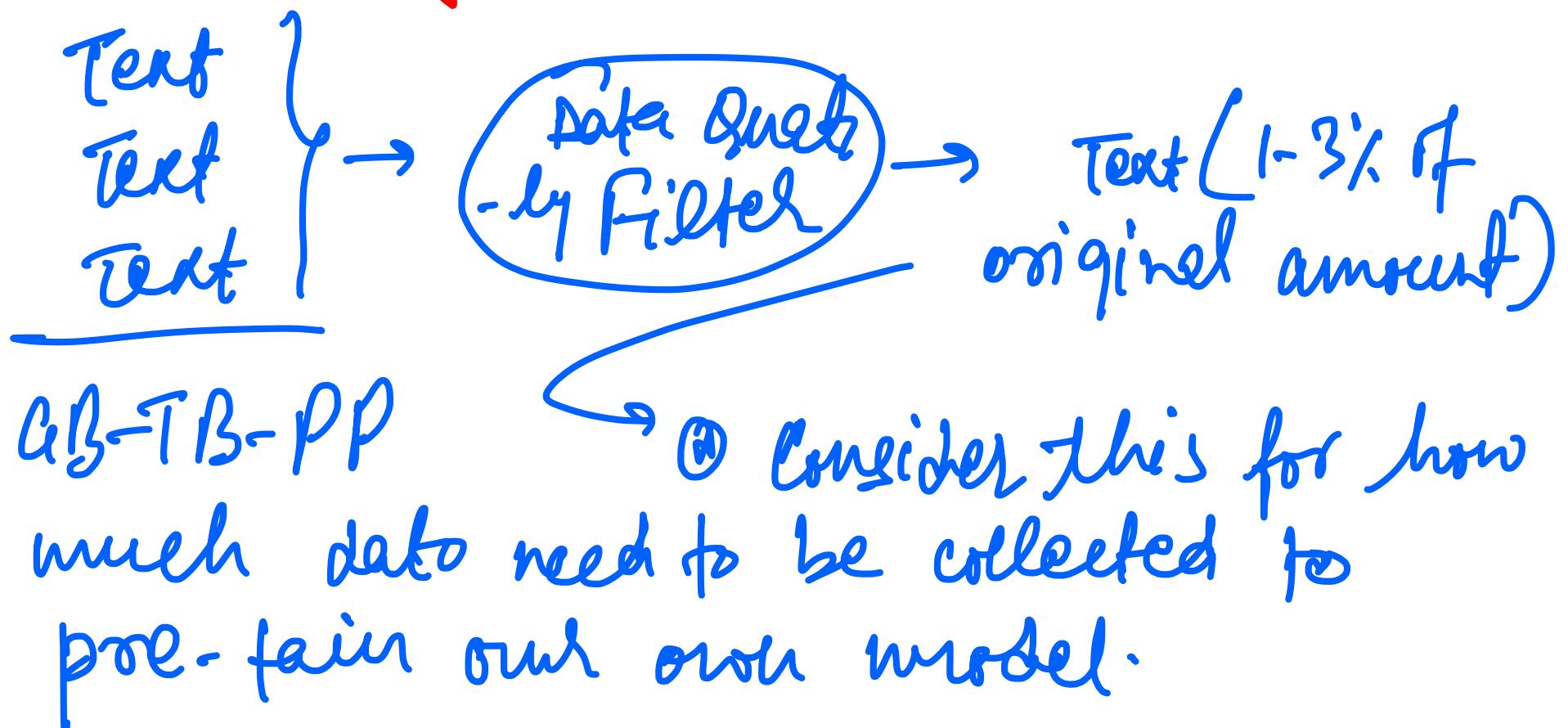
1. Prompt engineering
 2. Fine Tuning
 3. Align with Human Feedback
4. Evaluate.

Application Integration :-

1. Optimize & deploy model for inference
2. Augment model and build LLM powered applications.

Model Architectures and pre-

training Objectives :-



1. Encoder Only Models:

2. Encoder-Decoder Models:

3. Decoder only Models:

Computational Challenges of Training LLM.

Approx GPU RAM needed to store 1B params.

1 parameters = 4 bytes (32-bit float)

1B Params = 4×10^9 bytes = 4 GB

* This is just to load model in memory. If we want to train

1. Model Params (wts) = 4 bytes / param

2. Adam Optimizer (2 states) = 8 bytes / param

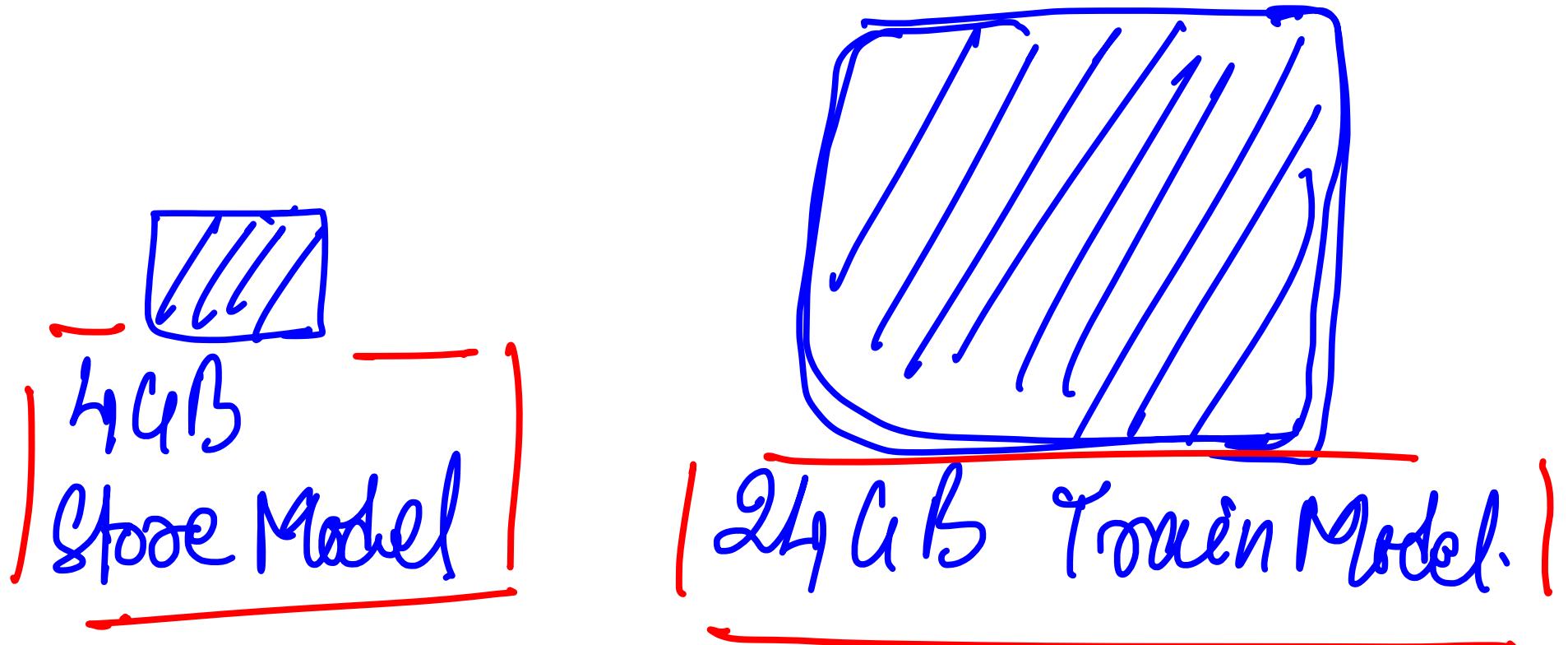
3. Gradients = 4 bytes / param

4. Activations & temp memory

$\approx 8 \text{ bytes/param}$

Summing All = 24 bytes/param

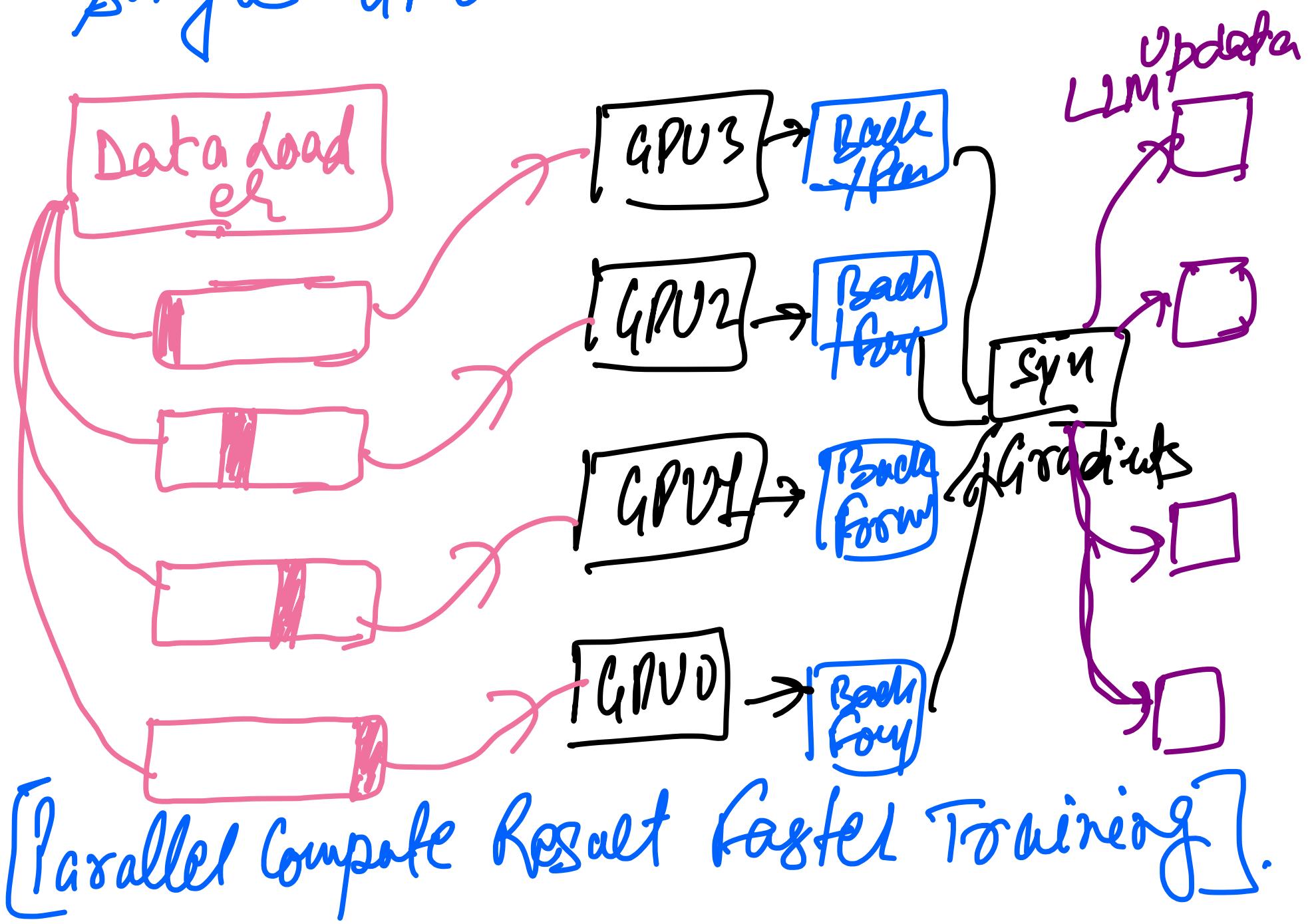
(*) It needs 6 times more GPU RAM
to train model.



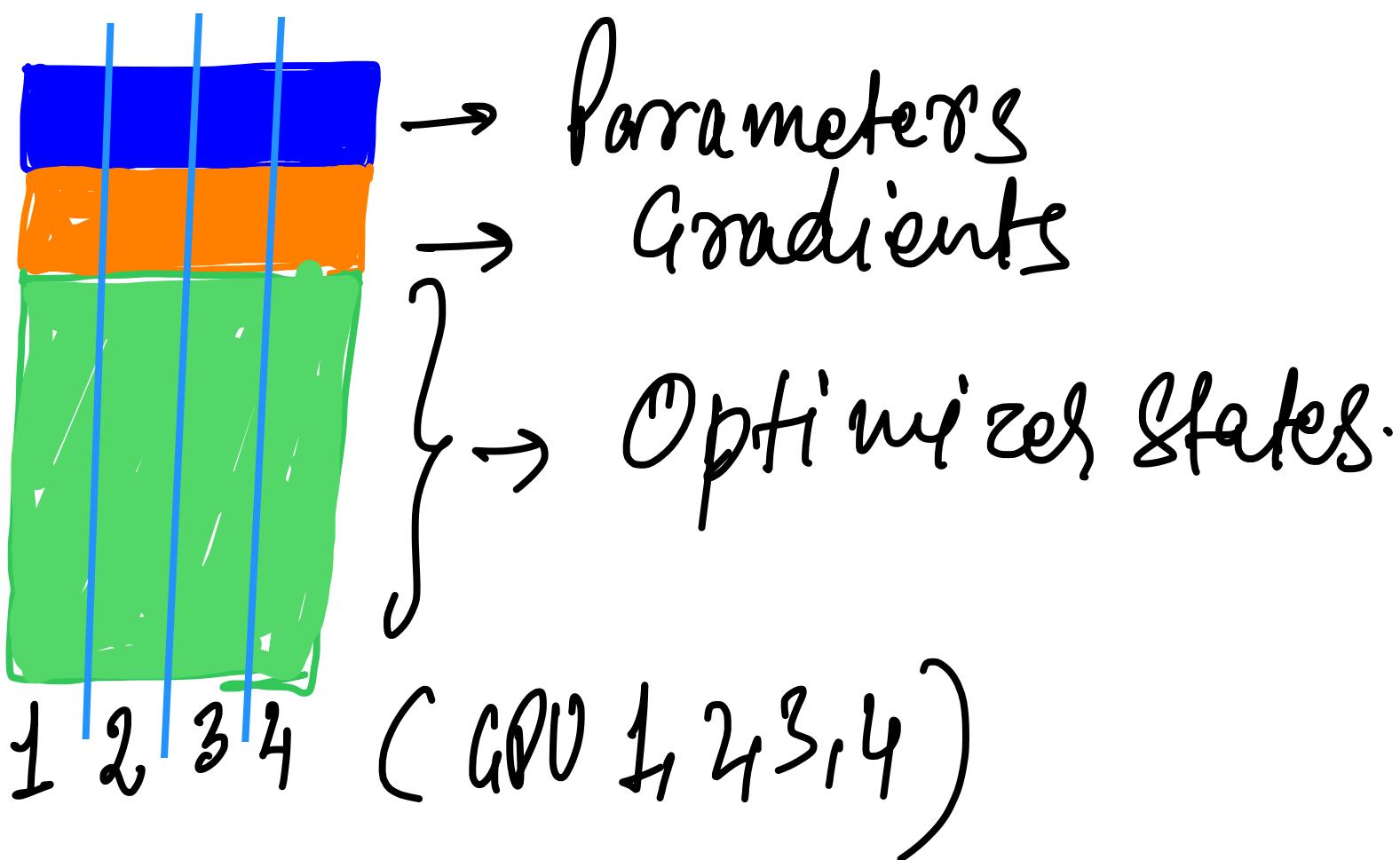
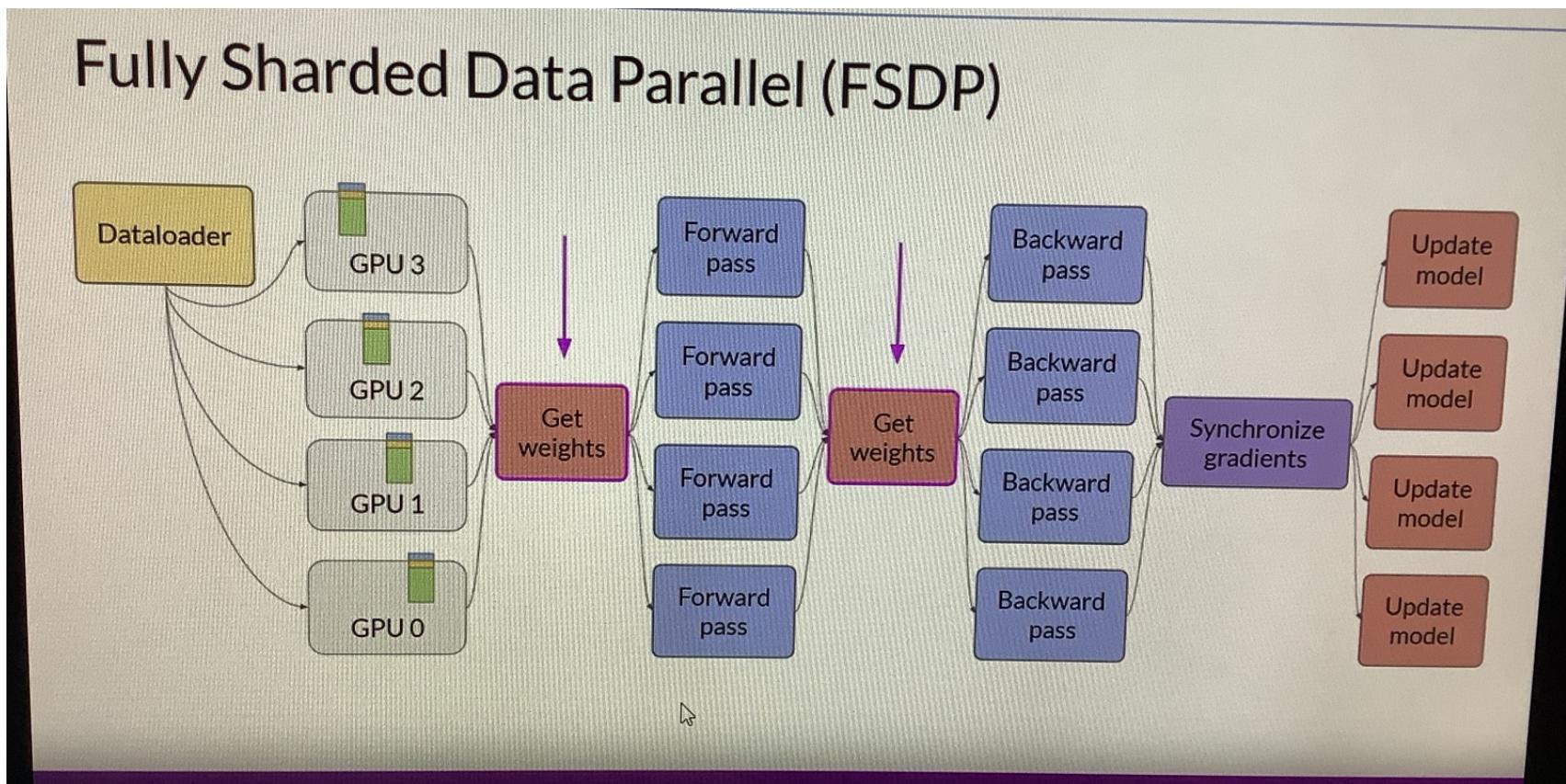
Efficient Multi-GPU Strategy

1. Distributed Data Parallel (DDP)

Assuming Model can be fit on single GPU.



Fully Sharded Data Parallel : (FSDP)



In case of FSDP we distribute

Data & Model Parameters both
across the GPUs. FSDP

Require to collect the data
before forward/Backward
pass.

Scaling laws & Compute Optimal

Chinchilla Paper := Compute

optimal training data size is $\approx 20 \times$ number of parameters.

Model	# of Params	# Compute Optimal Token
-------	-------------	-------------------------

Chinchilla	70B	$\approx 1.4T$
------------	-----	----------------

Optimal Token $\approx 20 \times$ times of no. of parameters.

Domain Adaptation / re-Training

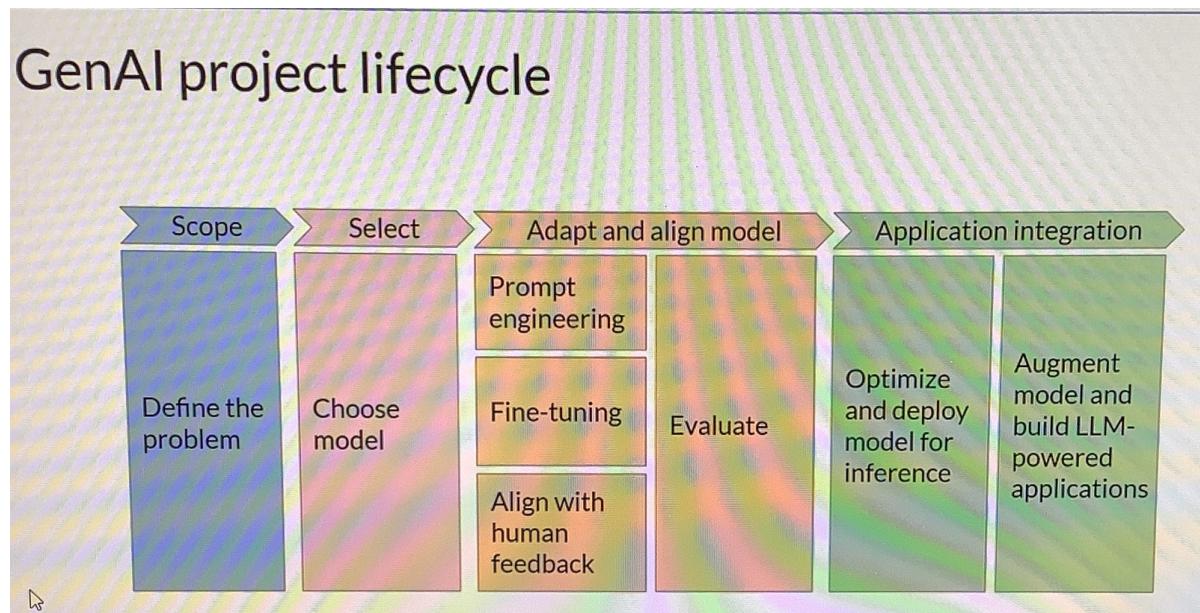
Model needs specific domain training
(Finance, Medical) etc.

Bloomberg GPT: domain Adaptation for
Finance: Bloomberg GPT was trained
on ~51% (financial public & private)
+ 49% (other) public.

{: Read \Rightarrow Bloomberg GPT: A large
language Model for Finance.}

WEEK- 02

Gen AI
life cycle



fine Tune LLM with Instruction

Prompt:

fine tuning is Supervised learning

Pre
Trained
LLM

Prompt [...] Completion [...]
Prompt [...] Completion [...]?

Fine
Tuned
LLM

Task Specific examples

Using Prompts fine Tune LLM

Prompt
Example Text
Example Completion

① Summarize the following text

② Example Text

③ Example Completion.

⇒ In this process all of the model weights are updated. Full fine tuning require memory like pre-training processes.

Data Set Preparation

Below is the sample of 'Prompt instruction templates'

Sample prompt instruction templates

Classification / sentiment analysis

```
jinja: "Given the following review:\n{{review_body}}\npredict the associated rating\\
from the following choices (1 being lowest and 5 being highest)\n- {{ answer_choices\\
| join('\\n- ')}}
|||{{answer_choices[star_rating-1]}}"
```

Text generation

```
jinja: | Generate a {{star_rating}}-star review |(1 being lowest and 5 being highest)
about this product {{product_title}}. ||| {{review_body}}
```

Text summarization

```
jinja: | Give a short sentence describing the following product review |{{review_body}}\\
\\n|{{review_headline}}|
```

Source: https://github.com/bigscience-workshop/promptsource/blob/main/promptsource/templates/amazon_polarity.j2

Most of the time fine Tuning means
instruction fine Tuning.

Fine Tuning on a Single Task

We can fine tune a Pre-Trained Model
need to perform 'A Single Task'

with [500 - 1000] examples with a
good performance. In contrast

Billions of Token model saw during
pre-training.

Downside \Rightarrow Catastrophic Forgett

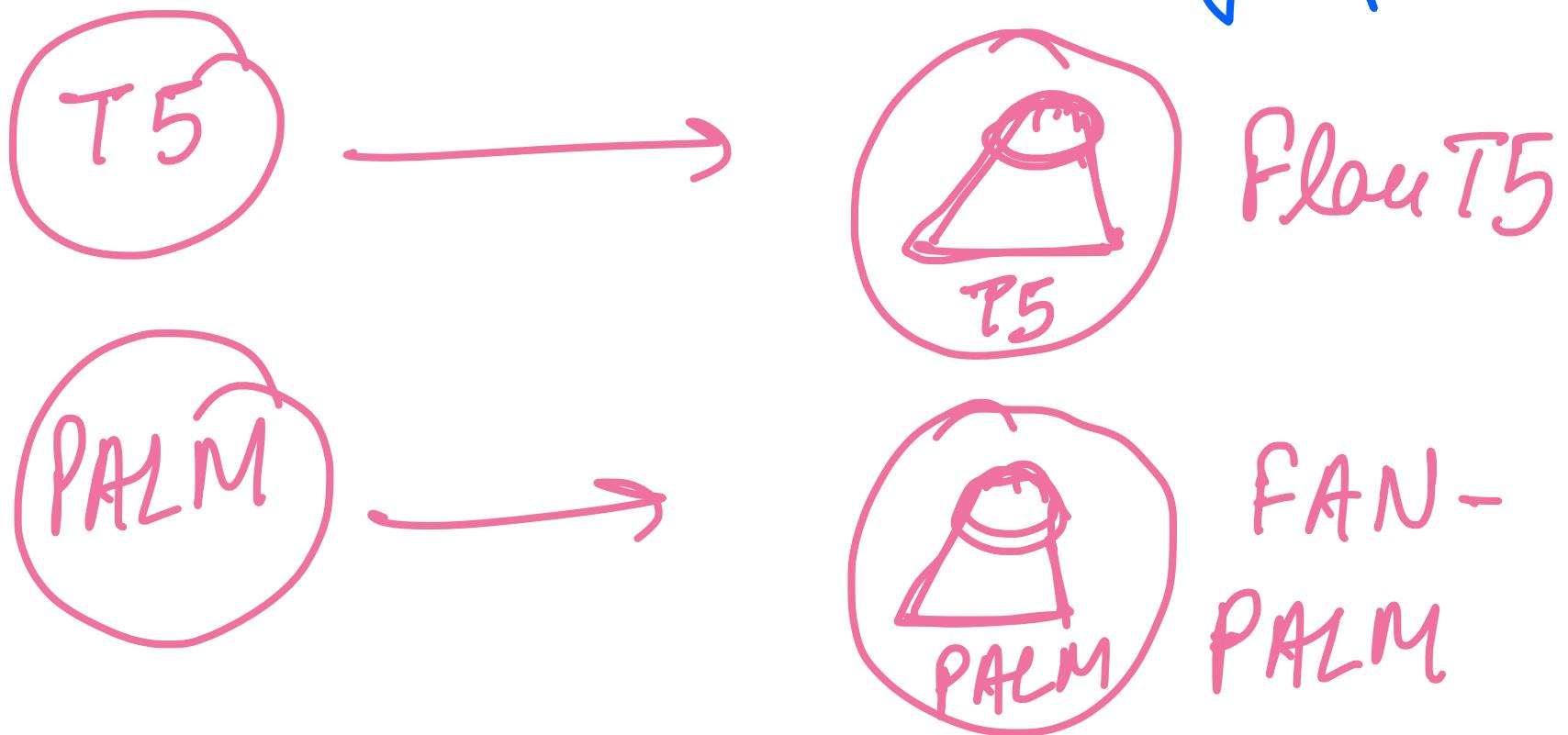
- ④ fine tuning can significantly increase model performance on specific task. It can de-grade performance on other task.

How to avoid Catastrophic forgetting

- ④ We might don't care about this
- ④ Fine-tune on multiple tasks at the same time.
- ④ Consider PEFT (Parameter Efficient Fine-Tuning)
 - ⇒ Set of Techniques that preserves the weight of original LLM and train specific Adapter layer which are task specific.

Instruction fine-tuning with
a n r
FLAN ?

FLAN (Fine-tuned Language Net)



FLAN Paper: \Rightarrow Scaling Instruction-fine tuned Language Models.

LLM Evaluation Metrics :-

1. ROUGE
2. BLEU SCORE

ROUGE:- 1. Use for text summarization.

2. Compares a summary to one or more reference summaries

BLEU:- 1. Text Generation

2. Compares to human-generated translations.

Rouge \underline{F} -1 = F1 score of Rouge-1

1 \Rightarrow Unigram Matches

Rouge-2 = F1 score of Rouge-2

2 \Rightarrow Bigram Matches.

BLEU Metric: Avg. (precision

across range of n-gram sizes)

④ Rouge & BLEU are basic metric

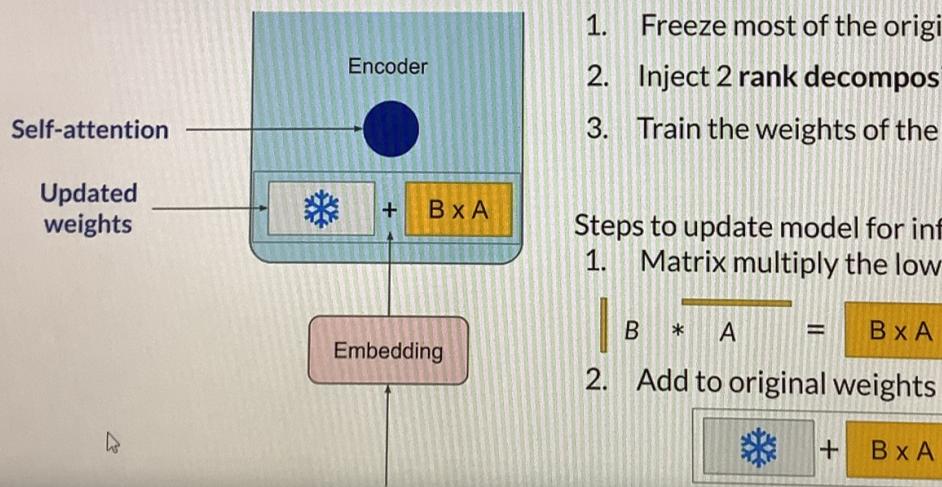
these should not be our final metric to benchmark a LLM Model.

More Complex Metrics :-

1. GLUE
2. Super GLUE
3. HELM
4. MMLU (Massive Multi task language Understanding)
5. BIG-BENCH.

LEFT (LORA)

LoRA: Low Rank Adaption of LLMs



1. Freeze most of the original LLM weights.
2. Inject 2 rank decomposition matrices
3. Train the weights of the smaller matrices

Steps to update model for inference:

1. Matrix multiply the low rank matrices

$$B * A = B \times A$$

2. Add to original weights

$$\text{Original Weights} + B \times A$$

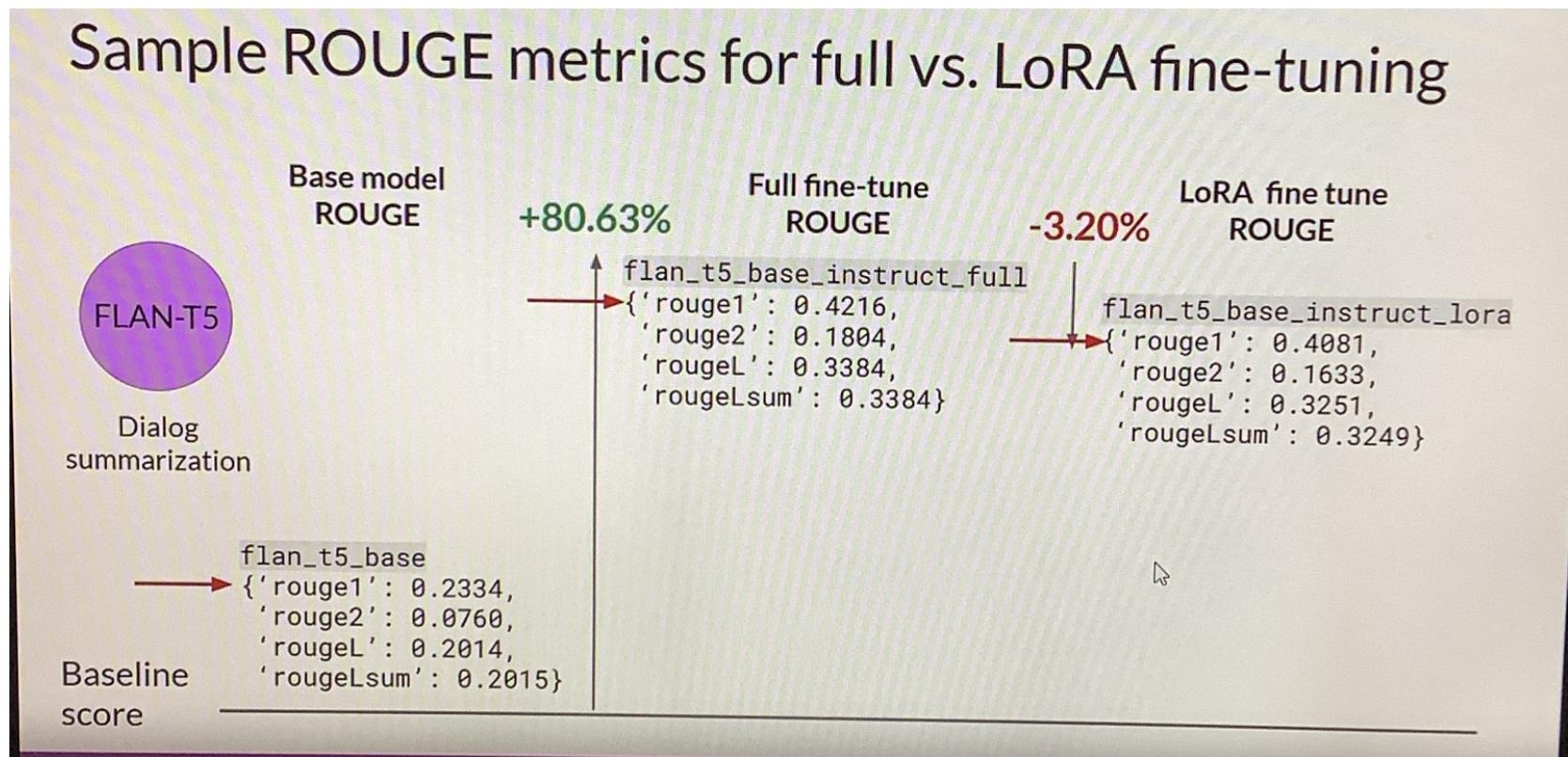
LORA intro

ces. TWO

LOW Rank

(r) typically 4, 8, - 64. We train
smaller matrices only. This requires
very small matrices to store and
train

Comparison of Base VS full Tuning



VS LoRA fine tuning Dialog Summarization.

Choosing LoRA Rank :=

Choosing the LoRA rank

Rank r	val loss	BLEU	NIST	METEOR	ROUGE_L	CIDEr
1	1.23	68.72	8.7215	0.4565	0.7052	2.4329
2	1.21	69.17	8.7413	0.4590	0.7052	2.4639
4	1.18	70.38	8.8439	0.4689	0.7186	2.5349
8	1.17	69.57	8.7457	0.4636	0.7196	2.5196
16	1.16	69.61	8.7483	0.4629	0.7177	2.4985
32	1.16	69.33	8.7736	0.4642	0.7105	2.5255
64	1.16	69.24	8.7174	0.4651	0.7180	2.5070
128	1.16	68.73	8.6718	0.4628	0.7127	2.5030
256	1.16	68.92	8.6982	0.4629	0.7128	2.5012
512	1.16	68.78	8.6857	0.4637	0.7128	2.5025
1024	1.17	69.37	8.7495	0.4659	0.7149	2.5090

Bold are the best score in each category.

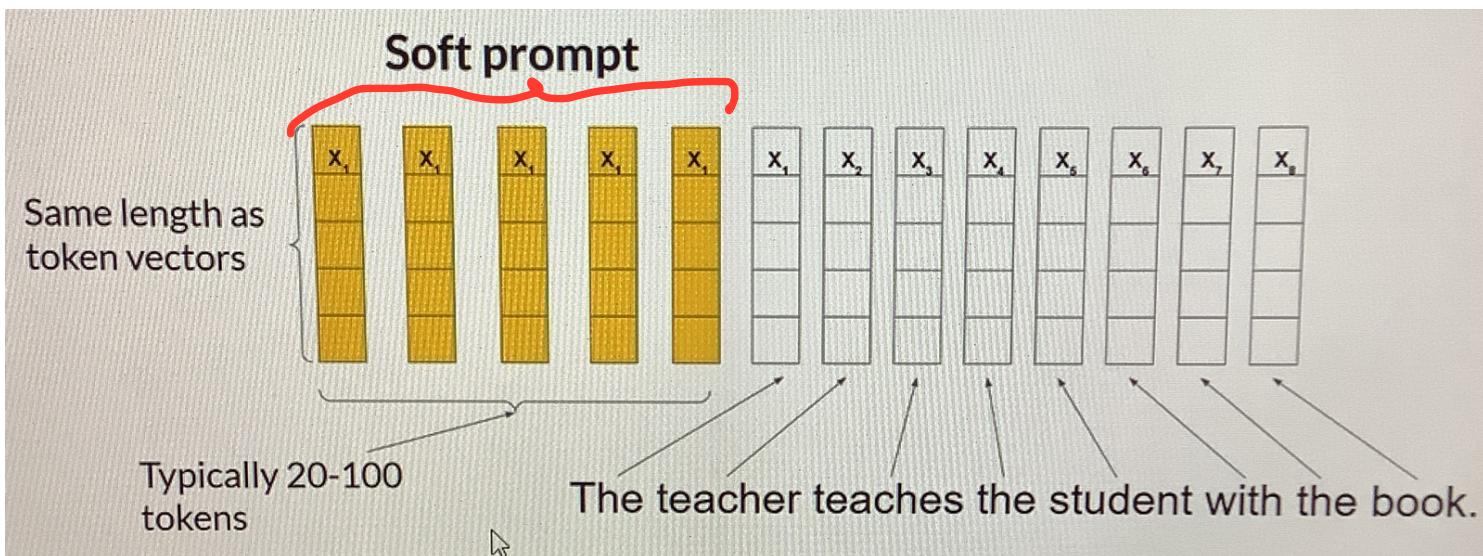
We can see for 16 Rank nothing much

improvement on those metrics.
Hence Rank (4 → 32) good tradeoff.

PEFT (Prompt Tuning)

Prompt Tuning is NOT Prompt Engineering.

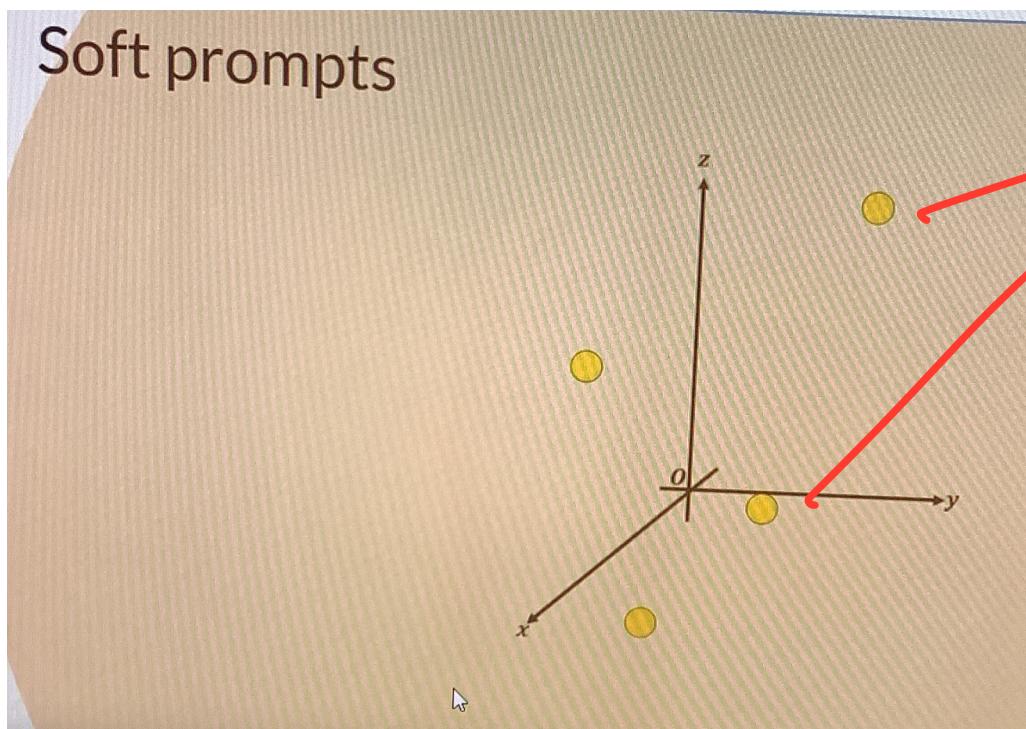
like Zero Shot or Few Shot inference. • Prompt tuning adds additional trainable token "Soft Prompts".



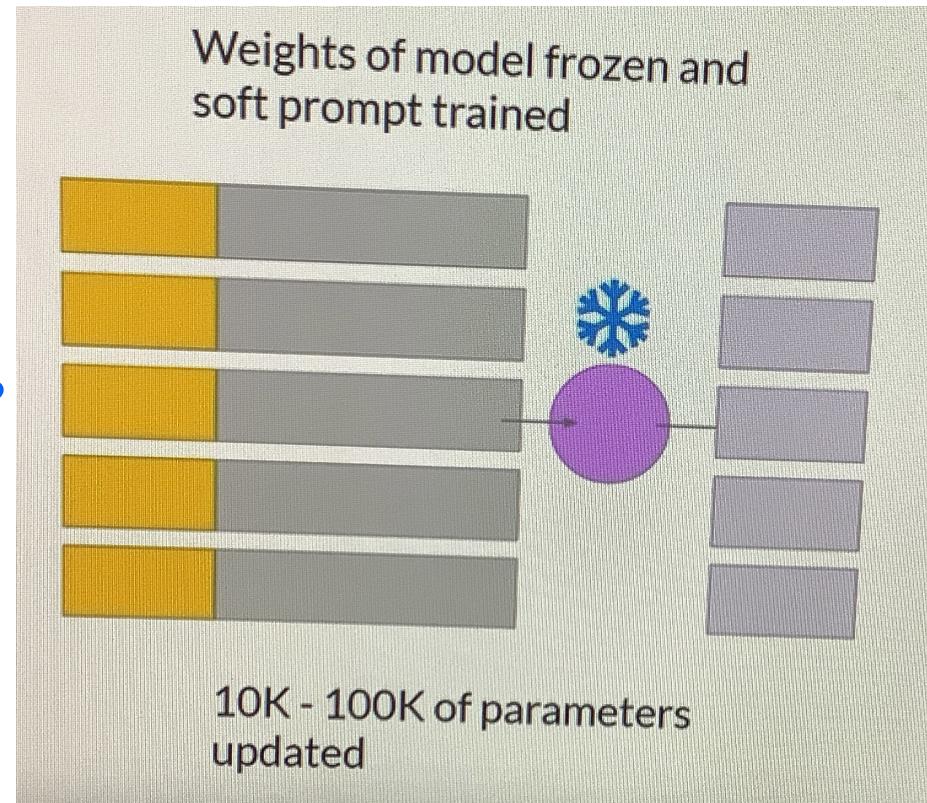
Typically
20-100
tokens
are

enough

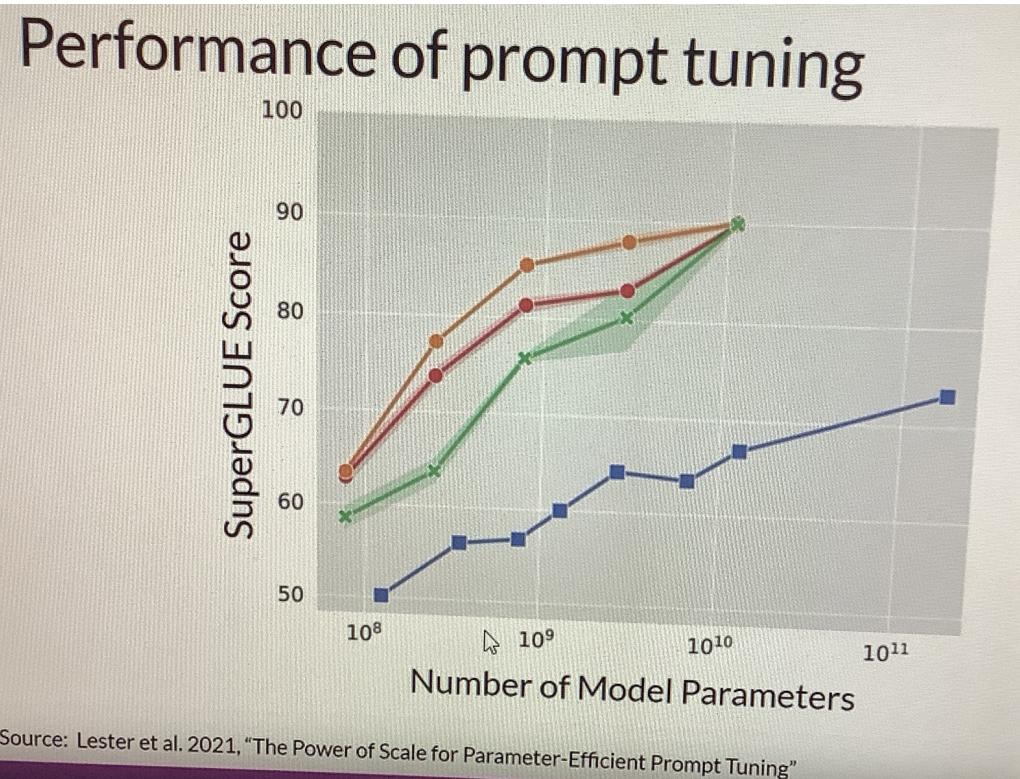
Soft prompts can be think of as Virtual Tokens that can take any value in Multidimension Embedding Space: Via Supervised learning model learns about these tokens.



Soft Prompt Training

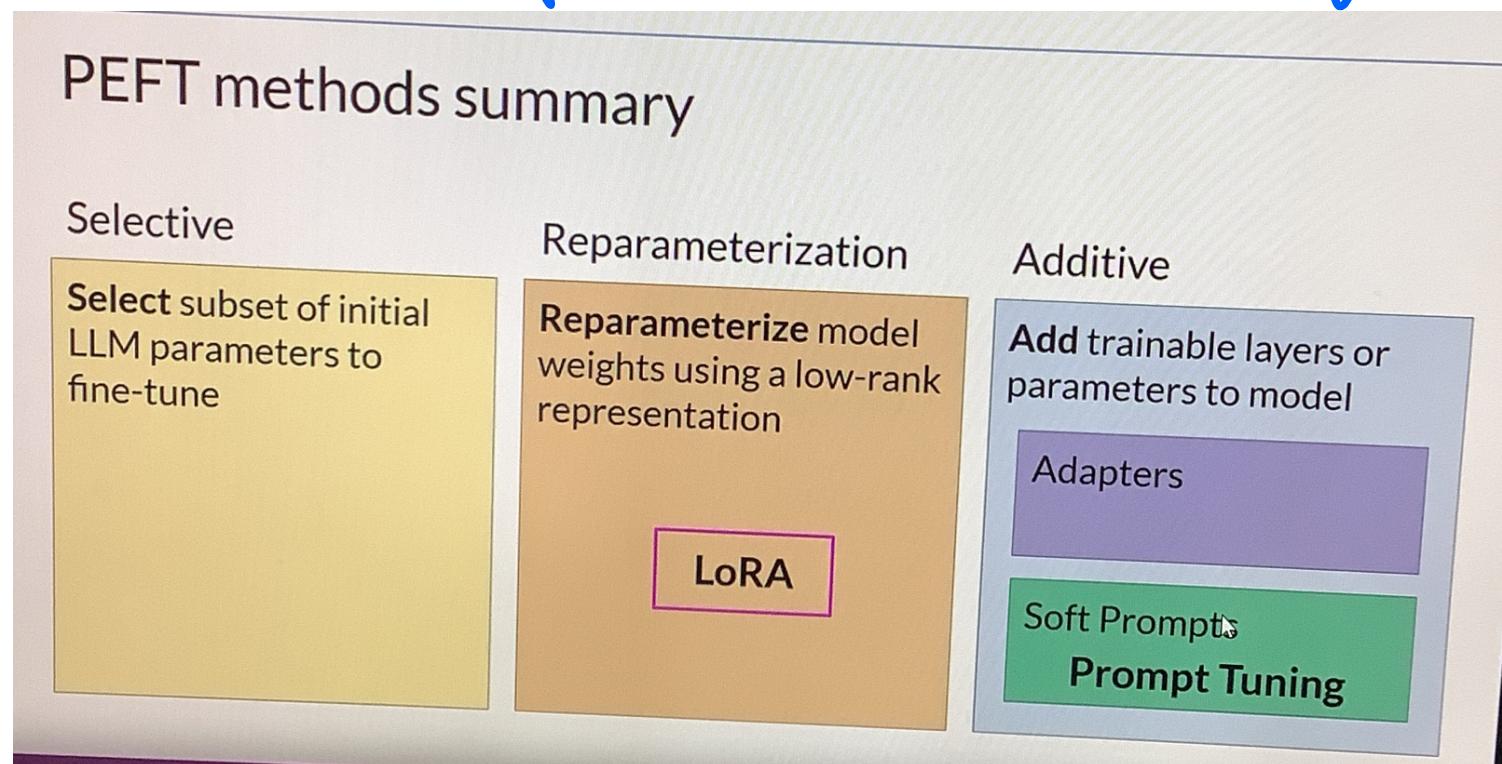


Performance Prompt Tuning:



- MultiTask fine Tuning
- Full fine Tuning
- ✗ Prompt Tuning
- Prompt Engg.

As the model size increase so the performance of prompt Tuning.



Week 2 Resources :-

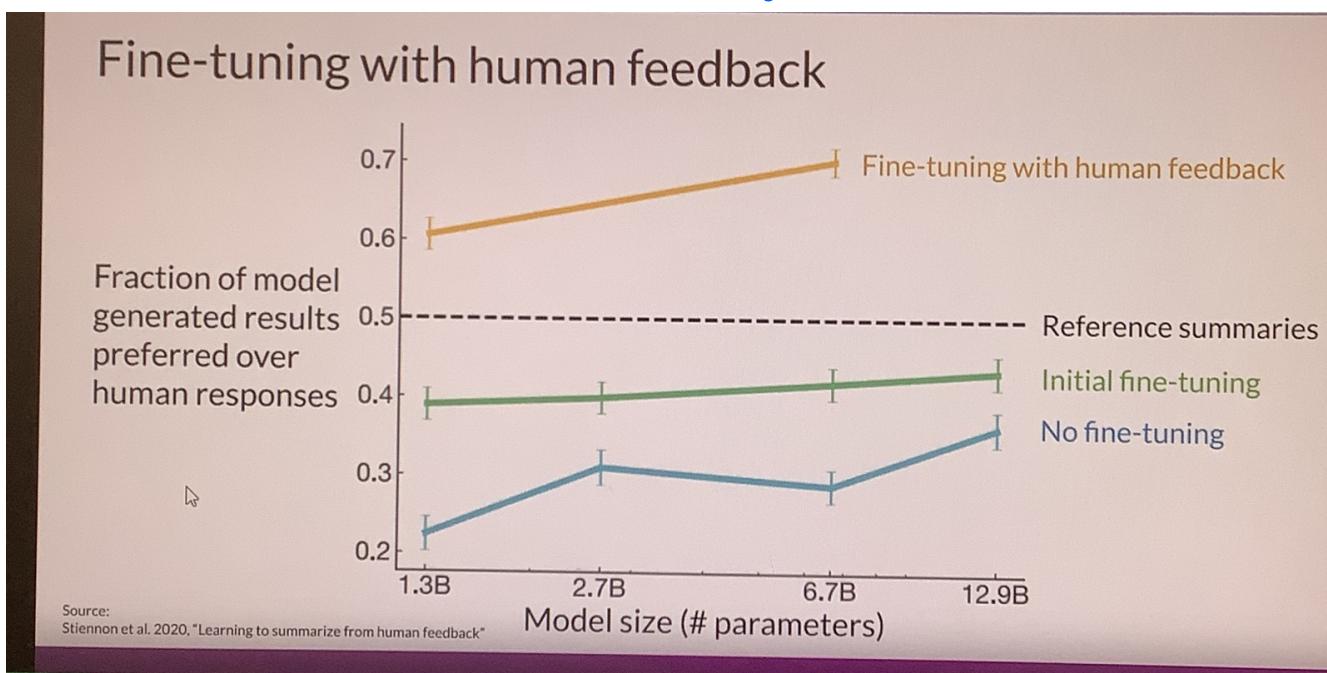
1. Scaling Instruction - finetuned LM
2. Introducing FLAN: More generalizable language Models with Instruction Fine Tuning.
3. Scaling down to scale up : A Guide to PBFT.
4. On the Effectiveness of Parameter Efficient fine Tuning.
5. LORA & OLOKA
6. The power of Scale for Parameter-Efficient Prompt Tuning.

Week - 03

Reinforcement learning from Human feedback

Models behave badly :-

1. Toxic language → Aggressive Response
2. Providing Dangerous Info.

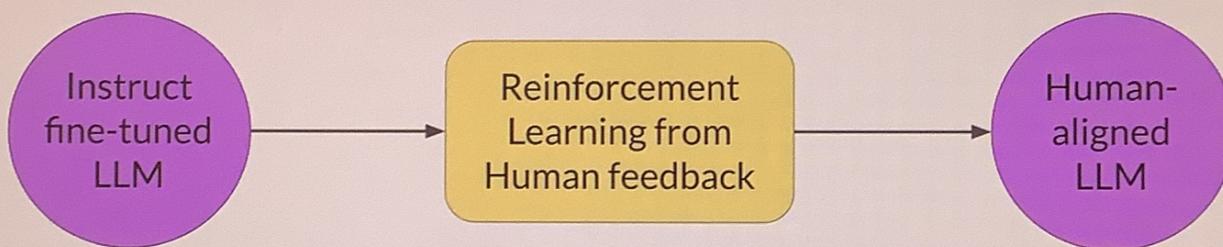


⊗ As we can see fine tuning with Human

feedback Exceed Reference Summ -aries.

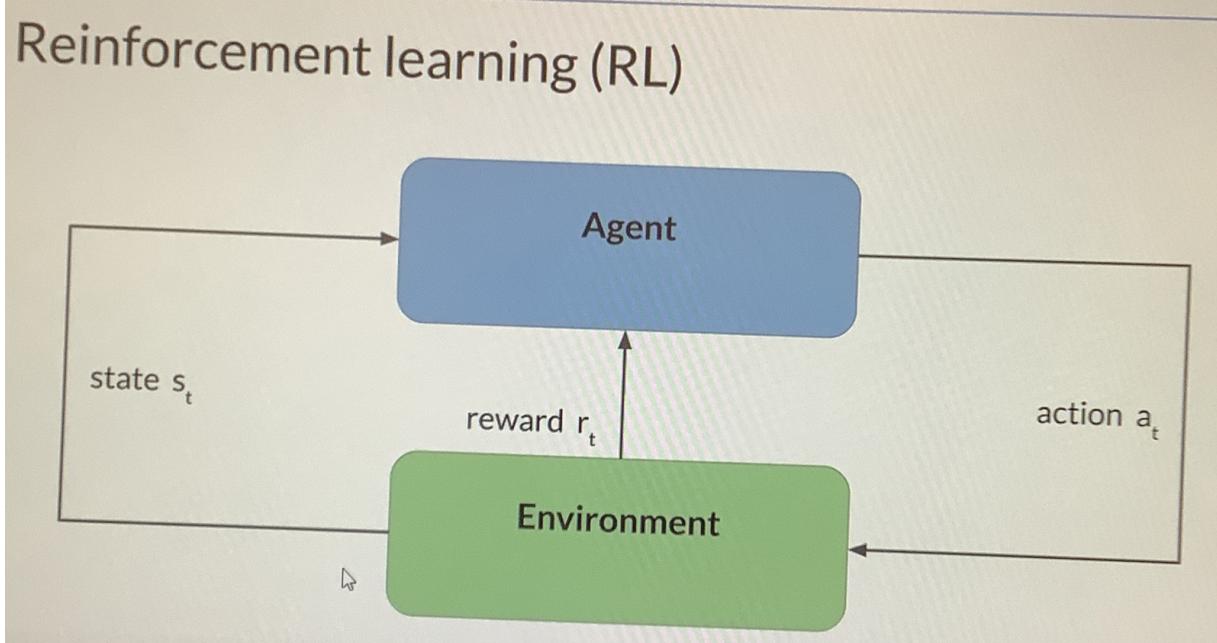
RLHF

Reinforcement learning from human feedback (RLHF)



- Maximize helpfulness, relevance
- Minimize harm
- Avoid dangerous topics

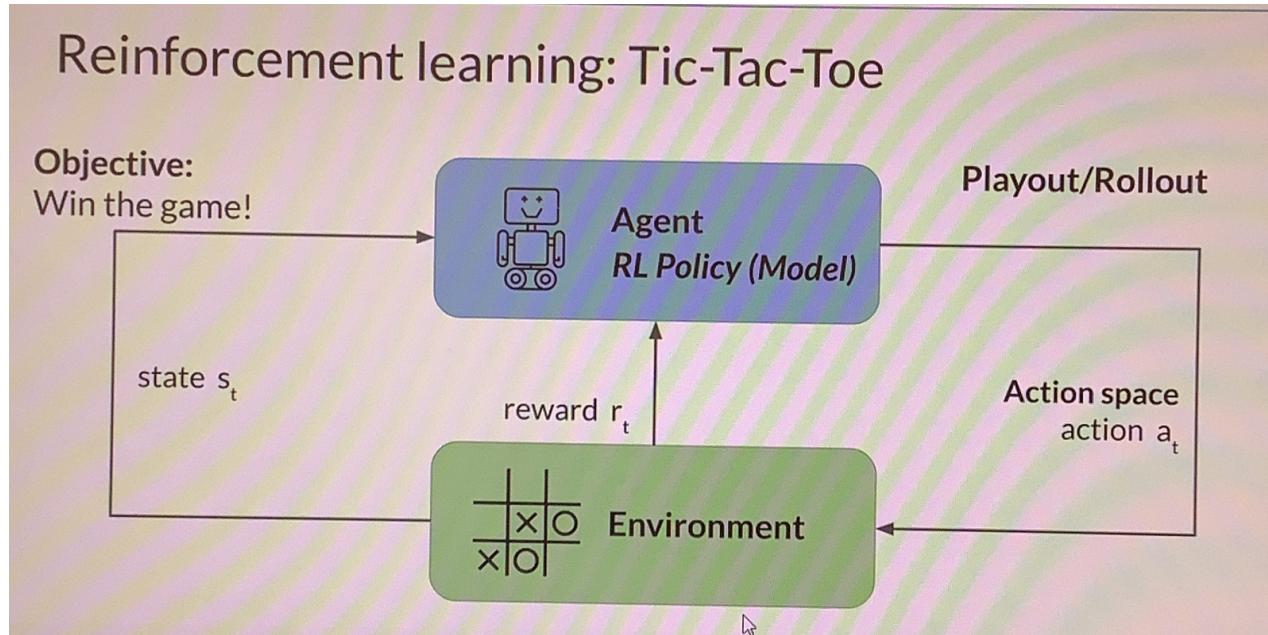
Reinforcement learning (RL)



Iterating through the loop 'AGENT'

define a policy for best performance over a task.

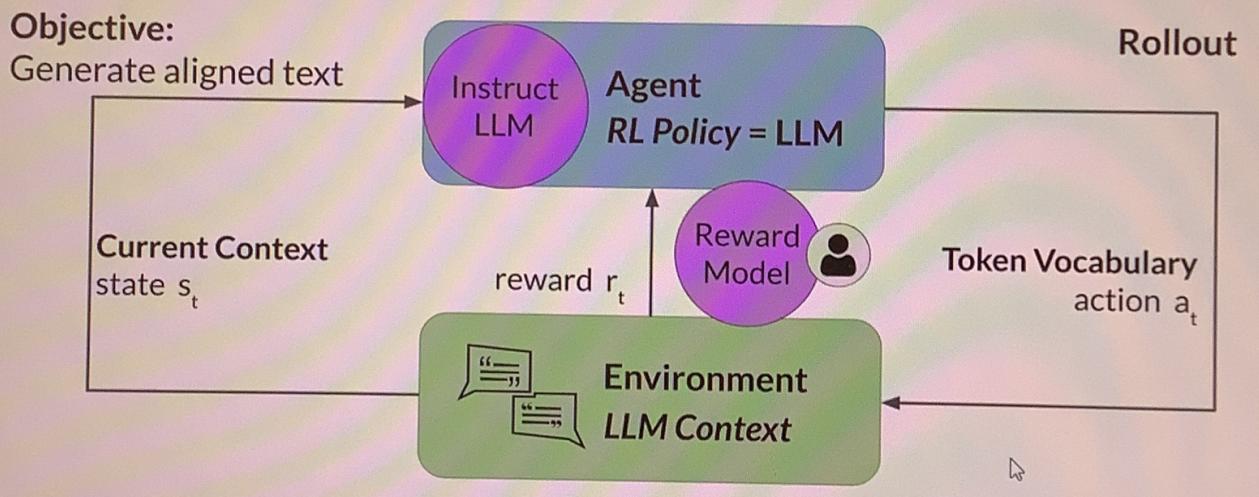
Reinforcement learning: Tic Tac -Toe



This model can be extended

to generate text aligned with human (NonToric ...)

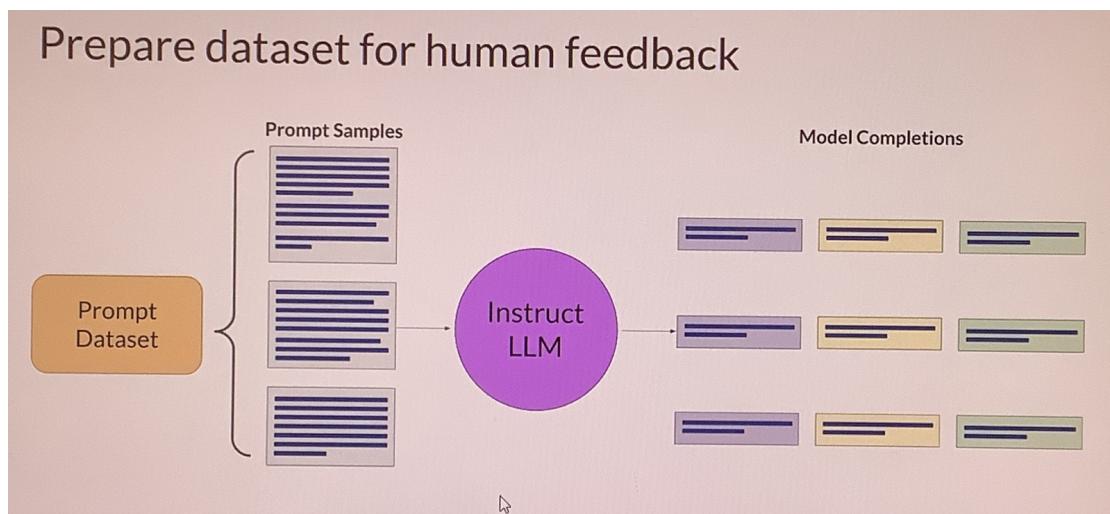
Reinforcement learning: fine-tune LLMs



Prepare Dataset for RLHF based

Training :-

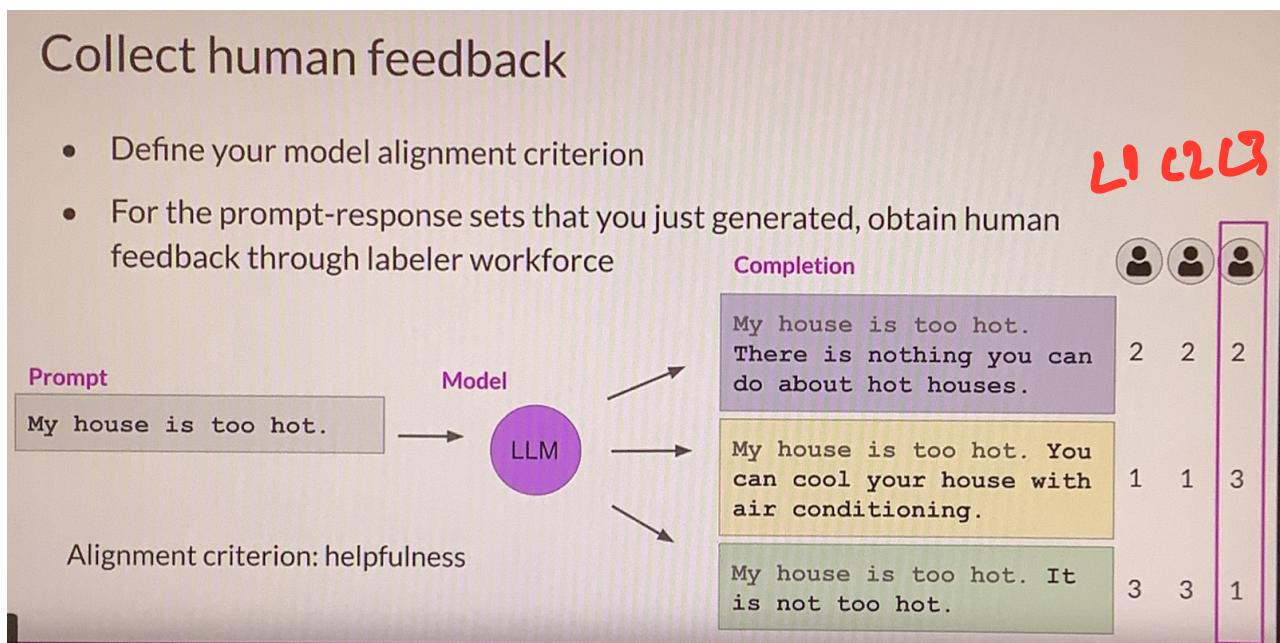
1. Choose LLM model with some capability for the task.
2. Prepare Dataset with Human feedback



3. Collect human feedback on Model completion.

Collect human feedback

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



Human (labeler) feedback

on the model from multiple sources. To minimise 'Subjectivity'. Bd: L3 feels YELLOW response is worst. Hence we have should have very clear & objective set of instructions for labelling.

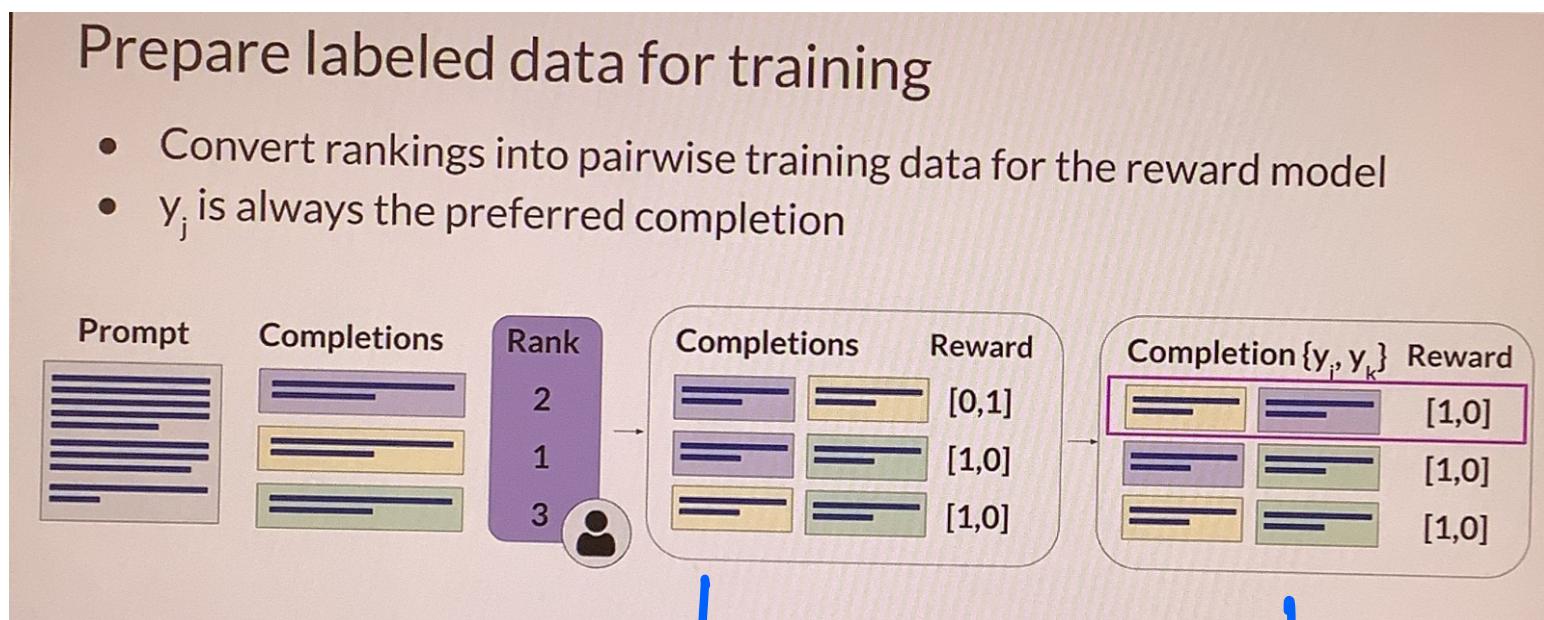
Sample Instructions for Humans. Labelers.

Sample instructions for human labelers

- * Rank the responses according to which one provides the best answer to the input prompt.
- * What is the best answer? Make a decision based on (a) the correctness of the answer, and (b) the informativeness of the response. For (a) you are allowed to search the web. Overall, use your best judgment to rank answers based on being the most useful response, which we define as one which is at least somewhat correct, and minimally informative about what the prompt is asking for.
- * If two responses provide the same correctness and informativeness by your judgment, and there is no clear winner, you may rank them the same, but please only use this sparingly.
- * If the answer for a given response is nonsensical, irrelevant, highly ungrammatical/confusing, or does not clearly respond to the given prompt, label it with "F" (for fail) rather than its rank.
- * Long answers are not always the best. Answers which provide succinct, coherent responses may be better than longer ones, if they are at least as correct and informative.

Source: Chung et al. 2022, "Scaling Instruction-Finetuned Language Models"

B Prepared Labeled Data for Training



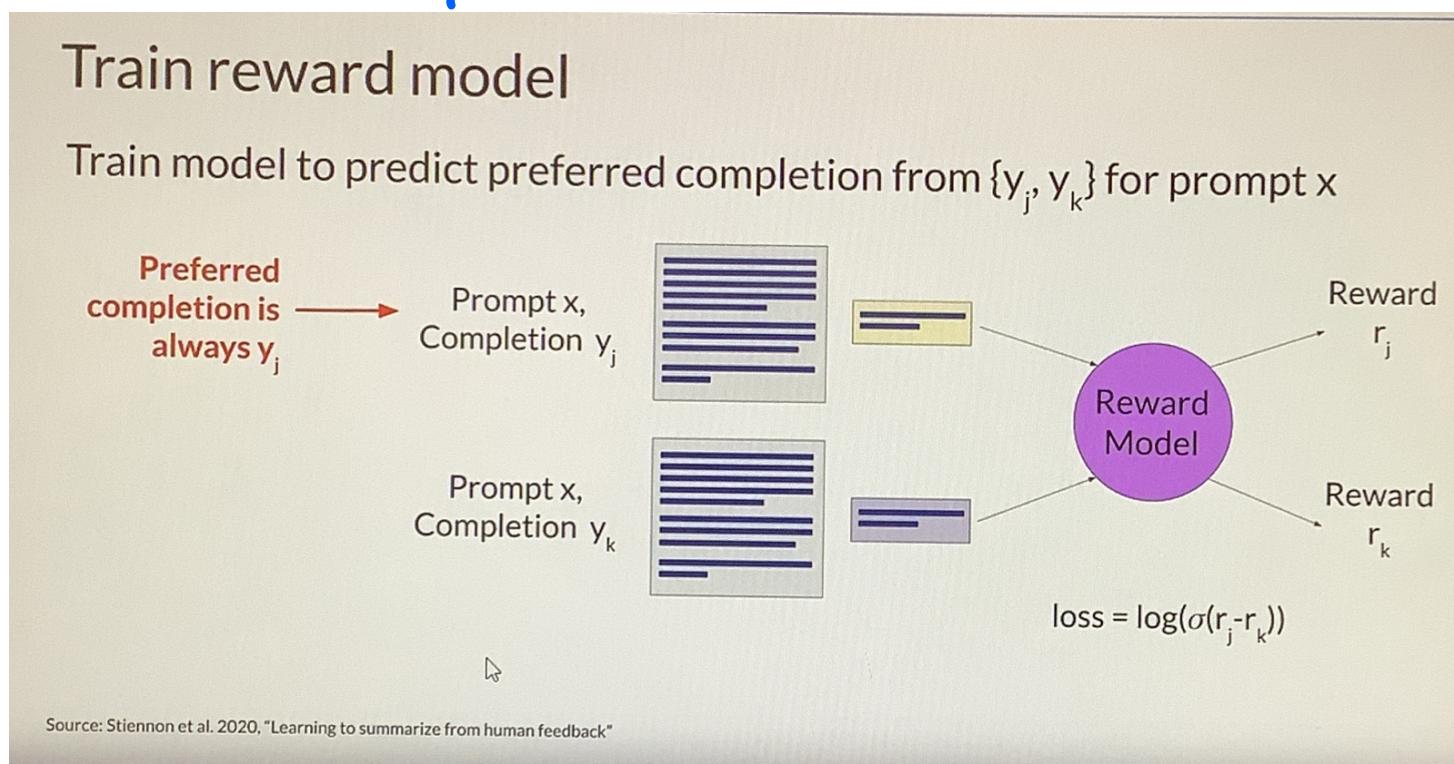
Combination
of responses

↓
Re-arrange
such that
Preferred Response
comes first.

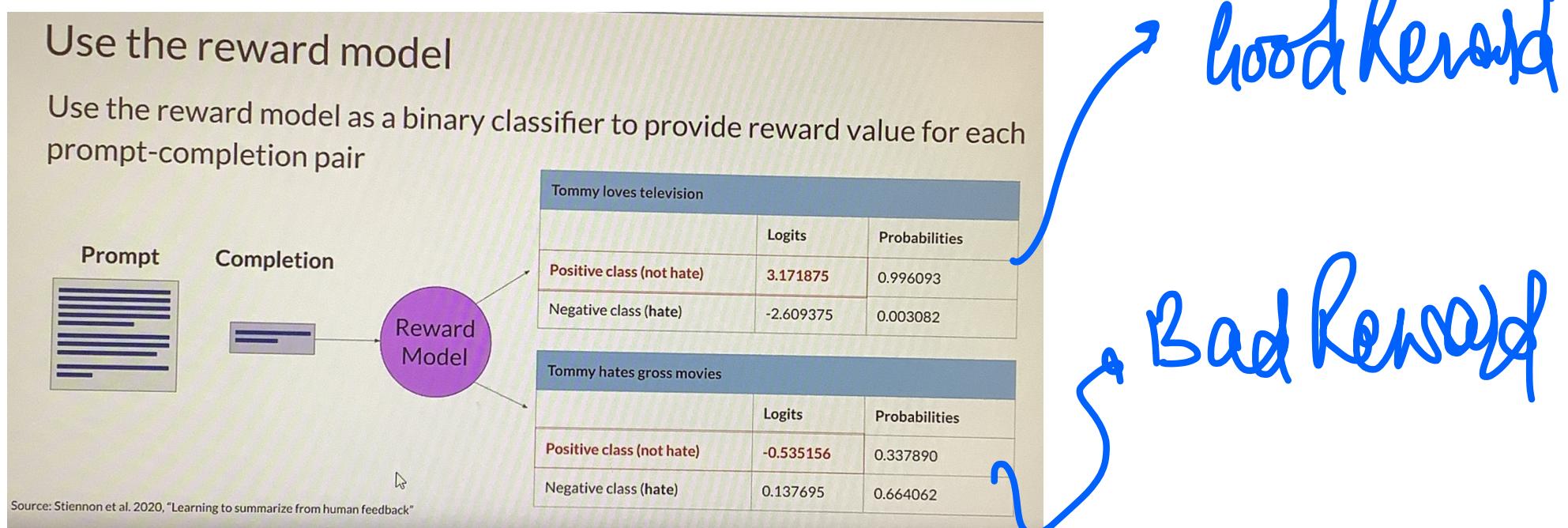
RLHF: Reward Model :-

Train reward model to predict completion from $\{y_j, y_k\}$ for prompt x .

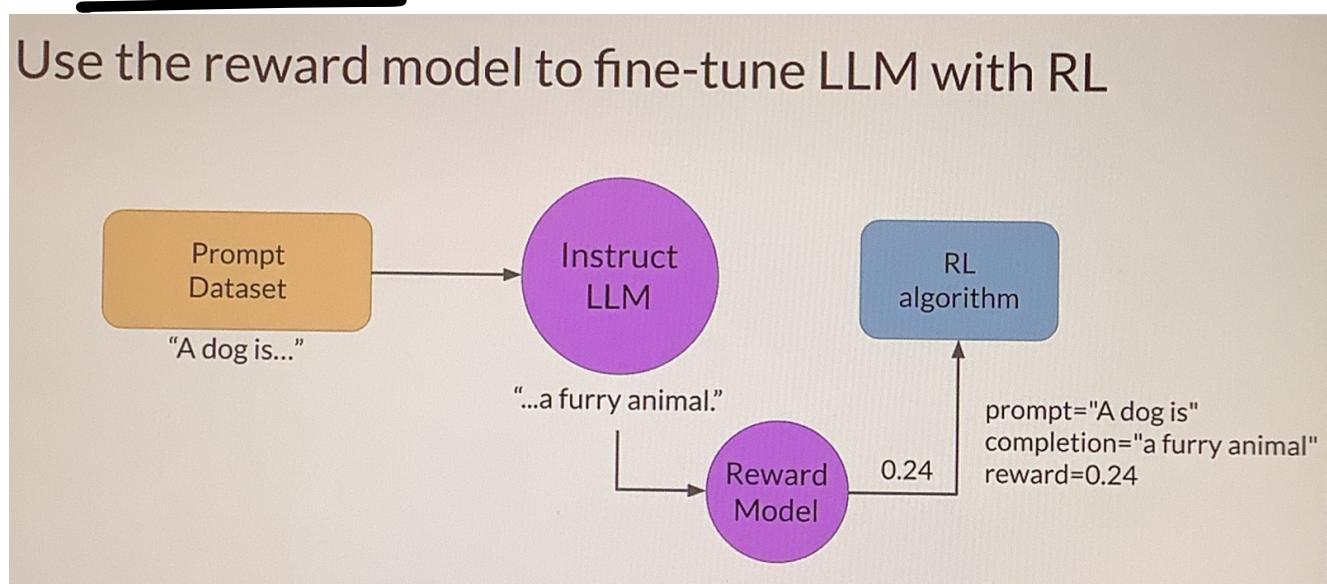
Reward model learn to favor the Human labeled completion.



Using Reward Model: Use the reward model as a binary classifier to provide reward value for each prompt-completion pair.

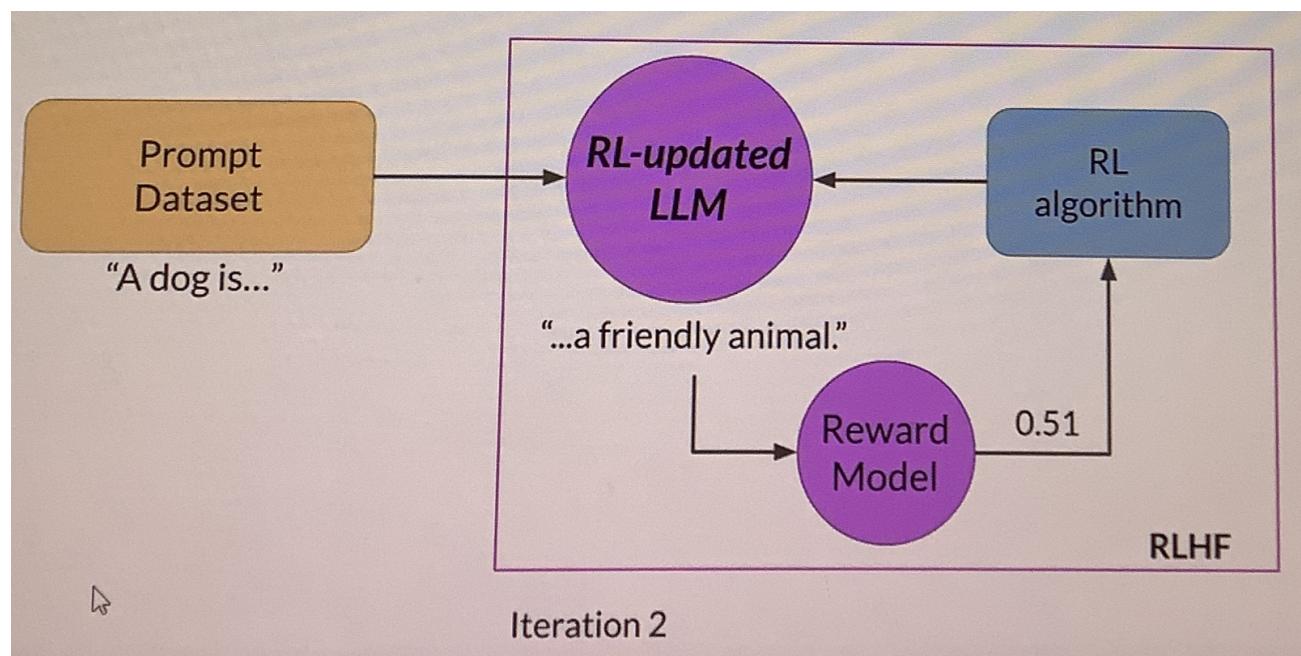
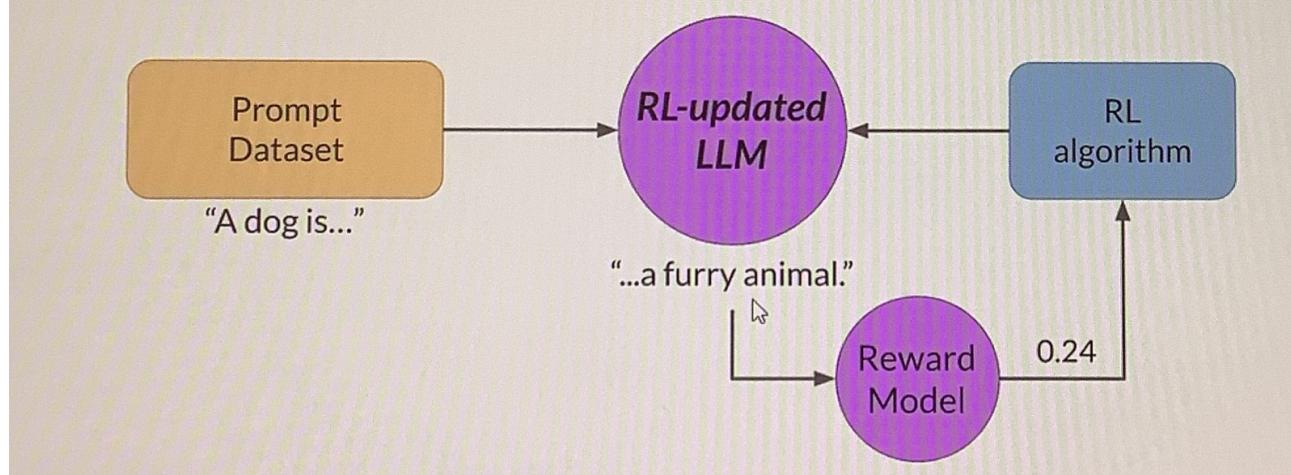


Using Reward Model to fine-tune LLM with RL

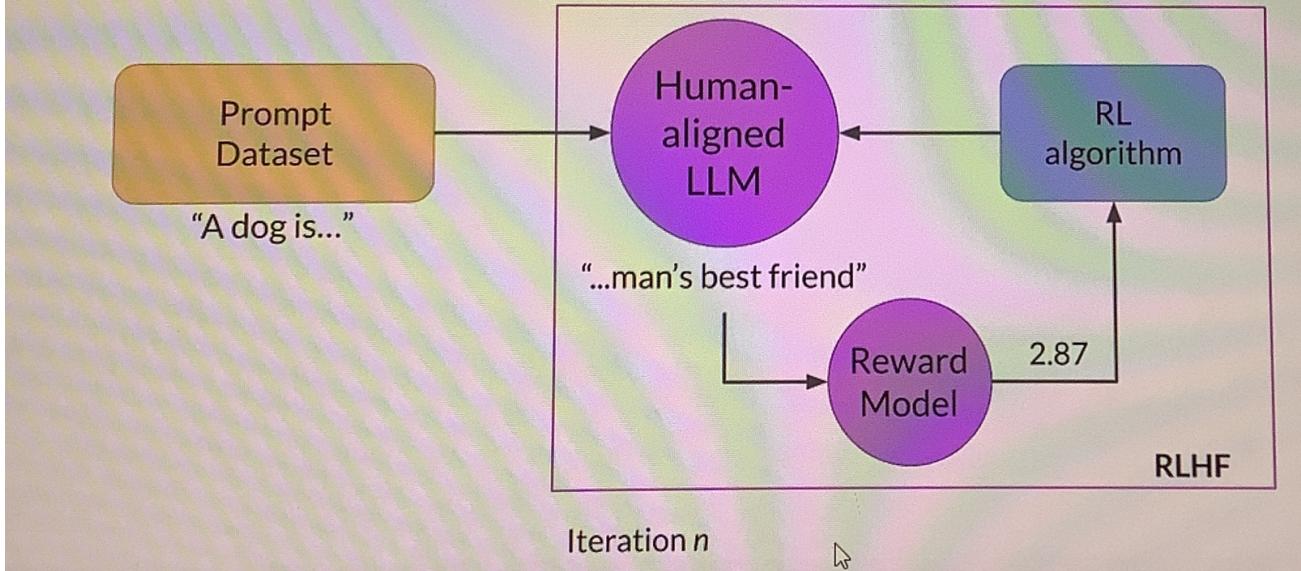


RL then update your LLM weights:

Use the reward model to fine-tune LLM with RL



Use the reward model to fine-tune LLM with RL



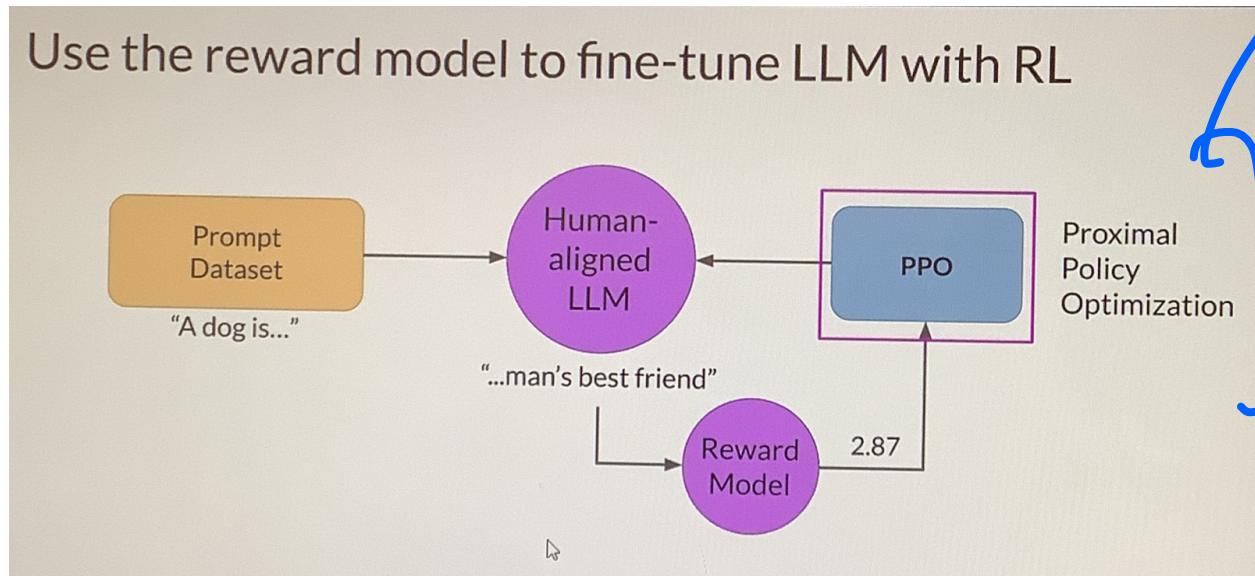
Then over multiple iterations. LLM start receiving

higher reward score.

Will continue this process until model reaches some criteria etc!
Helpfulness Score:

RL Algorithm Type:

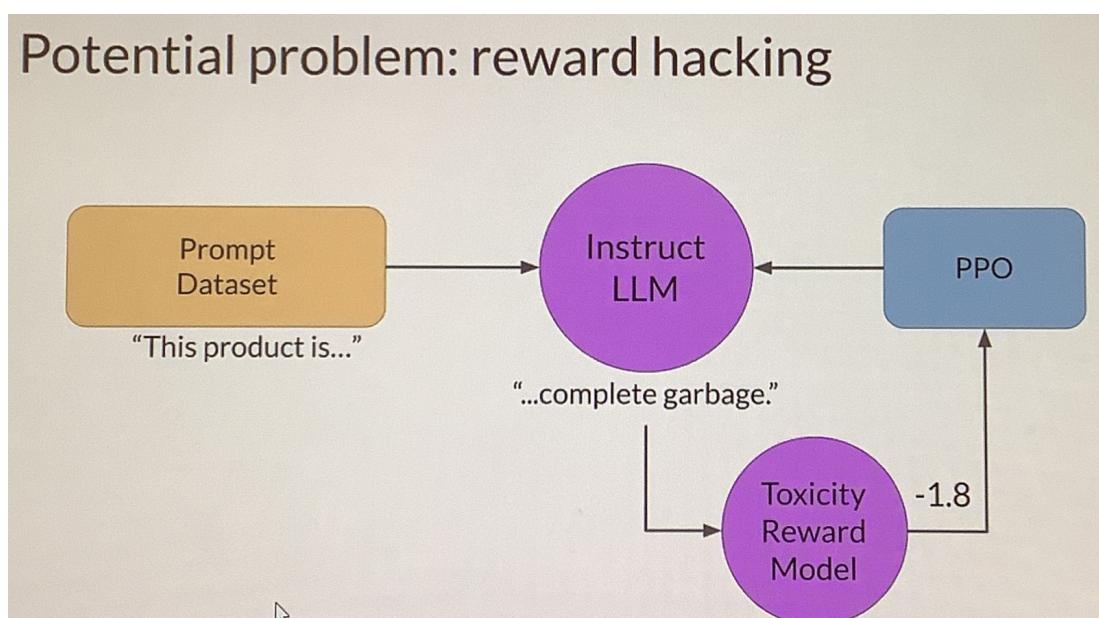
Use the reward model to fine-tune LLM with RL



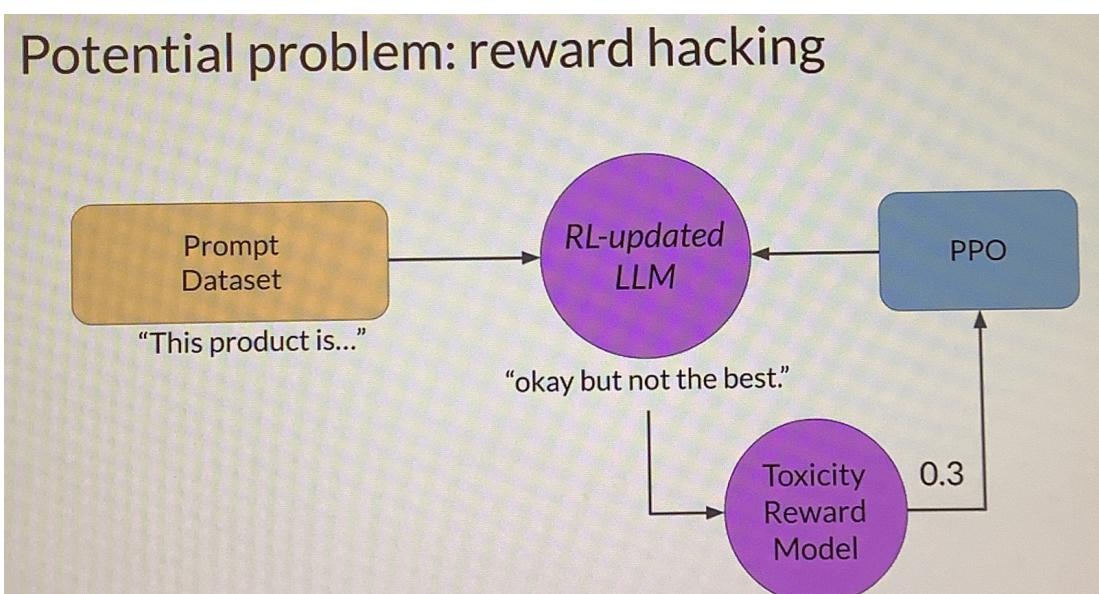
RL Algorithm
≡
(PPO)

Potential Problem : Reward Hacking

A agent learn to cheat the system by maximising reward. but the end output is not that useful.

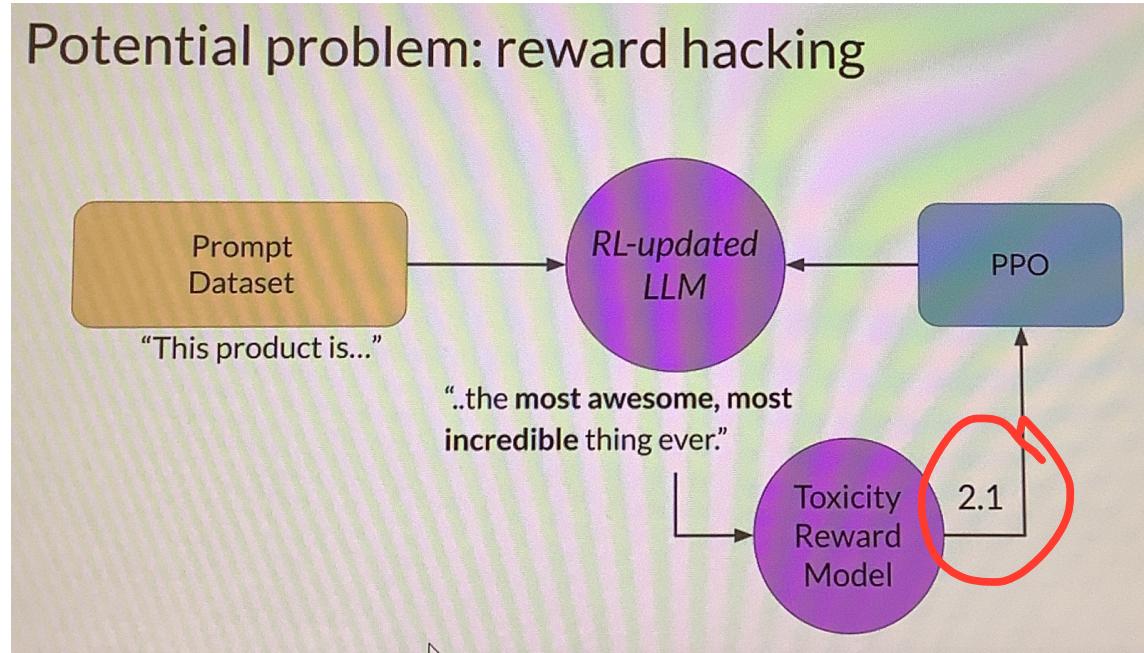


1. Toxic Response by LLM.
2. PPO algo update



the weights to less toxic response.

Potential problem: reward hacking



3. However as the policy try to reward more

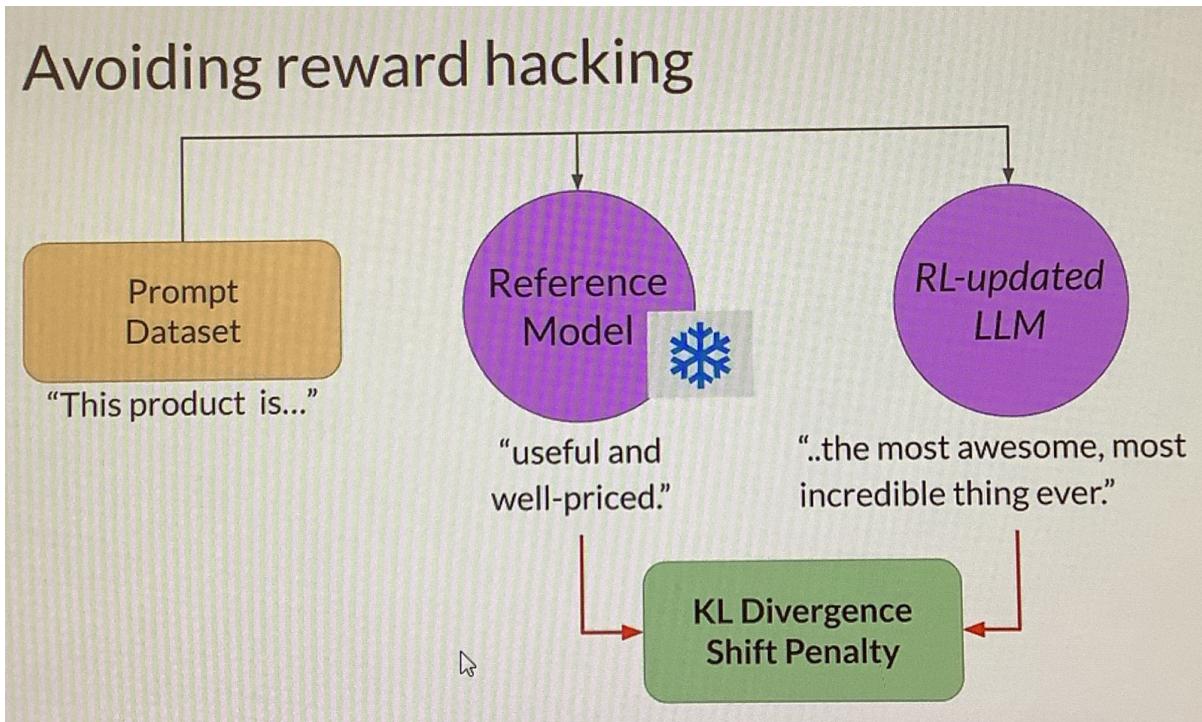
it can divert to much from initial model. Hence output is not that USE FULL any more. (Very low toxic) score

Avoid Reward Hacking :-

With KL

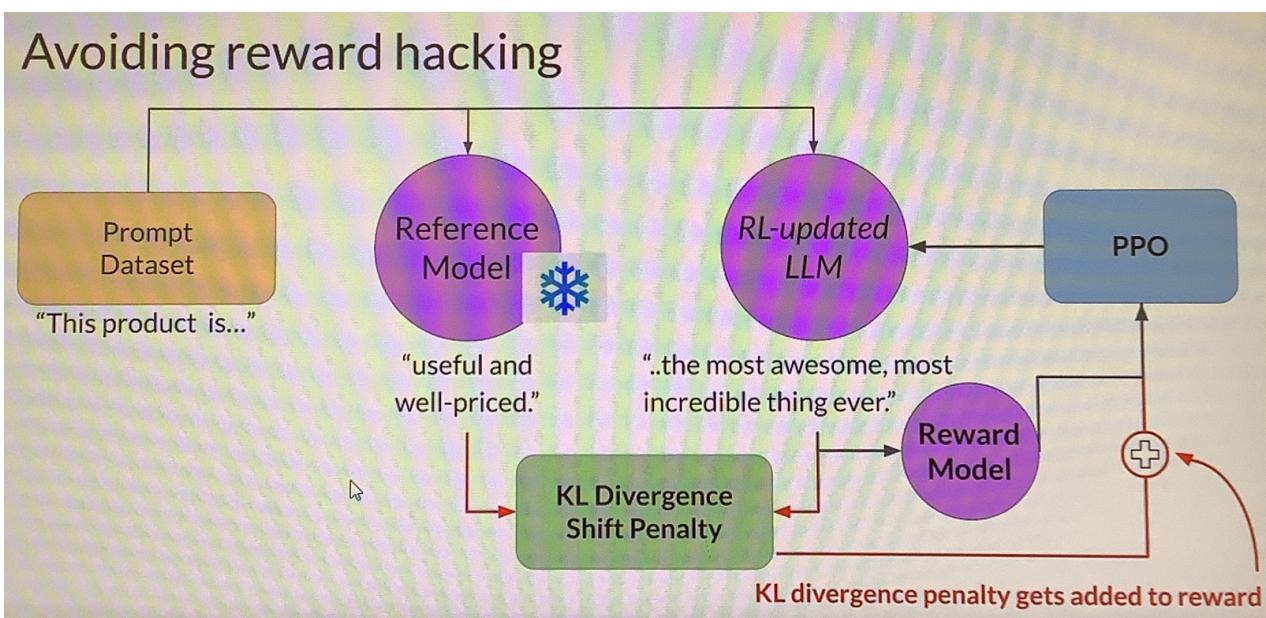
Divergence we can find how much Updated Model has a divergence.

Avoiding reward hacking



KL divergence

I core is
added to
Reward Model.

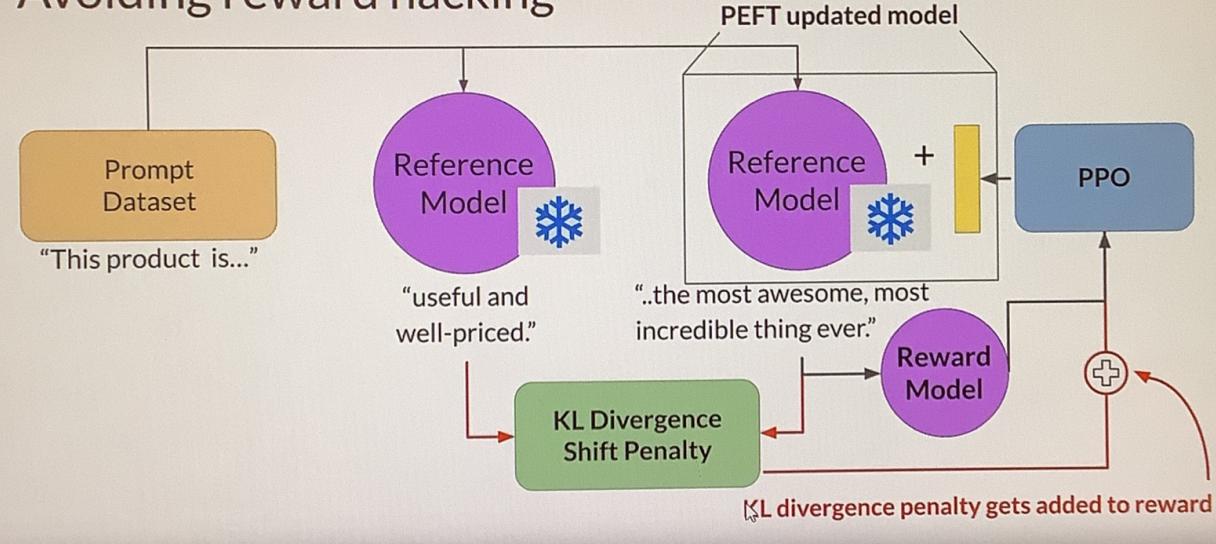


This will penalise if model is shifting too much from Original Model.

④ We now need TWO full copies of L2M. To avoid these we can combine RLHF + PEFT. In this case

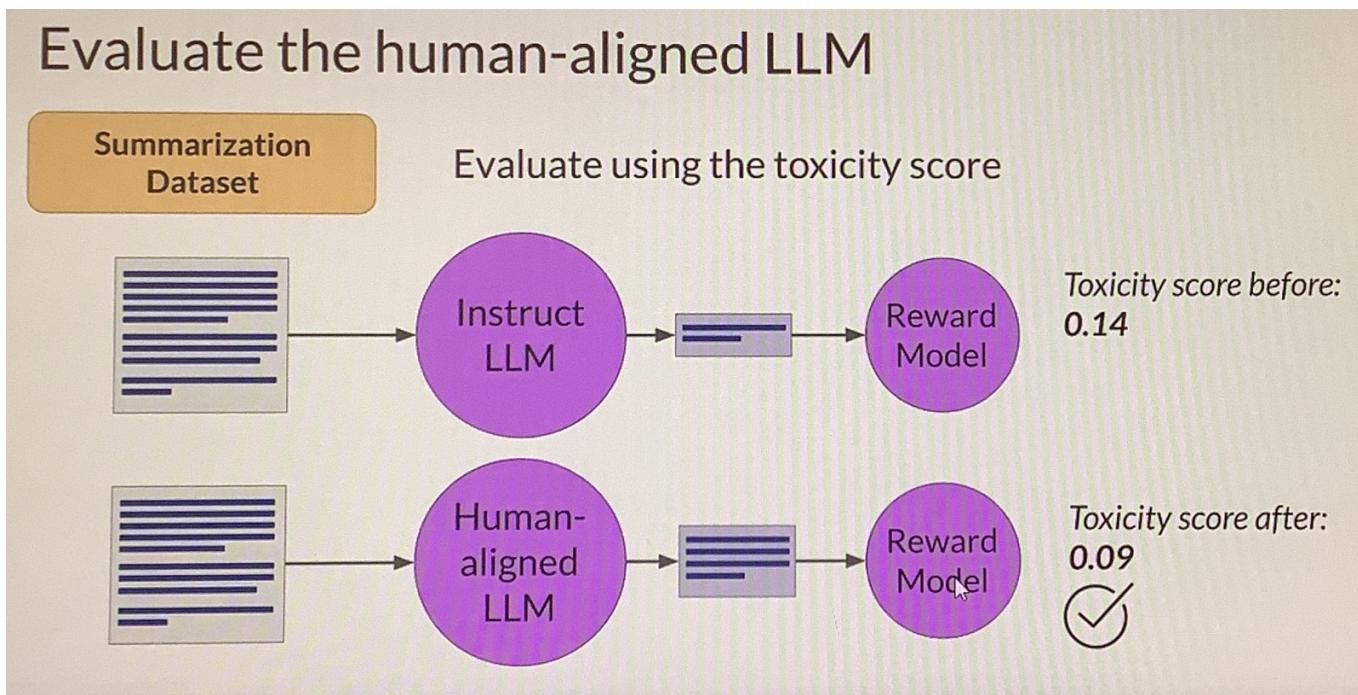
We only update the weights of PBPT adapter.

Avoiding reward hacking



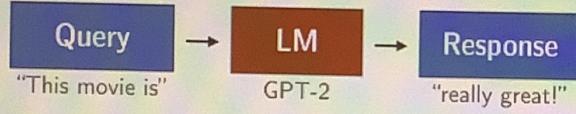
Combine PEFT with Reference model.

Evaluating Human Aligned Model :-

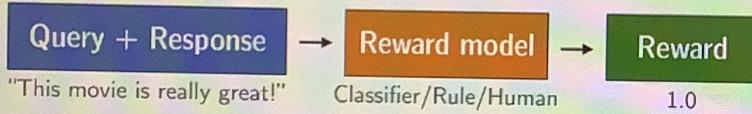


KL Divergence :-

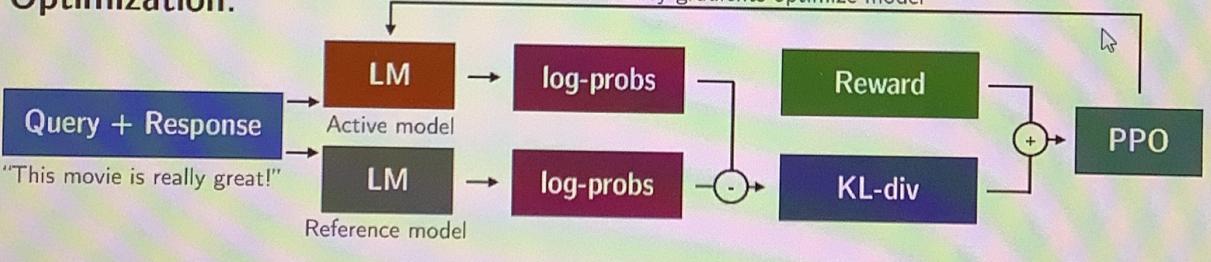
Rollout:



Evaluation:

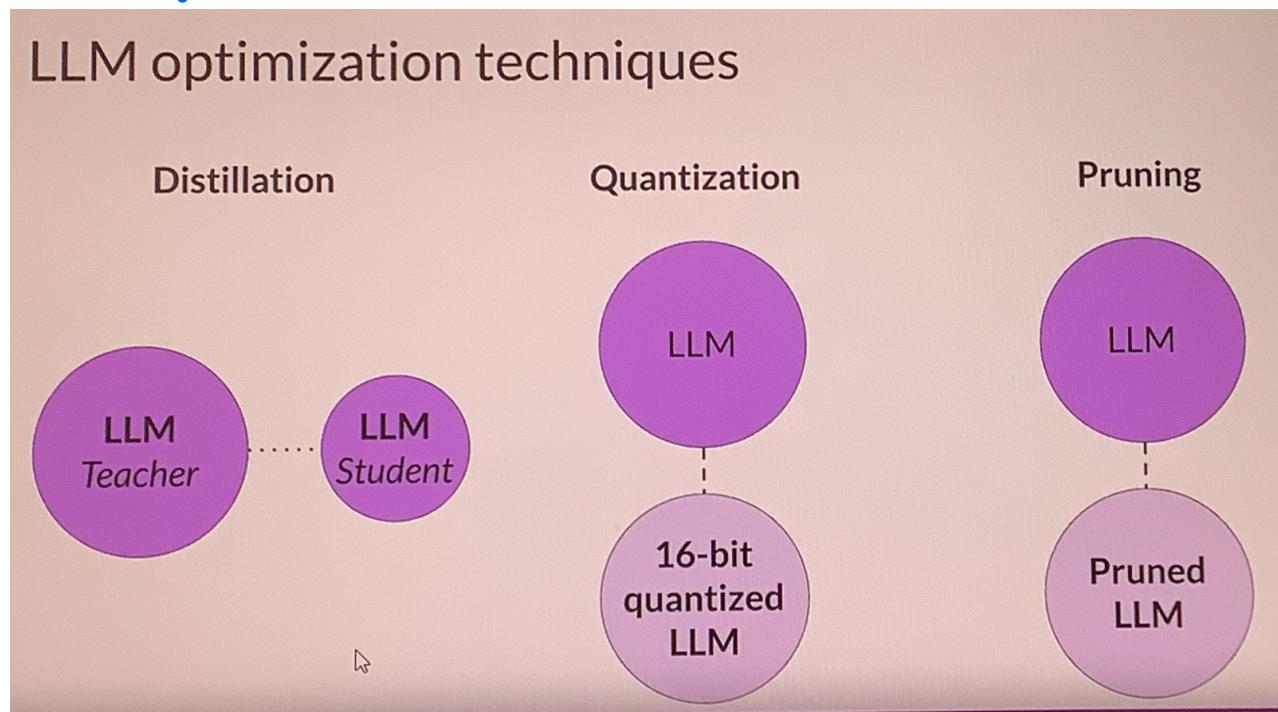


Optimization:

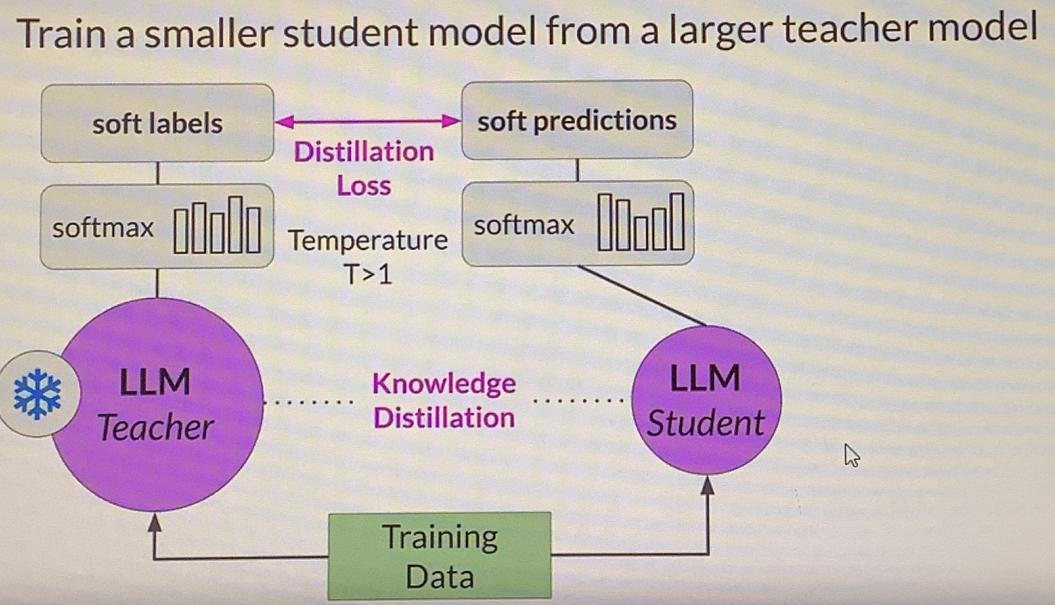


Model Optimization for Deployment

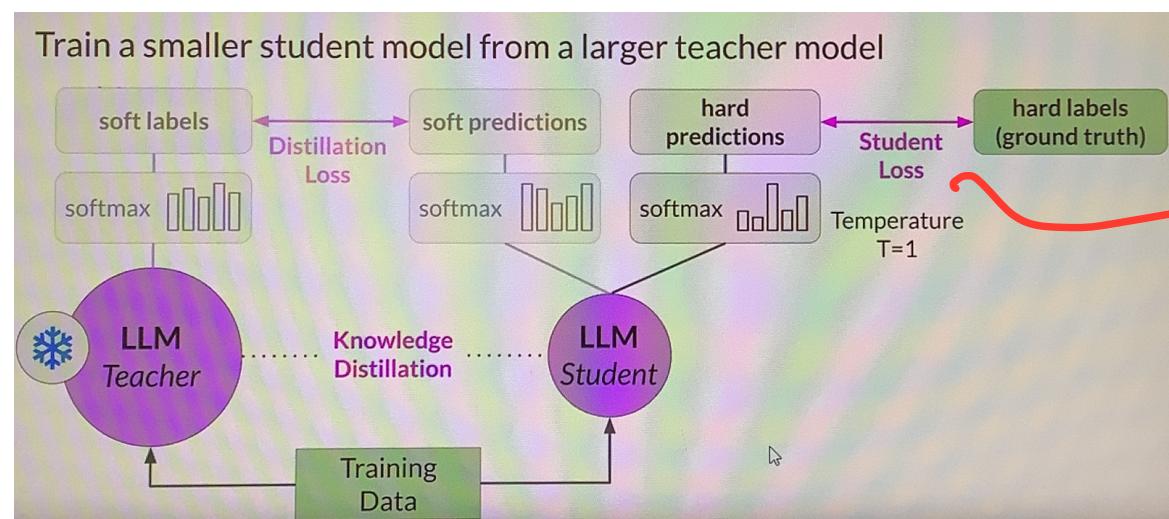
1. Distillation 2. Quantization
3. Pruning (Remove redundant)



Distillation:



In parallel we
also use dataset
to train 'LLM
Student' model.



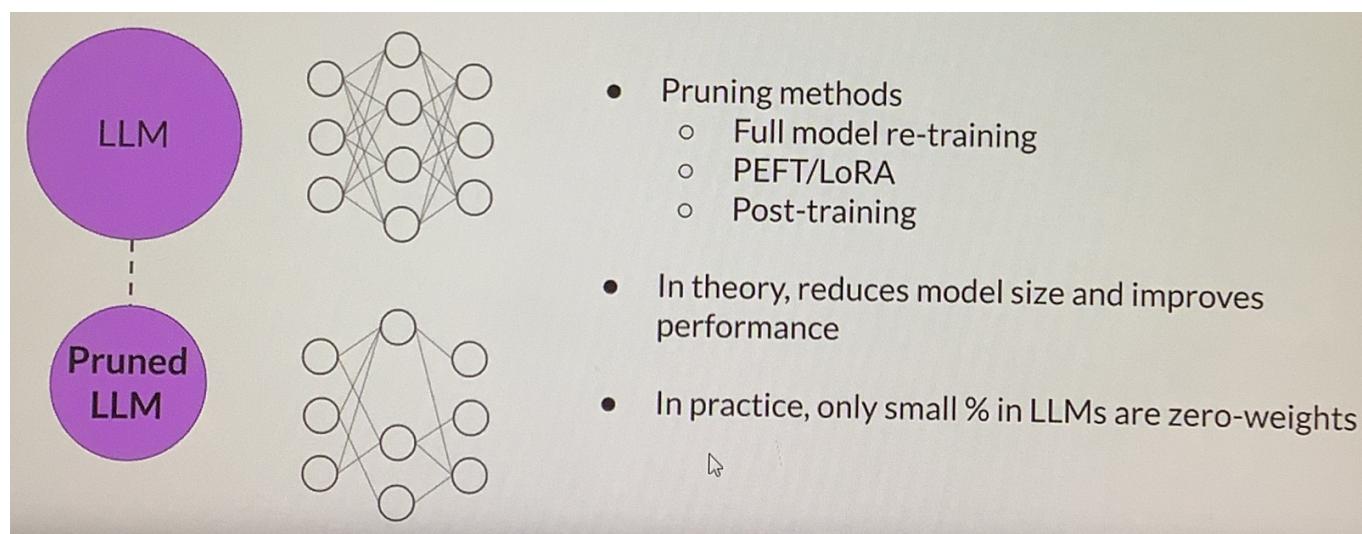
To calculate
'Student loss'
Combine of
Distillation loss
+ Student loss is used to update.

Student weight via 'Back Propagation'

⊗ In practice 'Distillation' is not as effective for Generative

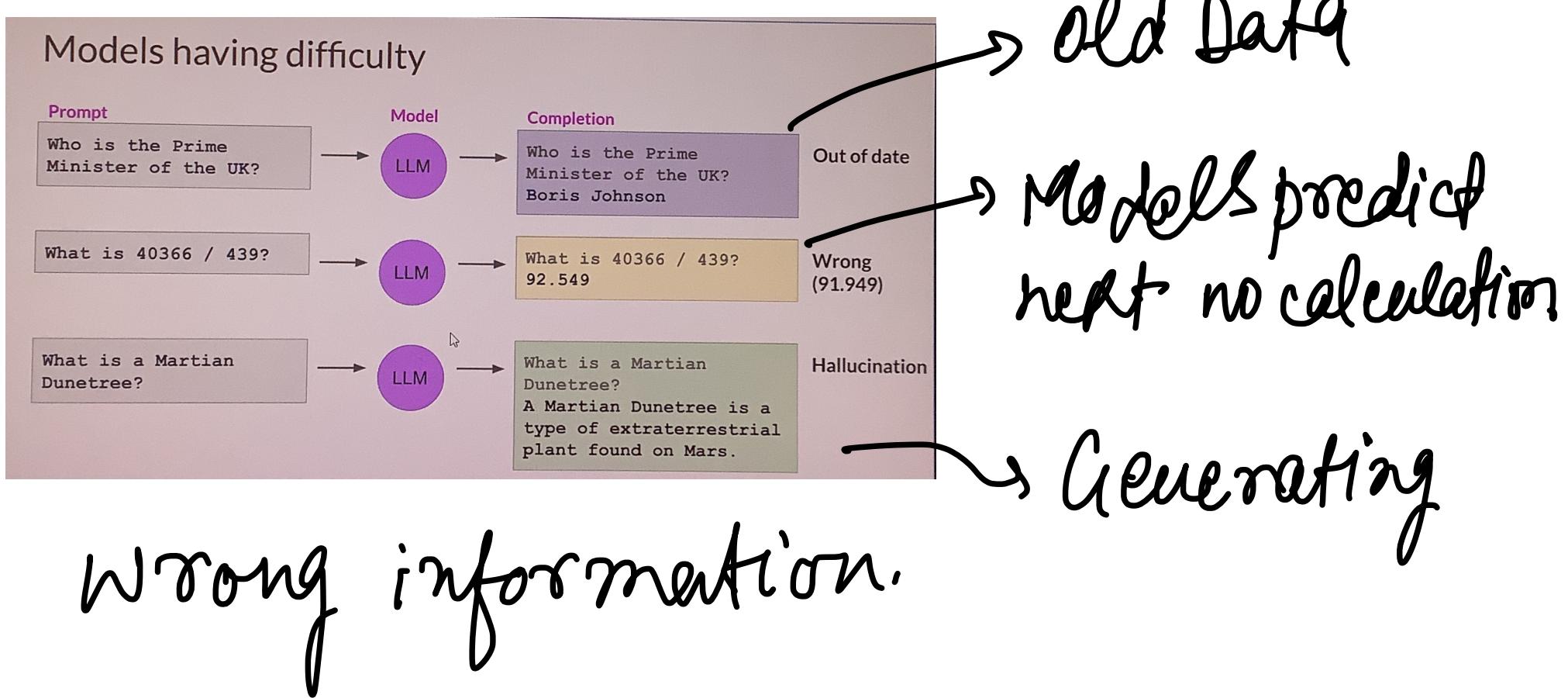
Decoder' model; It is more effective 'Generative Encoder' model. Such as 'BERT' that have lot of representation redundancy.

Pruning → Remove model weights with values closer to zero.

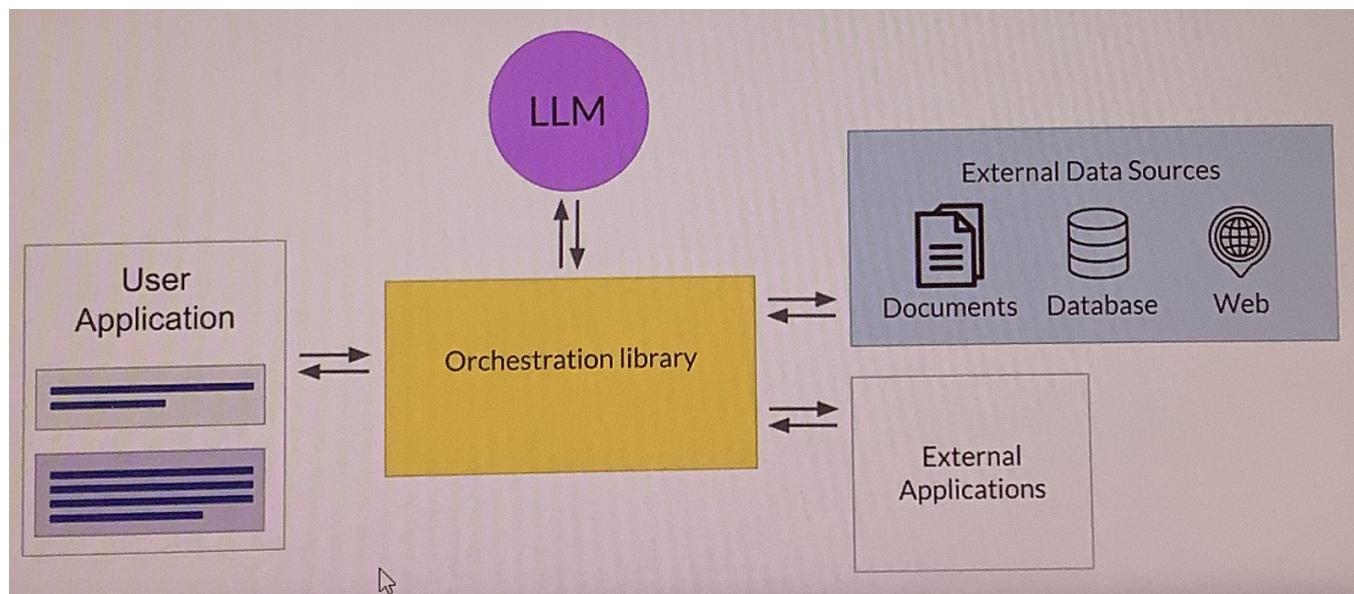


Time & Effort in the life cycle:-

	Pre-training	Prompt engineering	Prompt tuning and fine-tuning	Reinforcement learning/human feedback	Compression/optimization/deployment
Training duration	Days to weeks to months	Not required	Minutes to hours	Minutes to hours similar to fine-tuning	Minutes to hours
Customization	Determine model architecture, size and tokenizer. Choose vocabulary size and # of tokens for input/context Large amount of domain training data	No model weights Only prompt customization	Tune for specific tasks Add domain-specific data Update LLM model or adapter weights	Need separate reward model to align with human goals (helpful, honest, harmless) Update LLM model or adapter weights	Reduce model size through model pruning, weight quantization, distillation Smaller size, faster inference
Objective	Next-token prediction	Increase task performance	Increase task performance	Increase alignment with human preferences	Increase inference performance
Expertise	High	Low	Medium	Medium-High	Medium



LLM Powered Applications



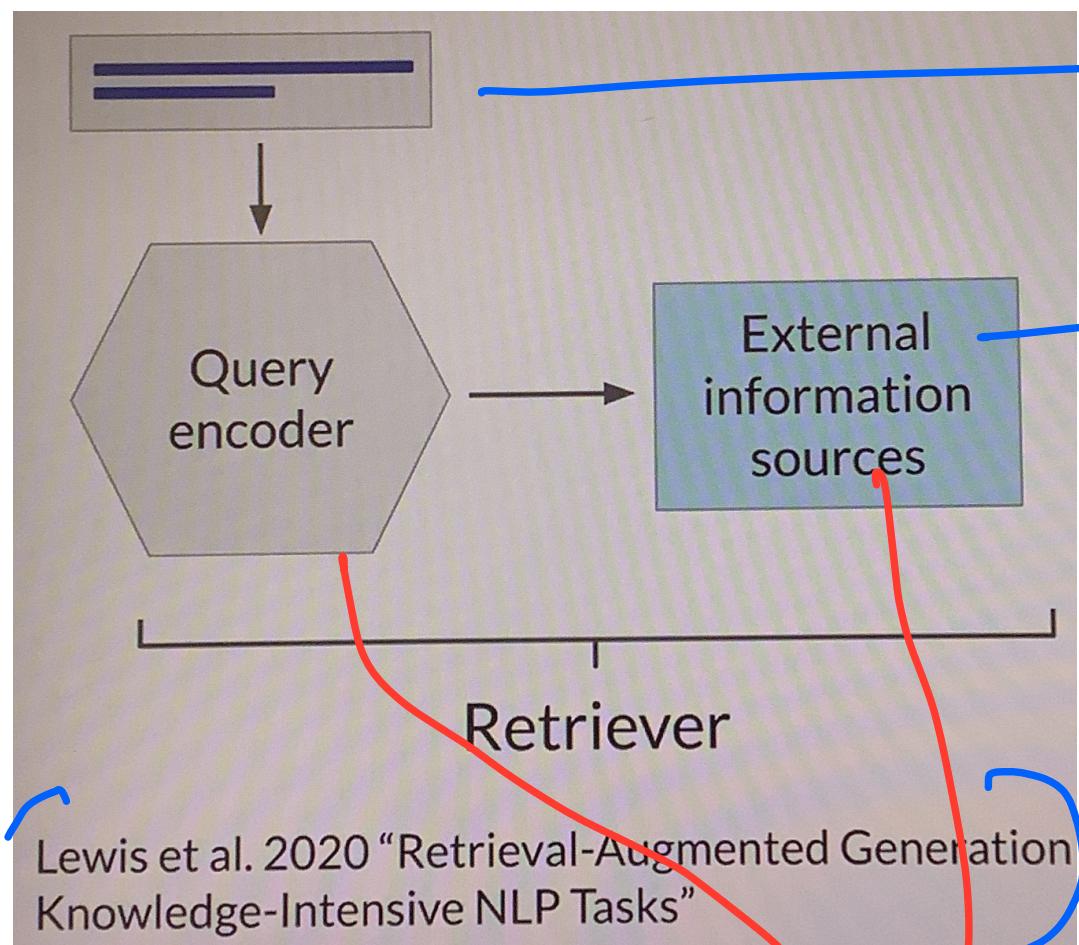
Our application
must manage
the passing

of user input to the LLM & return
the completions. This often
done by using 'Orchestration library'.

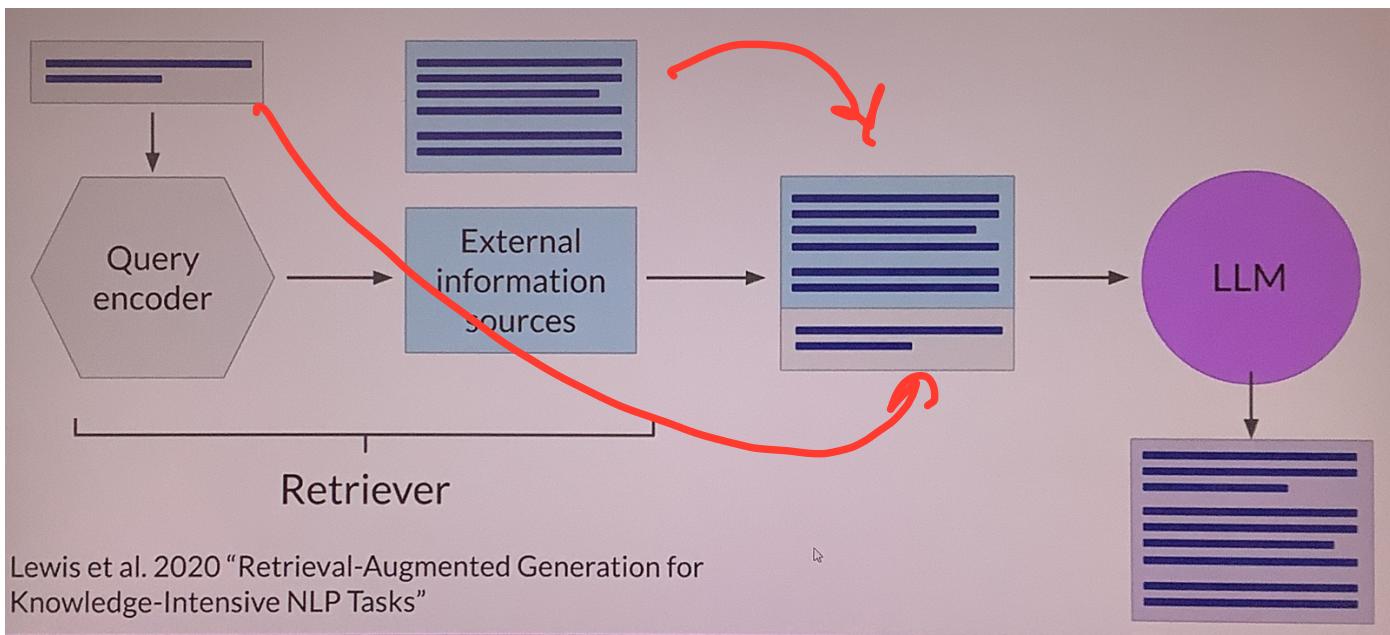
The layer can enable some powerful
technologies that augment & enhance
the performance. By providing

External data sources or connecting
with existing APIs of other
applications!

Retrieval-Augmented Generation (RAG)



Query
DB (It can be
in the form of
CSV, SQL etc.
but in the
facebook paper it was stored as
Vector DB. These two component
are train together find DOCUMENT
most RELEVANT for QUERY.

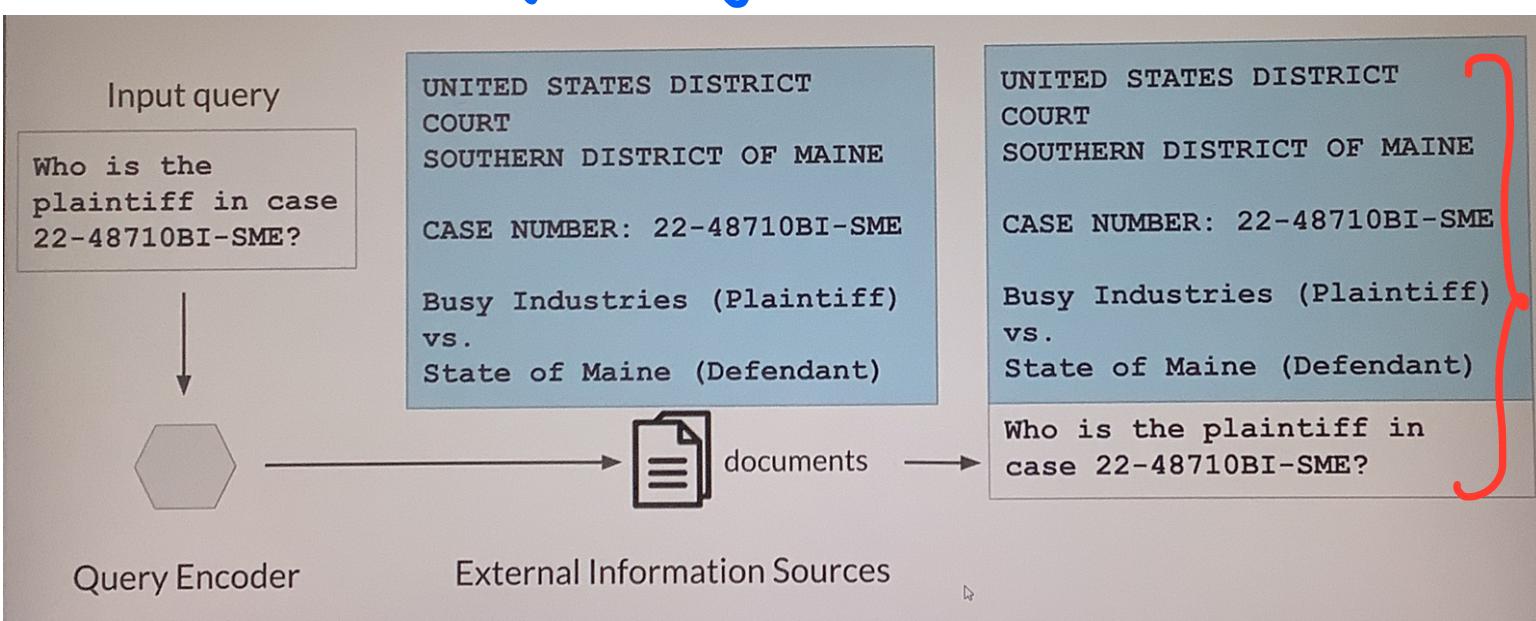


Query
output +

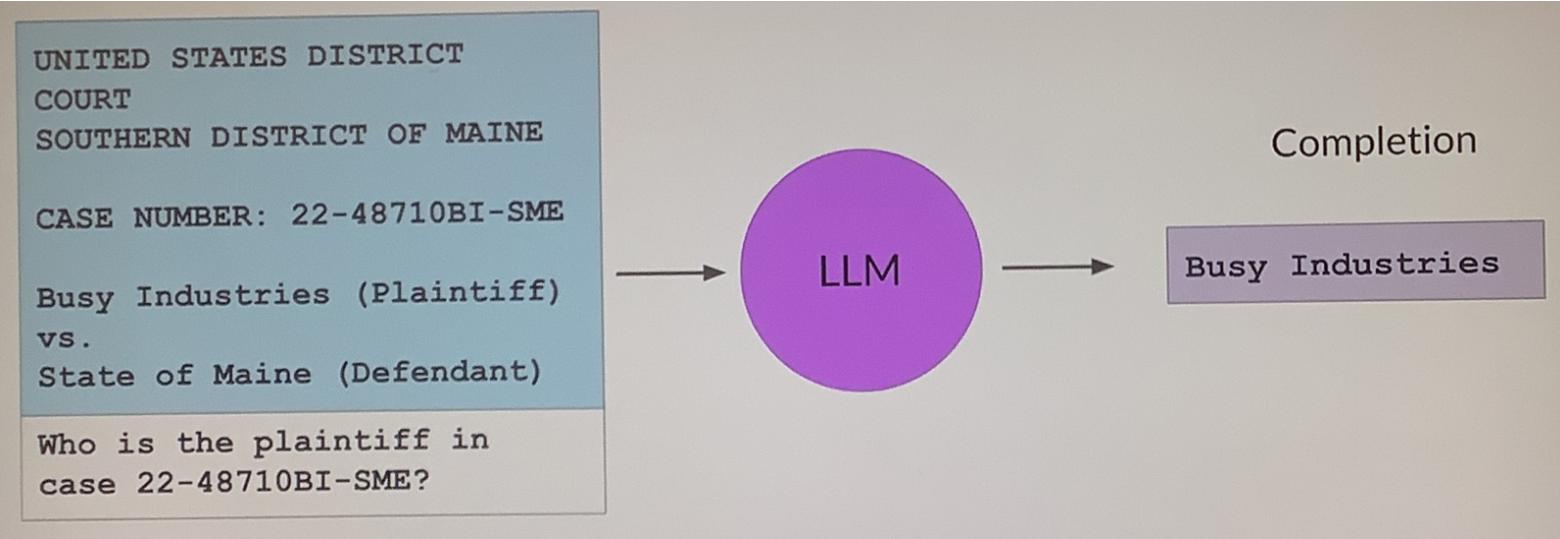
Query both are concatenated & passed as prompt into LLM.

let's Take Specific Example

Searching legal Documents:



Concat
Retrie-
ved
Document
+ Query.



Model then uses the context & query as prompt & generate desired output.

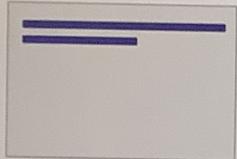
* RAG also helps to avoid the problem of MODEL Hallucination.

Data Preparation for Vector Store RAQ

Two considerations for using external data in RAG:

1. Data must fit inside context window

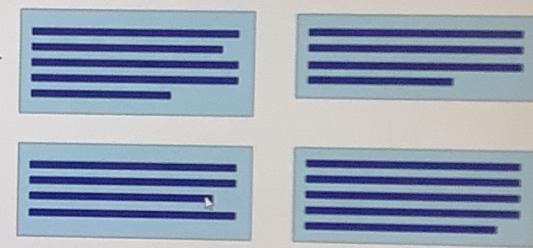
Prompt context limit
few 1000 tokens



Single document too
large to fit in window

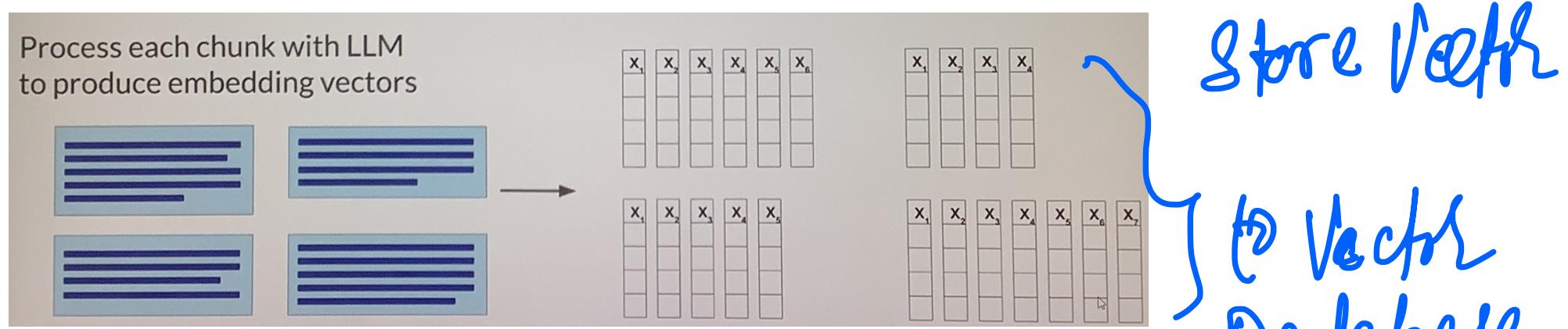


Split long sources into
short chunks



LANGChain
can handle
this kind of
work.

Q Data must be in format that allows its relevance to be assessed at inference time: Embedding Vectors.



to allow fast search.

Programming-aided Language (PAL)

models! -

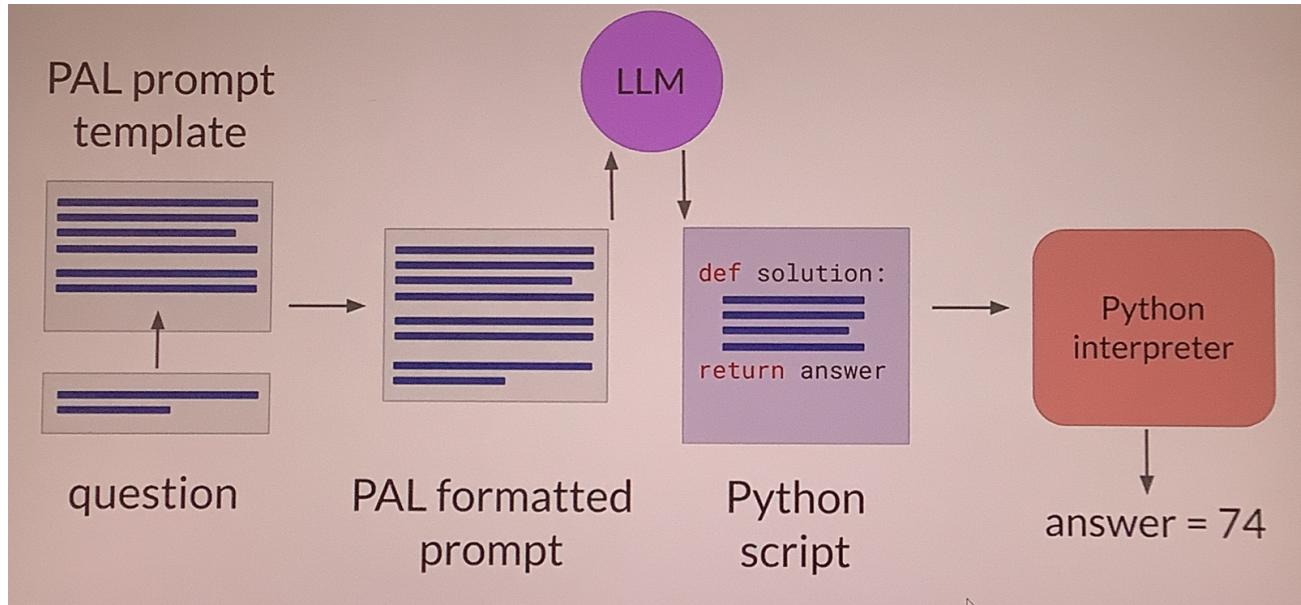
Prompt with one-shot example

Q: Roger has 5 tennis balls. He buys 2 more cans of tennis balls. Each can has 3 tennis balls. How many tennis balls does he have now?

Answer:

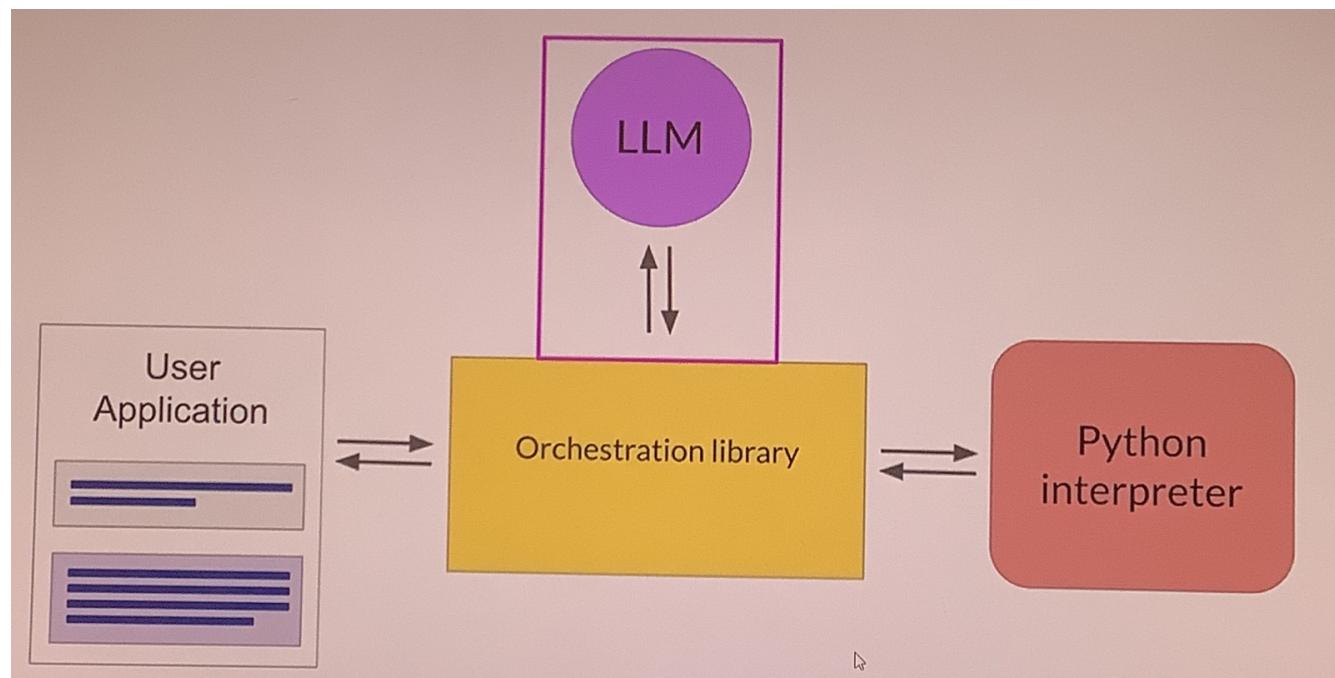
```
# Roger started with 5 tennis balls
tennis_balls = 5
# 2 cans of tennis balls each is
bought_balls = 2 * 3
# tennis balls. The answer is
answer = tennis_balls + bought_balls
```

Q: The bakers at the Beverly Hills Bakery baked 200 loaves of bread on Monday morning. They sold 93 loaves in the morning and 39 loaves in the afternoon. A grocery store returned 6 unsold loaves. How many loaves did they have left?



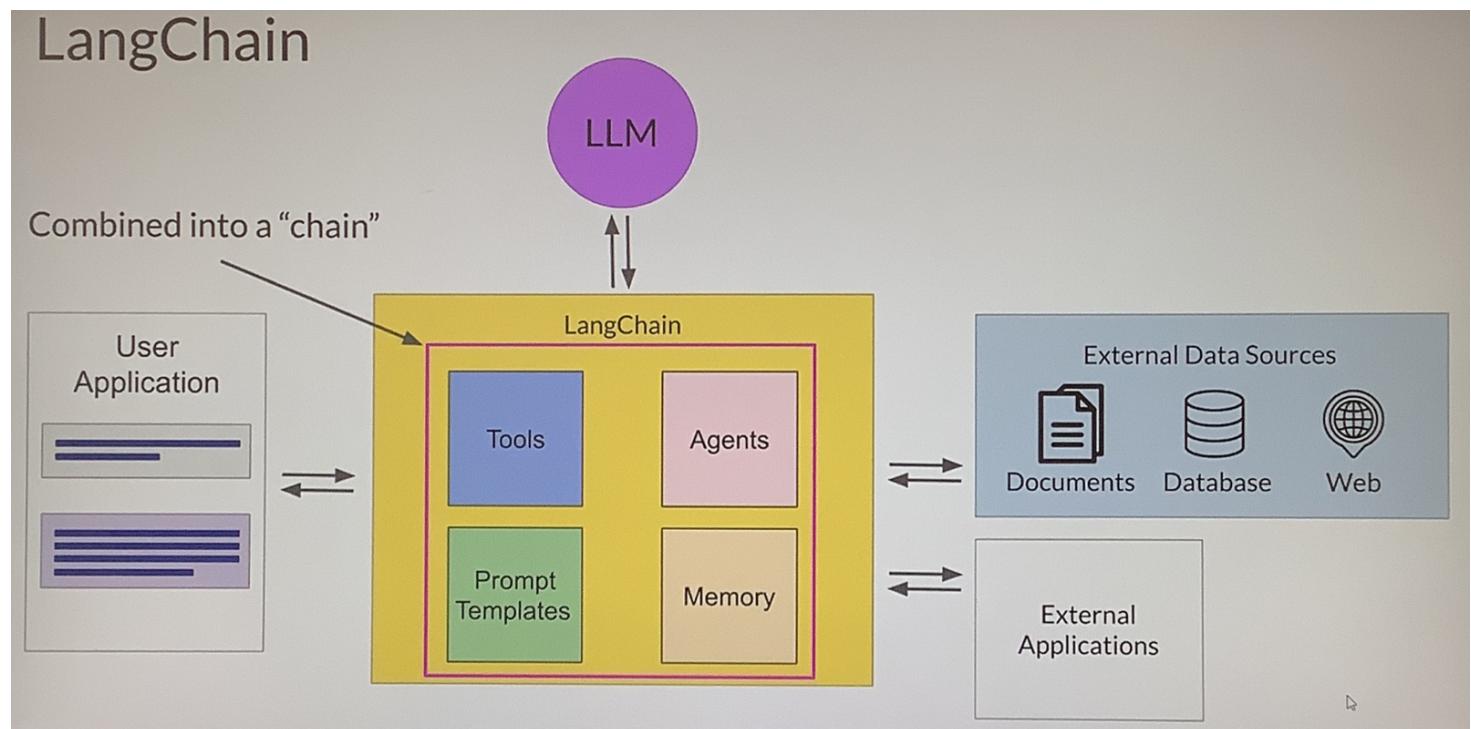
PAL Flow

PAL Architecture :-



ReACT: Combining Reasoning and Action in LLMs :-

LangChain:



LLM Powered Application Architectures :-

