

## Generative AI with LLM

(General purpose  
technology)

- useful for different application. (LLM)
- Generative AI is subset of Machine learning, learned their ability by finding statistical pattern in massive datasets (human)

### LLM

- GPT
- Bert
- Flan-T5

- LLAMA
- PALM
- BLOOM

(More parameters &  
more sophisticated tasks.)

Transformer Model → Neural Network that learns context & tracking relationship in sequential data.  
(foundation model)

### LLM use cases & tasks

(Next word prediction is basic task)

- Essay writer
- Summarize
- Translate

- Code AI
- Entity Extraction, word classification
- External API

Smaller Model can be fine tuned to perform well on specific focussed tasks.

LLM capability → Architecture powers them

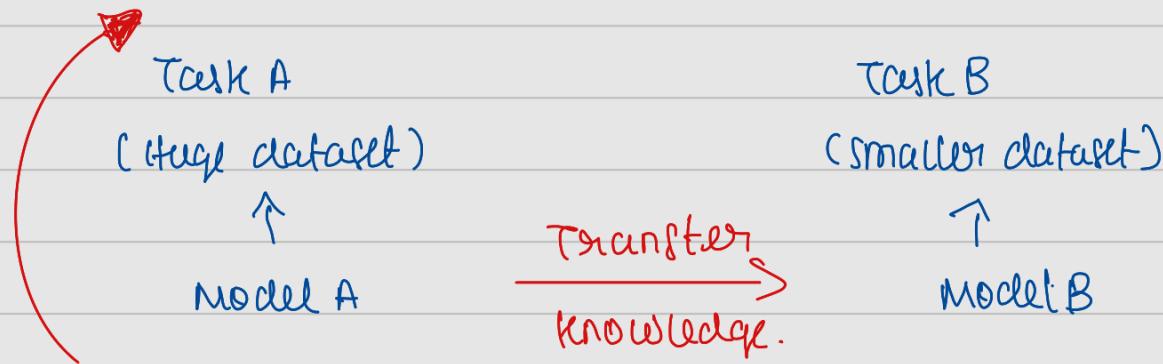
understanding language can be challenging, language is complex  
single word has multiple meaning. (homonyms)

Previous generation of language model made up of 'RNN'

### Self-supervised learning:

Type of training in which the objective is automatically computed from inputs of the Model. Means that humans are not needed to label the data.

### Transfer Learning



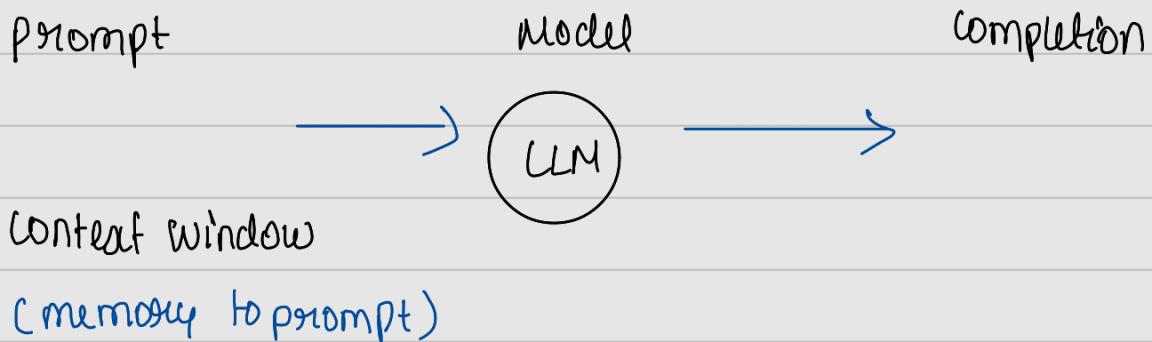
Act of initializing a model with another model's weight

Training from scratch requires more data and more compute to achieve comparable results.

Transformer Model are built with special layer called Attention layer

### Fine-tuning

- \* Full fine tuning
- \* parameter efficient fine tuning (PEFT)



using  
Inference : Model to generate text.

Generating feat: RNN

RNN is limited by the amount of compute & memory.

then, transformer Architecture had arrived  $\rightarrow$  'Attention is all you need'  
(Attention to input meaning)

"Attention is All you need"  $\rightarrow$  by Google in 2017.

$\rightarrow$  Neural Network Architecture replaced the RNN (CNN) with attention based mechanism.

layers.

$\rightarrow$  encoder & decoder, multi-head self attention mechanism & feedforward NN

$\rightarrow$  it uses self attention to capture input sequence. capture long term dependencies

Transformer

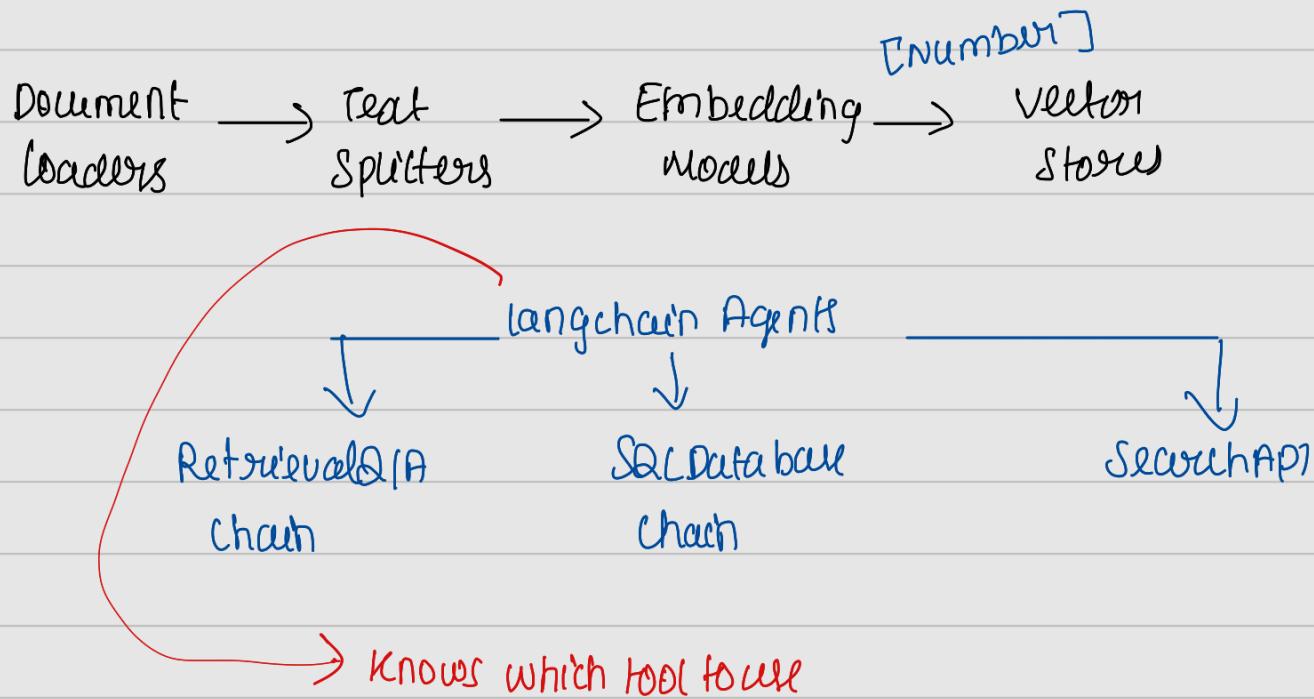
$\rightarrow$  learn the relevance & context in all the words in the sentences.  
and apply attention weights. (attention map  $\rightarrow$  each word with every word)

Model process the Number not the text.

## word Embedding.

↳ technique where individual words are represented as real valued vectors in a lower-dimensional space and captures inter word semantic.

## langchain



Output parsers → response from LLM follow certain format.

### Generative Models ↳

Learning patterns and distribution in unstructured data. So,

Reinforcement learning through human feedback is also applied.

LLM → transformer, a type of Neural Network architecture.

## language model

↳ A language Model is a Model which understands language more precisely how words occur together in Natural language , A language Model is used to predict what word comes next.

temperature = 0  $\Rightarrow$  argmax sampling.

lower temperature means less randomness, vice versa.

top-p (Nucleus Sampling)

- ↳ It controls how many random results should be considered for completion as per the temperature

### Generative AI project life-cycle

(own)

- ① foundation Model | pre-training Model (decision)
- ② fine-tuning - custom data.

Self Supervised Learning to optimize these Models.

Issue with RNN for Generative part is, understanding language can be challenging (homonyms)

then transformer architecture has arrived

- scale efficiently
- parallel process
- Attenuation to input meaning.

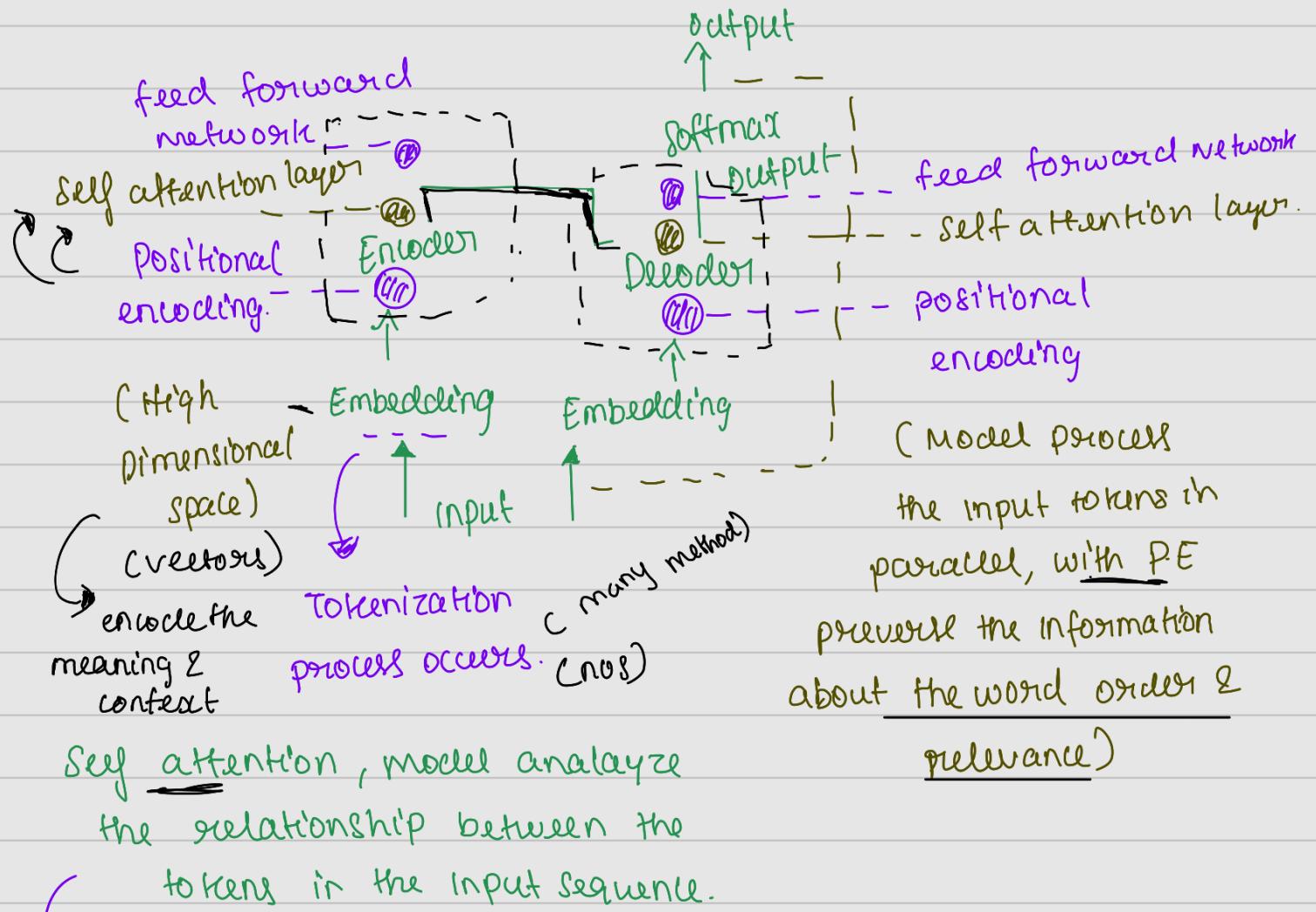
increases the performance of Natural Language task.

(self-attention)

power of transformer architecture, its ability to learn the relevance and context of all of the words in a sentence, every other word in a sentence. (Attention map)

Apply attention weights to relationship for relevance of each word to other word.

Transformer



Self attention, model analyze the relationship between the tokens in the input sequence.

(layer's weights are stored.)

↓  
 (to capture the contextual dependencies between the words  
 (Reflects the importance of each word in the input sequence))

transformer has Multi-headed Self attention (12-100)

means that multiple sets of self-attention weights / heads are learned in parallel independently of each other.

- 1. Relationship 2. Activity 3. other properties

Each self-attention head will learn different aspects of language, weights of each head is randomly initialized.

Output is processed through feed forward network, proportional to probability score for every token then passed to softmax layer → P(CE) for each word. (higher)

## Language Model

what should be the next words == word with highest probability.

Given a sequence of words

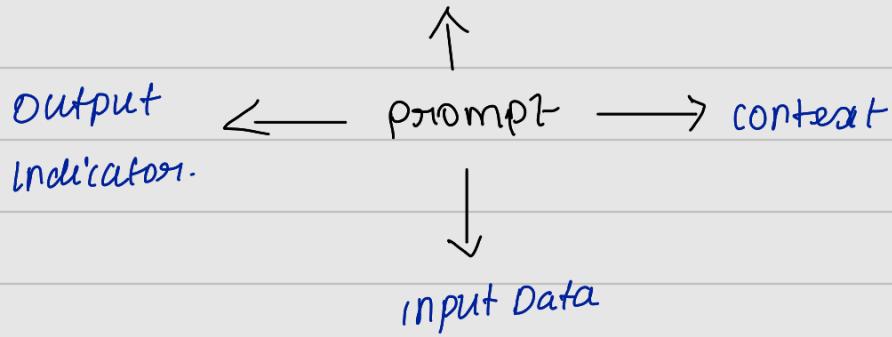
$x_1, x_2, \dots, x_t$

compute the probability of the next word

$$P(x_{t+1} | x_1, x_2, \dots, x_t)$$

LM  $\rightarrow$  trained on huge amount of data.

Instructions



zero-shot prompt (popular)  $\rightarrow$  accuracy (imitation)

prompt does not contain any explicit instructions or example for the model to follow. Instead it relies on the model's ability to understand & interpret natural language.  
(without specific training)

## Embeddings

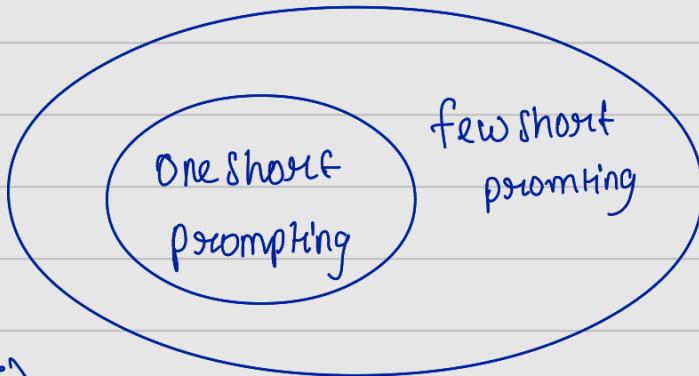
classic technique in NLP. The idea is to create a vector space from the text such that the distance between the vectors in the space have a certain meaning.

vector  $\rightarrow$  simply a sequence of numbers. ('represent complex Model')

vector database is a database that's saving those embeddings.  
(pinecone)

### few short prompt

Model is given a few examples of the result you want.  
useful where there is limited data available for a given task.



(COT)

### chain of thought prompting (google)

↳ improve reasoning abilities of LLM, it decomposes the multi-step problems into intermediate steps, allowing to solve a complex problem.

→ prompt engineering technique.

↳ if's how human solves the problem.  
(problem solving)

(major step in prompt engineering)

② few short  
chain of thought  
(similar to few  
shot learning)

① zero shot chain of  
thought

↳ let think step by  
step (prompt)

### Embedding in transformer (high dimensional space)

↳ each token represented as a vector & occupies  
a location. (wordvec)

## fine tuning

fine tuning is taking these general purpose models like GPT-3 & specializing them into something chatGPT

lets the model to learn the data, rather than just get access to it.

Reduces hallucinations & biases is similar to the model's earlier training.

## generating text with transformer

The attention mechanism is used in both the encoder & decoder layers of the transformer mechanism.

The data leaves the encoder is, a deep representation of the structure and meaning of the input sequences. This representation is inserted into the middle of the decoder to influence decoder's self-attention mechanism.

Encoder ↗ Encodes inputs ('prompt')  
with contextual understanding  
& produces one vector per input token.

Decoder ↗ Accepts input tokens and generates new token.

]  
] closes this  
in a loop

Encoder only → Sentiment analysis,  
models BERT

Encoder-decoder  $\rightarrow$  sequence to sequence tasks,  
model translation.  
 $\text{len}(\text{input}) \neq \text{len}(\text{output})$   
 $\rightarrow$  Google Bard

$\rightarrow$  Decoder only  $\rightarrow$  GPT, BLOOM, Llama (models)  
model

## Prompting & prompt Engineering

The develop & improve the prompt is known as  
prompt engineering.

Providing examples inside the context window is called  
in-context learning (ICL)  $\rightarrow$  there's a limit to context window.

$\rightarrow$  zero shot inferences.

$\rightarrow$  few shot inferences.

$\rightarrow$  one-shot inferences

$\rightarrow$  engineering your prompts to encourage the model to  
learn by examples.

Large models are good at zero shot inferences.

## Generative configuration (inference parameters)

\* max-new tokens  $\rightarrow$  number of tokens that model will  
generate.

The output from the transformers' Softmax layer is a probability distribution across the entire dictionary of words.

Most LMs by default will operate with so-called *greedy decoding*,

(greedy):

The word [token with the highest probability is selected.

So, bring random sampling.

Random (weighted) Sampling, Select a token using a random-weighted strategy across the probabilities of all tokens.

---

top-K & top-P Sampling techniques.

↳ limit the random sampling & increase the output to be sensible

top-K → select an output from the 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  $p(E)$  and with a cumulative  $p(E) \leq P$

Specify no. of tokens to randomly choose from.

↳ specify the total  $p(F)$  want model to choose from.

temperature

↳ higher temperature, higher randomness.

applied in the final Softmax layer.

temperature  $< 1$

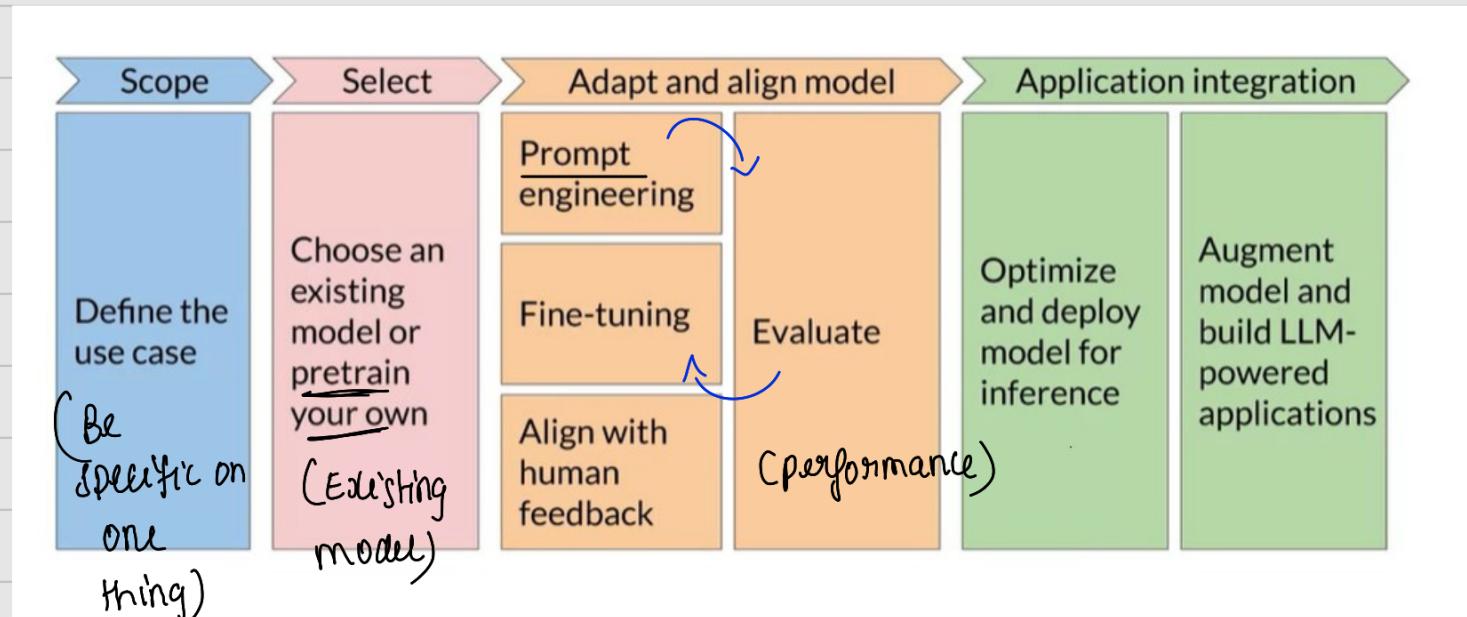
strongly peaked probability distribution

temperature  $> 1$

broaden, flatter probability distribution.

Generative AI project lifecycle

fine tuning  $\hookrightarrow$  supervised learning.



pre-training large language model

consideration for choosing a Model.

→ foundation model  
(pretrained LLM)

→ Train your own Model  
(Custom LLM)

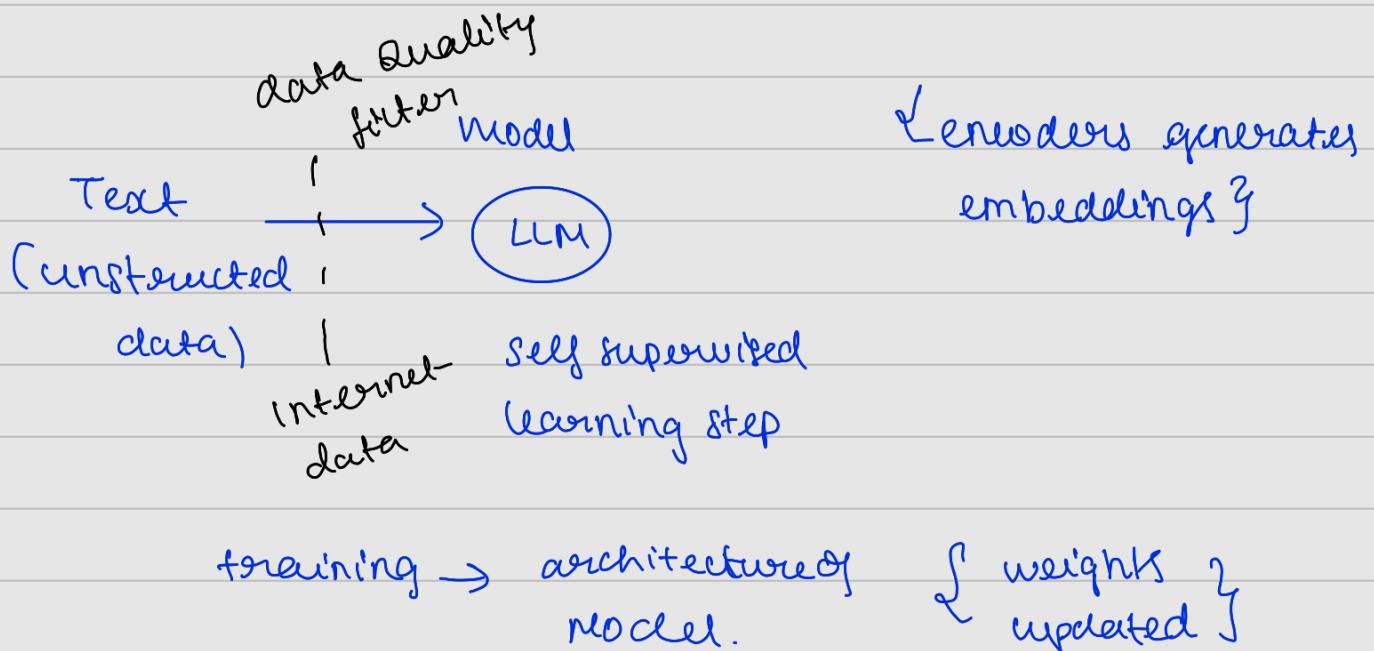
sticking face  $\rightarrow$  Have Model cards  
(Hub) (imp details & best use cases for each model)

BloombergGPT  $\rightarrow$  Decoder only language model, pretrained using an extensive financial datasets.

The goal of pre-training is to teach the model the general pattern and structure of language, allowing it capture the relationship b/w words.

The goal of fine-tuning is to adjust the model's parameters to make it better at performing the desired task.

LLM encode a deep statistical representation of language developed during model's pretraining phase.



pre-training also requires a large amount of compute & use of GPUs

Encoder only models are also called Autoencoding models, they are pre-trained using masked language modelling.

masked language modelling (MLM)

The teacher teaches the student

(mask) (randomly Masked)

training objective is to predict the mask token in order to reconstruct the original sentence.

Objective: Reconstruct text ('denoising')

The teacher teaches the student.

→ ← Bidirectional context

model has understanding ↘

of full context of a token.

Good use cases:

- sentiment analysis
- named entity recognition
- word classification

[Good for tasks that require understanding of the input, such as sentence classification & named entity recognition]

Example models:

- BERT
- ROBERTA

Decoder only model | Autoregressive models , which are pre-trained using Causal Language Modelling (CLM)

Here training objective is to predict the next tokens based on the previous sequences of tokens.

↳ Fully language modelling.

Objective: Predict next token.

The ? (CLM) ↳ no knowledge of the end of sentences  
↳ The teacher

↳ iterate one by one tokens for next word, context is unidirectional.

Good use case

Example models

→ Text generation

→ GPT

→ Other emergent behaviour

→ BLOOM.

↳ depends on model size

[Good for generative tasks such as text generation]

Encoder - decoder model | Sequence to Sequence models , good for generative tasks that require an input, such as translation or summarization.

Span  
corruption

Original  $\rightarrow$  Encoder-Decoder  
text LLM

Span corruption

The teacher  $\xrightarrow{\text{mask}}$   $\text{mask} \rightarrow$  student  
The teacher  $\xrightarrow{x}$  student  
 $\hookrightarrow$  sentinel token

Objective: Reconstruct Span

$\xrightarrow{x}$  teaches the

[Sequence-to-sequence models are best suited for tasks revolving around generating new sentences depending on a given input such as summarization]

Models:

$\xrightarrow{\text{TS}}$

$\xrightarrow{\text{BART}}$

Growth of Models powered by;

- $\rightarrow$  Introduction of transformer
- $\rightarrow$  Access to massive datasets
- $\rightarrow$  more powerful compute resources.

computational challenges of  
training LLM

$\rightarrow$  CUDA out of memory

(Compute Unified Device Architecture)

PyTorch | tensorflow  
use this.

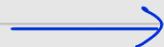
(S tools developed for  
Nvidia GPUs

Approximate GPU RAM needed to store 1B parameters.

1 parameter  $\rightarrow$  4 bytes.

1B para  $\rightarrow$   $4 \times 10^9$  bytes  $\rightarrow$  4 GiB

memory needed to  
store model



4GiB @ 32-bit  
full precision

Single  
processor.

memory needed to  
train model.

24GiB @ 32-bit  
full precision.

① Reduce the memory,

Reduce the memory required for training is called  
quantization

↳ reduce the memory required to  
store weight by reducing their  
precision from 32-bit floating to 16-bit  
floating point numbers. / 8-bit  
the model weights.

As model sizes get larger, you will need to split your  
model across multiple GPUs for training.

Approximate GPU RAM needed to train 1B parameters. )

## Additional GPU RAM needed to train 1B parameters

Bytes per parameter	
Model Parameters (Weights)	4 bytes per parameter
Adam optimizer (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

Quantization: Summary.

	Bits	Exponent	Fraction	Memory needed to store one value
FP32	32	8	23	4 bytes
FP16	16	5	10	2 bytes
<b>BFLOAT16</b>	16	8	7	2 bytes
INT8	8	-/-	7	1 byte



- Reduce required memory to store and train models
- Projects original 32-bit floating point numbers into lower precision spaces
- Quantization-aware training (QAT) learns the quantization scaling factors during training
- BFLOAT16 is a popular choice

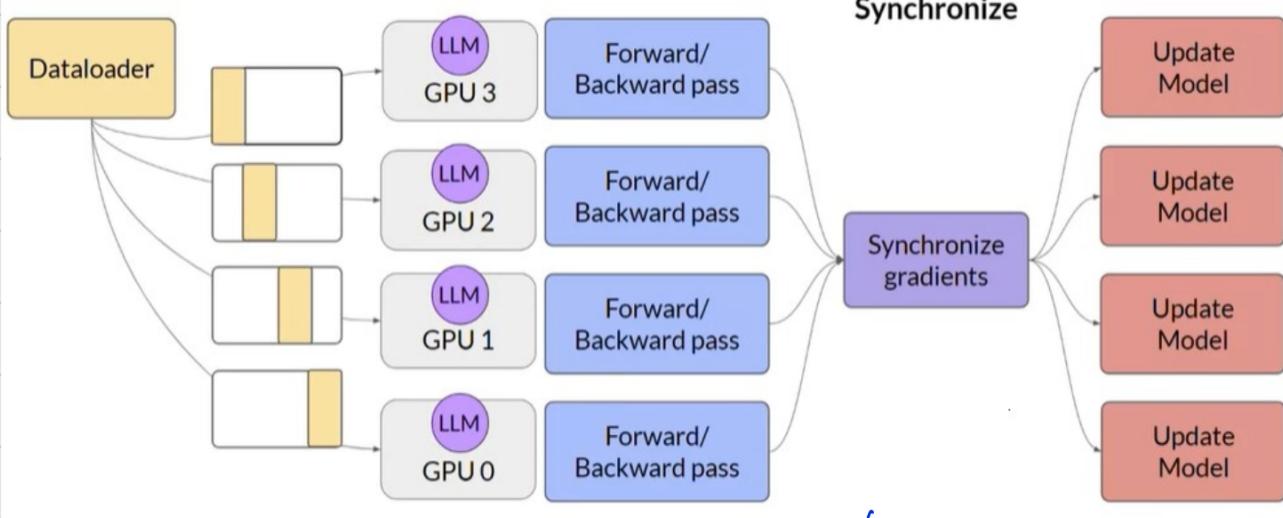
By quantization, we can reduce memory consumption required to store the model parameters.

Efficient multi-GPU compute Strategies

This implementation allows parallel computations across all GPUs that results in faster training.

# Distributed Data Parallel (DDP)

(dataset is distributed)



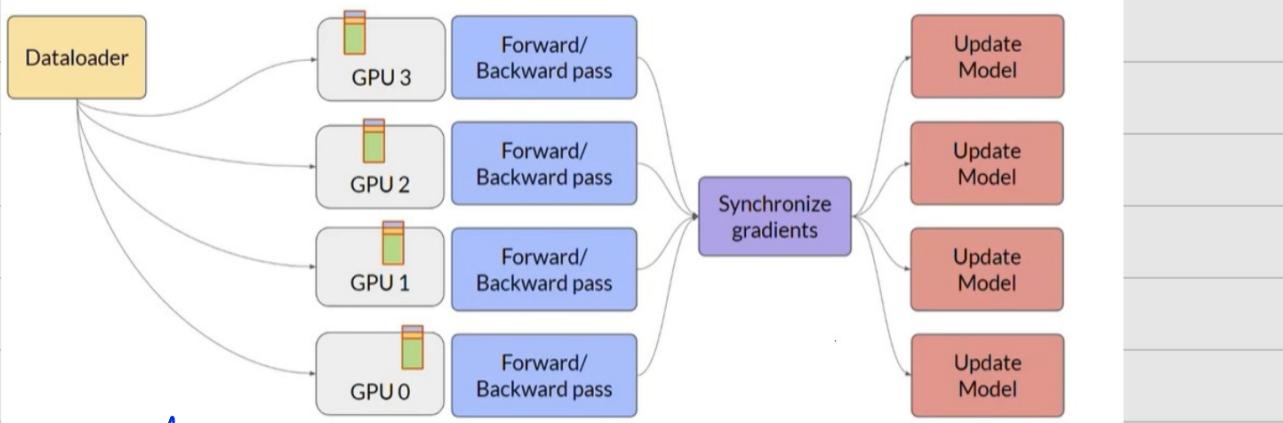
↳ requires model weights  
and additional parameters.

If model is too big for  
this, look for  
model sharding technique.

fully sharded Data parallel (FSDP), motivated by the  
'zero' paper - zero data overlap between GPUs

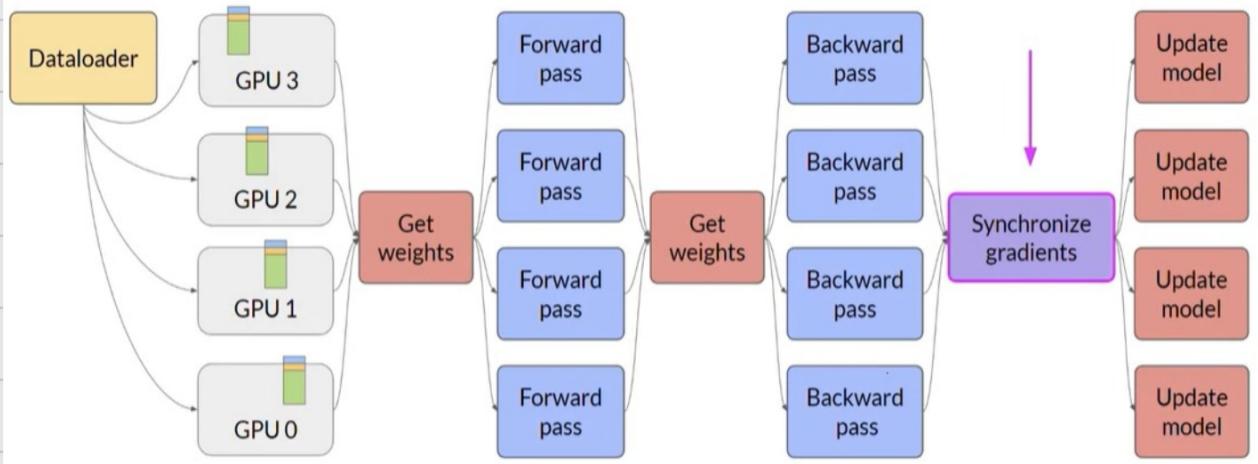
Zero → zero redundancy optimizer. (reducing the  
memory by distributing [sharding])

## Fully Sharded Data Parallel (FSDP)



↳ also distribute the model parameter, gradients &  
optimizes the states across the GPU nodes.  
using GPU

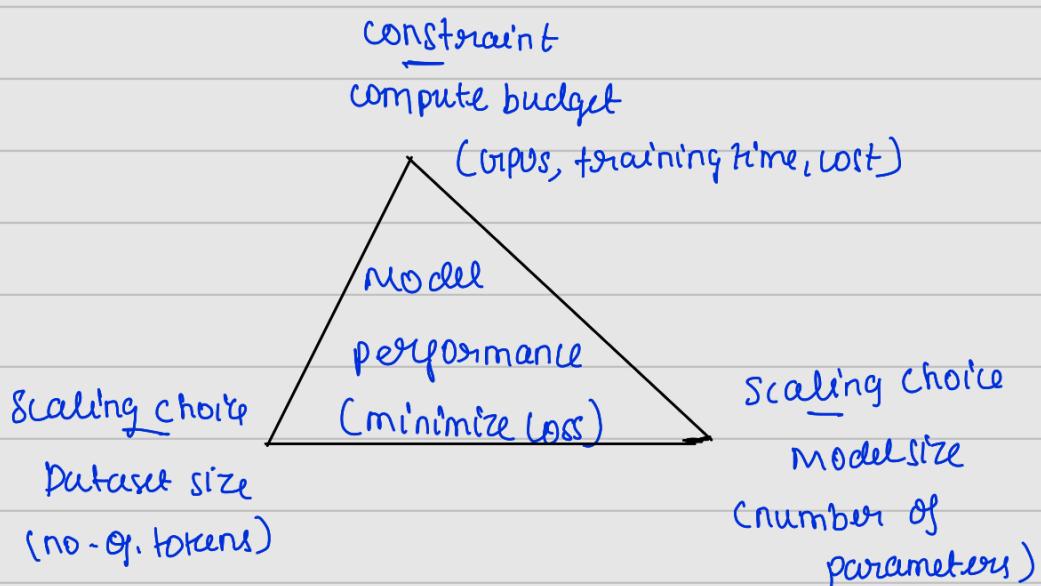
# Fully Sharded Data Parallel (FSDP)



- Helps to reduce overall GPU memory utilization
- Supports offloading to CPU if needed.
- Configure level of sharding via sharding factor.

## Scaling laws and compute-optimal models

Pre-training  
a model  $\Rightarrow$  maximize model  
performance.



'petaflop/s-day' = floating point operations performed at rate of 1 petaflop per second for one day.

- increase the compute budget to achieve better model performance. (Hardware, project timeline, financial budget)
- dataset  $\uparrow$  performance  $\uparrow$  } model size  $\uparrow$  performance  $\uparrow$

### compute optimal model

- very large models may be over-parameterized and under-trained
- smaller models trained on more data could perform as well as large models.

### Chinchilla scaling laws for model and dataset size

Model	# of parameters	Compute-optimal* # of tokens (~20x)	Actual # tokens
Chinchilla	70B	~1.4T	1.4T
LLaMA-65B	65B	~1.3T	1.4T
GPT-3	175B	~3.5T	300B
OPT-175B	175B	~3.5T	180B
BLOOM	176B	~3.5T	350B

Compute optimal training datasize is ~20x number of parameters

Sources: Hoffmann et al. 2022, "Training Compute-Optimal Large Language Models"  
Touvron et al. 2023, "LLaMA: Open and Efficient Foundation Language Models"

\* assuming models are trained to be compute-optimal per Chinchilla paper

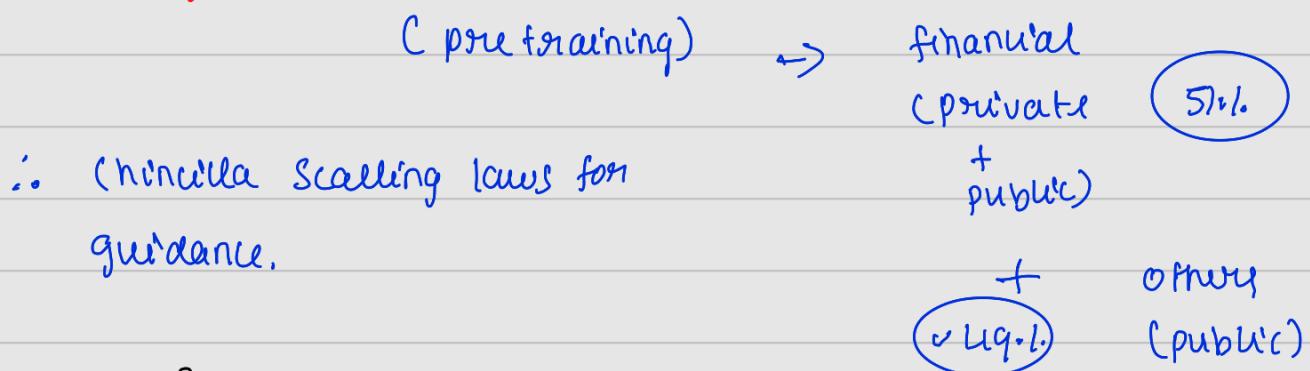
### pre-training for domain adaptation

words are rarely used, which means that they are unlikely to have appeared widely in the training text of existing LLMs.

↳ Models have difficulty in understanding them, using correctly.

Pre-training your model from scratch will result in better models for highly specialized domains like law, medicine, finance.

Bloomberg L.P. → domain adaptation for finance.



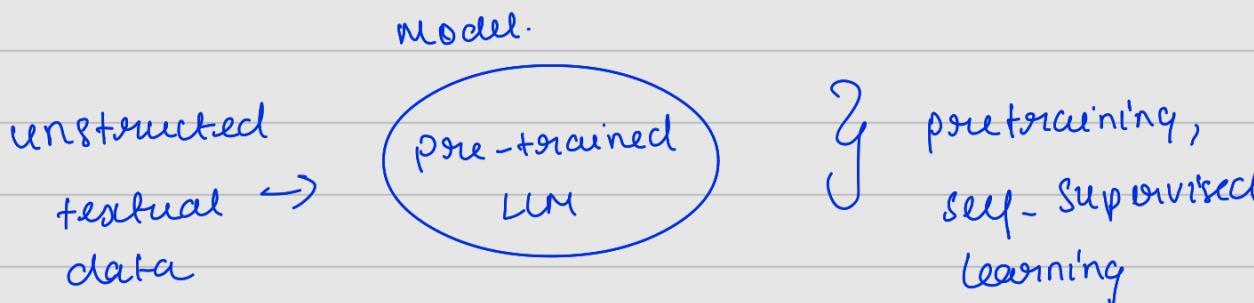
## Week-2

### instruction fine-tuning

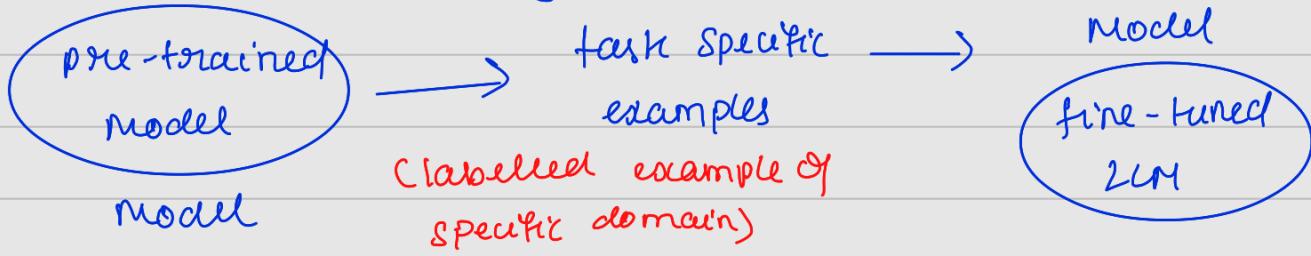
fine-tuning an LLM with instruction prompts.

- In-context learning may not work for smaller models (LLM)
- Examples take up space in the context window. So, instead try fine-tuning the model.

### LLM pretraining



LLM fine-tuning, supervised learning.



( prompt - completion pair. )  
examples  
Improved performance

Output of LLM is a probability distribution across tokens.

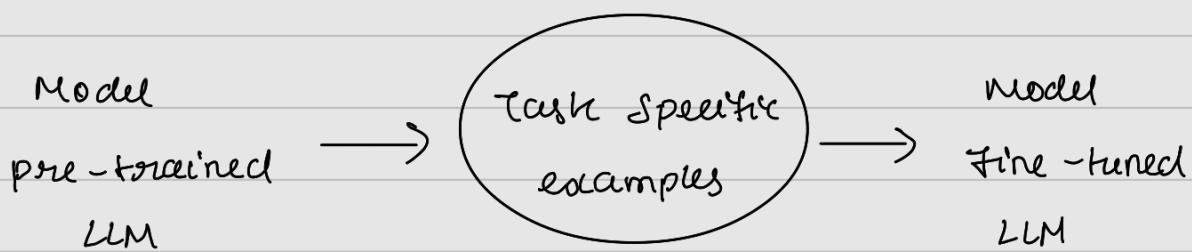


The fine tuning process results in a new version of the base model is Instruction model. (common way)

Fine-tuning with instruction prompts is common way, fine-tuning always means Instruction fine-tuning

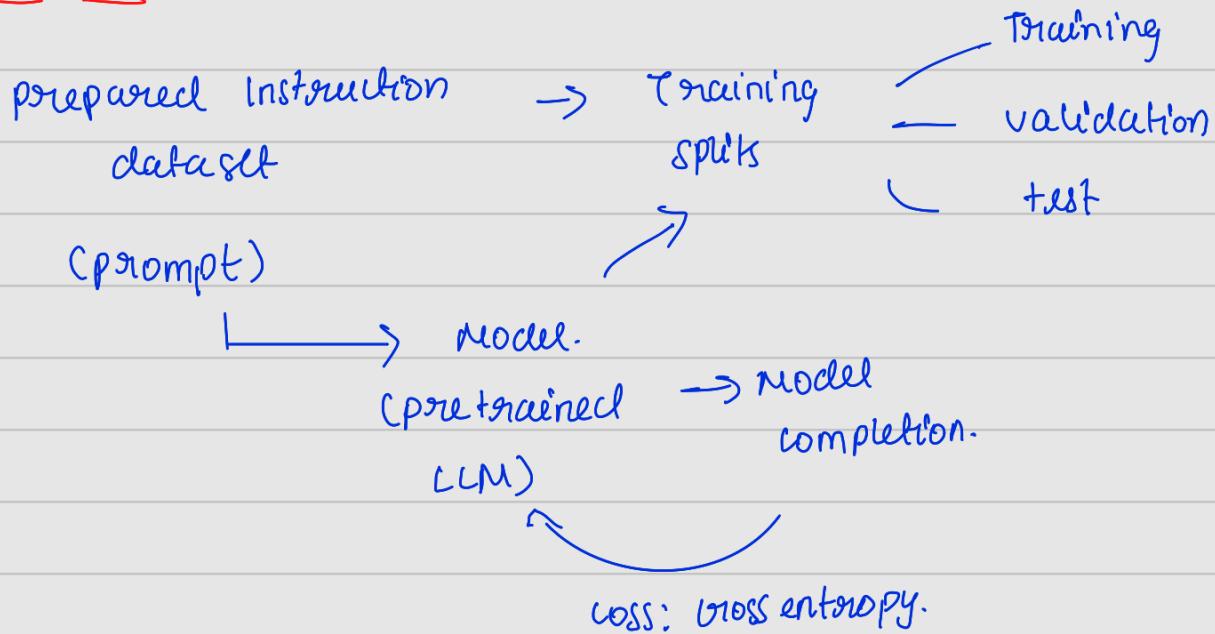
Instruction fine-tuning trains the model using examples that demonstrate how it should respond to a specific instruction.

Juul-fine-tuning updates all parameters.



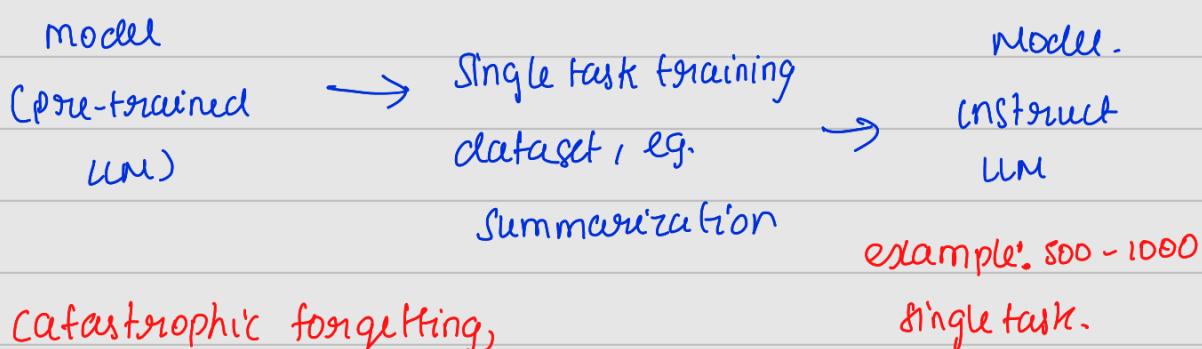
+ training data → formatted as instruction.

## CLM-fine-tuning



fine-tuning on a single task

parameter efficient }  $\rightarrow$  Set of techniques that preserves  
fine-tuning } the weights of the original  
LM (most of pre-trained  
weights are left unchanged)



## Catastrophic forgetting,

sample: 500 - 1000

single task.

↳ full fine tuning can significantly increase the performance of a model on a specific task.  
(modifies original LLM weights)

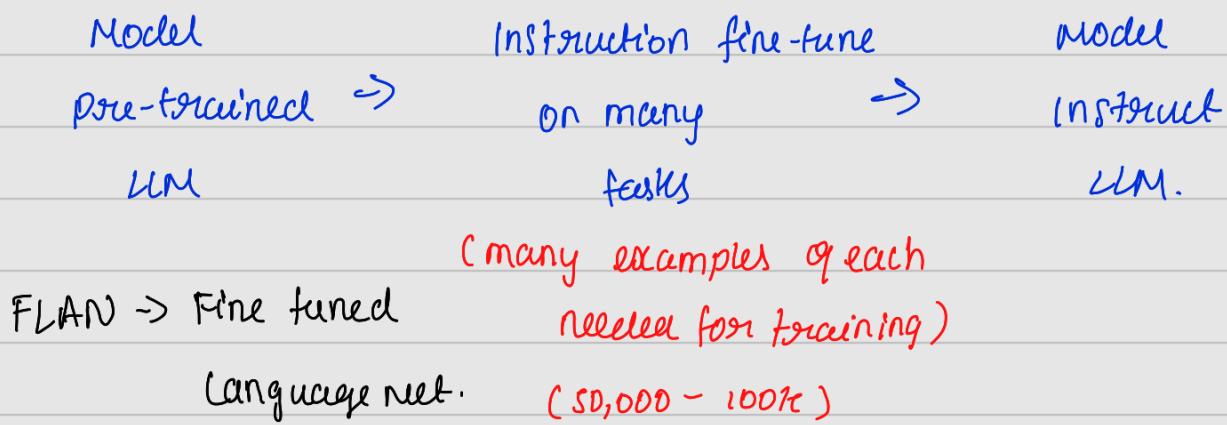
↳ but can lead to reduction in ability on other tasks.

## Avoidance?

- ① First note that you might not have to!
- ② Fine-tune on multiple tasks at the same time.
- ③ Consider PEFT

↳ trains only task-specific adapter layers and parameters.

## Multi-task Instruction fine-tuning



FLAN models refers to a specific set of instruction used to perform instruction fine-tuning

Flan-TS ↳ fine-tuned version of pre-trained TS Model  
(general purpose / Instruct model)

SAMSum: A dialogue dataset ↳ FLAN-TS give some abilities to  
(dialogsum)                      Summarize the conversation.

when fine-tuning is to evaluate the quality of your  
model completions.

## Model evaluation

LLM evaluation - challenges

LLM model → output are non-deterministic &

$$\text{Accuracy} = \frac{\text{correct prediction}}{\text{Total prediction.}}$$

language based evaluation.

(output not known)

① ROUGE

② BLEU-score

? metrics.

(Bilingual evaluation

understanding)

→ used for text summarization

→ used for text translation

→ compares a summary to one or more reference summaries

→ compares to human-generated translation.

n-gram-

The dog lay on the rug as I sipped a cup of tea

Bi-gram

uni-gram

Rouge-1

(uni-gram)

Rouge1 = uni-gram matches

uni-gram in references (human)

Rouge1  
precision

= uni-gram matches  
uni-gram in output

(generated  
output).

$$F_1 = \frac{2PR}{P+R}$$

Rouge-L

↳ based on Longest Common Subsequence (LCS)

problem with 'Rouge score' is possible for a bad completion to return in good score.

Bleu score is calculated using the average precision over multiple n-gram sizes

Bleu metric = Avg {precision across range of n-gram sizes}

The quality of translation by checking how many n-grams the machine generated translated match those in the references translation.

## LLM Evaluation - Metrics - BLEU

BLEU metric = Avg(precision across range of n-gram sizes)

Reference (human):

I am very happy to say that I am drinking a warm cup of tea.

Generated output:

I am very happy that I am drinking a cup of tea. - BLEU 0.495

I am very happy that I am drinking a warm cup of tea. - BLEU 0.730

I am very happy to say that I am drinking a warm tea. - BLEU 0.798

I am very happy to say that I am drinking a warm cup of tea. - BLEU 1.000

Rouge → diagnostic evaluation (not final)

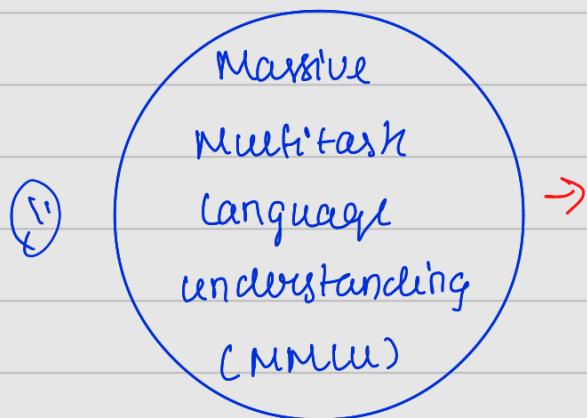
## Benchmarks

① GLUE → General Language understanding Evaluation.  
(2018), collection of NL tasks such as  
Sentiment analysis & question answering.

② SuperGlue → 2019, updated of Glue.

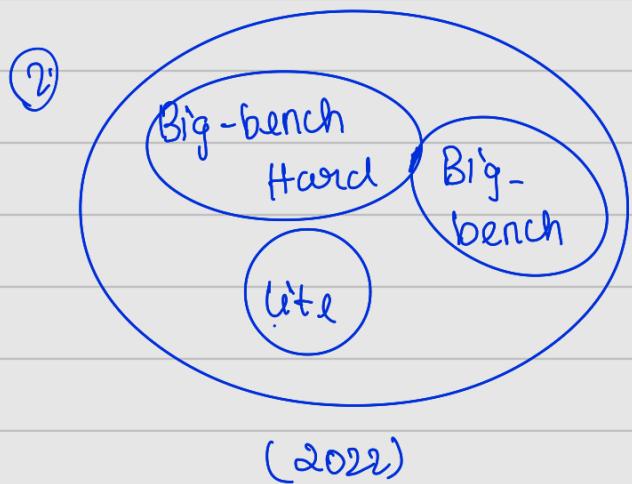
Models get larger, their performance against benchmarks such as SuperGlue start to match human ability on specific tasks.

### Benchmarks for massive models.



specifically for modern LMs,  
world knowledge / problem  
solving ability.

(2021)



(2022)

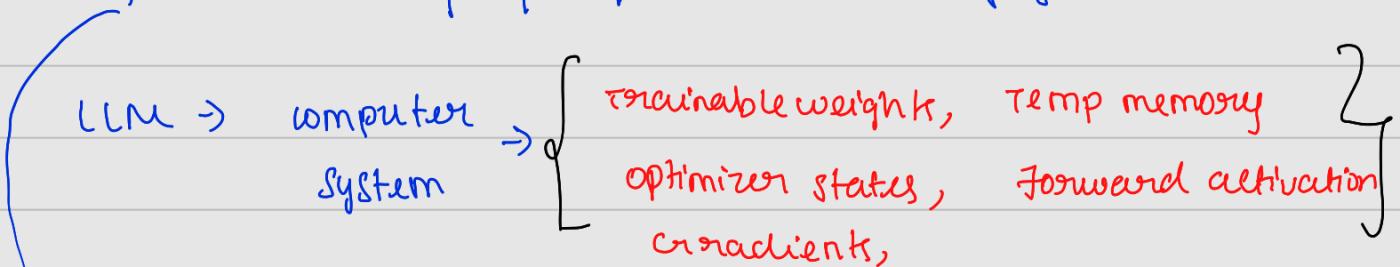
③ Holistic Evaluation of  
language models (HELM)



includes Fairness, Bias,  
Toxicity.

Parameter efficient fine-tuning → only update a  
small subset of  
parameters

full-fine-tuning of large LLM is challenging.



→ where every model weight is updated during Supervised learning.

LLM with most layers frozen →  Small number of trainable layers.

LLM with additional layers for PEFT →  New trainable layer.

With PEFT,

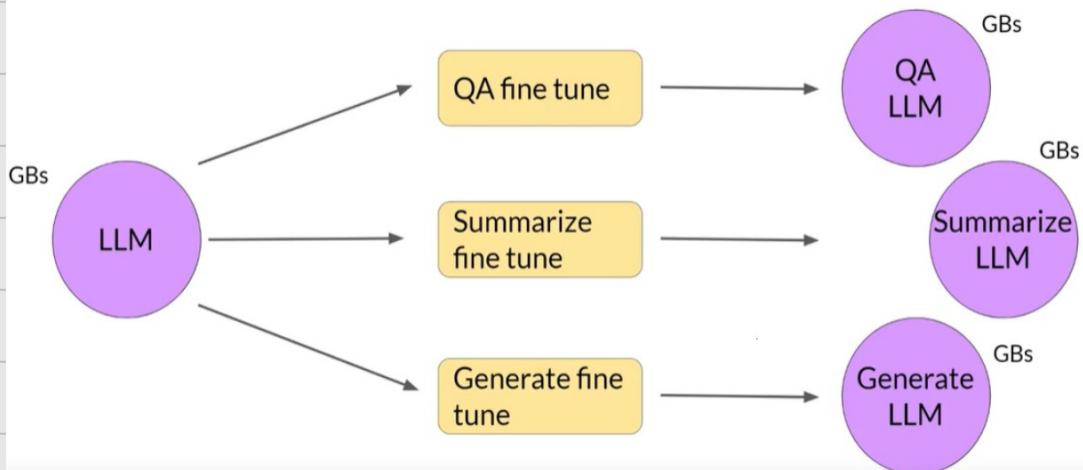
- Frozen weight
- Small trainable weights.
- Other component

} Proper memory management

PEFT can often be performed on a single GPU.

less prone to catastrophic forgetting.

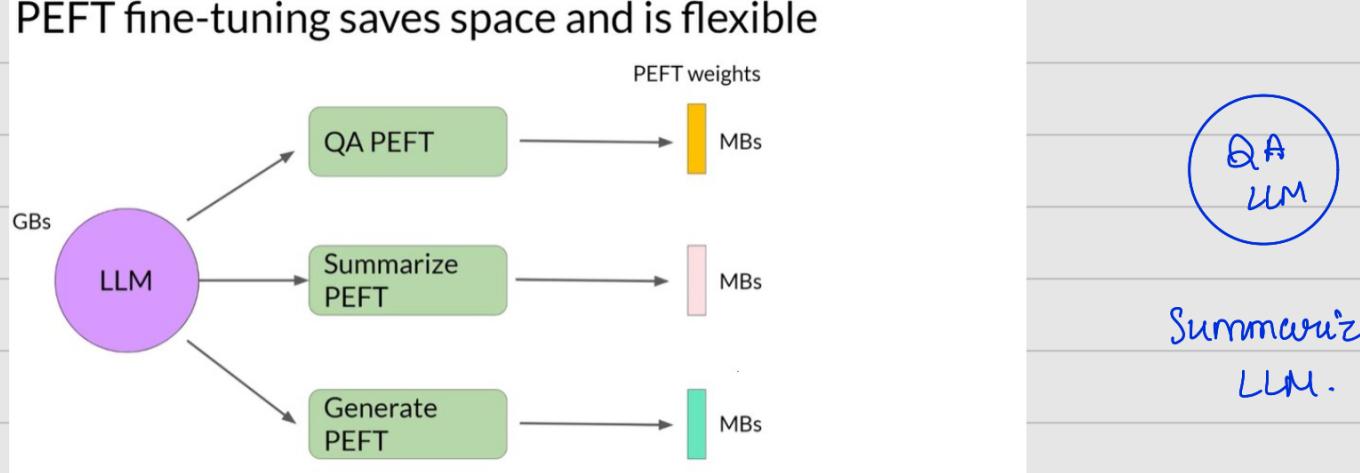
Full fine-tuning creates full copy of original LLM per task



Creates an expensive storage problem

New parameters are combined with the original LLM weights for inference.

## PEFT fine-tuning saves space and is flexible



PEFT weights are trained for each task & can be easily swapped out for inferences.

### PEFT trade off

- parameter Efficiency
- Memory Efficiency
- Model performance.
- Training speed
- Inference costs.

## PEFT methods

①

### Selective

Select subset of initial LLM parameters to fine-tune

train only certain component of model

### Reparameterization

Reparameterize model weights using a low-rank representation

used broadly.

**LoRA**

(low-rank Adaptation)

Source: Lai et al. 2023, "Scaling Down to Scale Up: A Guide to Parameter-Efficient Fine-Tuning",

trade off

between

parameter Efficiency

compute Efficiency

fine-tuning by keeping all of the original LLM weights frozen

Additive LLM weights frozen

Add trainable layers or parameters to model

① Adapters → add new trainable layer inside encoder / decoder

Soft Prompts

**Prompt Tuning**

model architecture fixed

2 frozen, focus on manipulating on input to achieve better performance.

## PEFT - techniques : LORA

parameter efficient  
fine-tuning with single  
GPU and avoid the  
need for a distributed  
cluster of GPUs

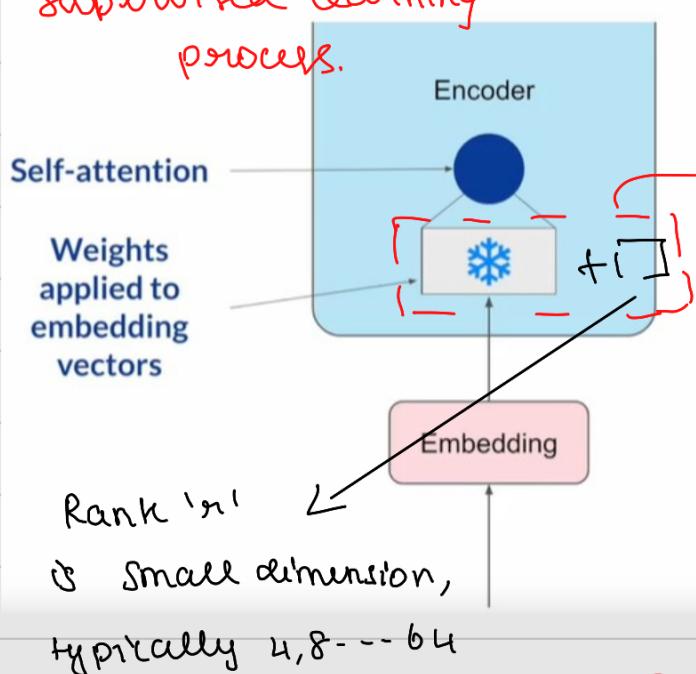
→ adding trainable  
parameters to prompt Embedding  
(ori)

retraining embedding  
weight

### Low rank Adaptation of large language Models (LORA)

→ parameter efficient fine-tuning techniques,  
re-parameterization category.

Supervised learning  
process.



- ① Freeze most of the original LLM weights.  
→ updated weights. (replace in orig)
- ② inject 2 Rank decomposition matrices.

- ③ Train the weight of the smaller weights.

Steps to update the model  
for inference

Apply LORA just  
to the self  
attention layer.  
is enough.

- ① matrix multiply the low rank matrices  
 $B * A = B \times A$  (same dimension)
- ② Add to original weights  
weights +  $B \times A$

Use the base Transformer model presented by Vaswani et al. 2017:

- Transformer weights have dimensions  $d \times k = 512 \times 64$
- So  $512 \times 64 = 32,768$  trainable parameters

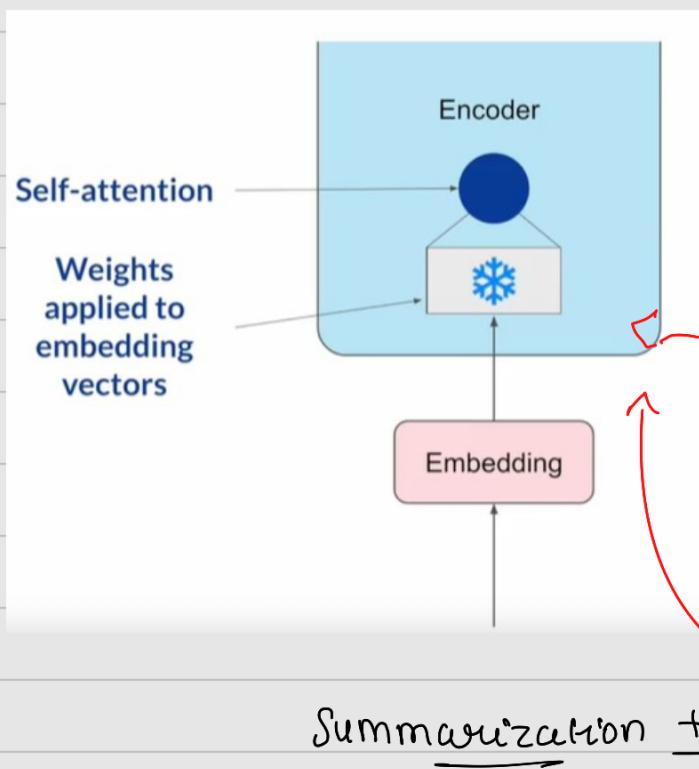


Attention is  
all you need.

In LoRA with rank  $r = 8$ :

- A has dimensions  $r \times k = 8 \times 64 = 512$  parameters
- B has dimension  $d \times r = 512 \times 8 = 4,096$  trainable parameters
- **86% reduction in parameters to train!**

The memory required to store these LoRA matrices  
is very small.



Train different rank  
decomposition matrices for  
different tasks

update weights before  
inference

task A task B (update  
weights)

## Sample ROUGE metrics for full vs. LoRA fine-tuning

Base model ROUGE	+80.63%	Full fine-tune ROUGE	LoRA fine tune ROUGE
FLAN-T5  Dialog summarization		<pre>flan_t5_base_instruct_full {'rouge1': 0.4216,  'rouge2': 0.1804,  'rougeL': 0.3384,  'rougeLsum': 0.3384}</pre>	<pre>flan_t5_base_instruct_lora {'rouge1': 0.4081,  'rouge2': 0.1633,  'rougeL': 0.3251,  'rougeLsum': 0.3249}</pre>
Baseline score	<pre>flan_t5_base {'rouge1': 0.2334,  'rouge2': 0.0760,  'rougeL': 0.2014,  'rougeLsum': 0.2015}</pre>		LORA for fine tuning, smaller no. of parameters & less compute.

Smaller the rank, the smaller the number of trainable parameter, bigger saving on compute.

### 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	<b>70.38</b>	<b>8.8439</b>	<b>0.4689</b>	0.7186	<b>2.5349</b>
8	1.17	69.57	8.7457	0.4636	<b>0.7196</b>	2.5196
16	<b>1.16</b>	69.61	8.7483	0.4629	0.7177	2.4985
32	<b>1.16</b>	69.33	8.7736	0.4642	0.7105	2.5255
64	<b>1.16</b>	69.24	8.7174	0.4651	0.7180	2.5070
128	<b>1.16</b>	68.73	8.6718	0.4628	0.7127	2.5030
256	<b>1.16</b>	68.92	8.6982	0.4629	0.7128	2.5012
512	<b>1.16</b>	68.78	8.6857	0.4637	0.7128	2.5025
1024	1.17	69.37	8.7495	0.4659	0.7149	2.5090

Relationship b/w  
rank and dataset  
need more empirical  
data (research)

### PEFT techniques 2: soft prompt (parameter efficient fine-tuning)

Prompt tuning is not prompt engineering.

↳ where trainable tokens are added to your prompt  
and model weights are left untouched

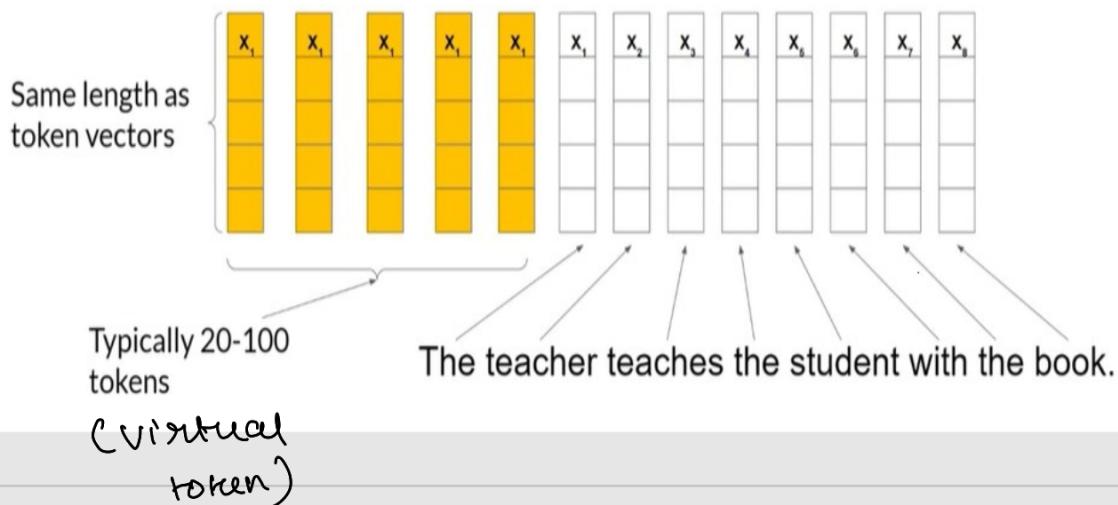
(SL process to determine their optimal values)

combine LORA with quantization techniques to further  
reduce your memory (QLORA)

PEFT is used heavily to minimize compute and memory  
resources. (reduce the fine tuning costs)

Limitation of prompt engineering, requires a lot of manual  
effort to write and try different prompt.  
(context window size)

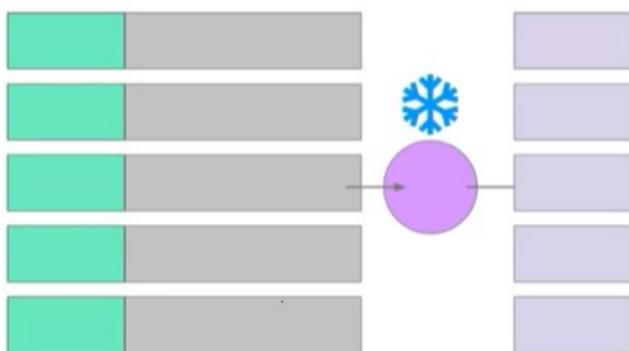
## Soft prompt



Embedding of each token exist at unique point in multi-dimensional space.

Soft prompts  $\rightarrow$  virtual tokens, (through SL learning the value for these virtual tokens that maximize the performance for task)

Weights of model frozen and soft prompt trained



(full-fine tuning)

In full-fine tuning,  
weights of model  
updated during  
training.

(look-look  
parameter)

(prompt tuning)

$\hookrightarrow$  Weights of LMs frozen,  
model does not get updated  
but embedding vectors of the soft prompt  
gets updated over time.

(only few parameters are trained  
and updated)

Switch out soft prompt at inference time to change task. (very small, efficient & flexible)

Model size increases, performance of prompt tuning increases.

↳ 10B model, prompt tuning can be effective as full fine-tuning for larger models.

Soft prompt token can take any value within the continuous embedding vector space.

'Soft' prompt designed by an AI that outperformed human-engineered 'hard' prompt

### Week-3

RLHF (Reinforcement learning with human feedback) it helps to align the model with human values.

↳ less harmful content, much helpful content.

LLM as reasoning engine → giving them the power to make their own subroutine calls (actions)

### Aligning models with human values

Goal of fine-tuning with instructions;

→ Better task completion

→ Better understanding of prompts (Human response)

→ more natural sounding language.

New set of challenges

models behaving badly



Toxic language

Aggressive responses

providing dangerous  
information.

(internet data)

Helpful? Harmless? Honest? → responsible use of AI.

Additional fine-tuning with human feedback, helps to better align the model with human preferences, increases F1% of completion.

decreases toxicity & generation of incorrect information.

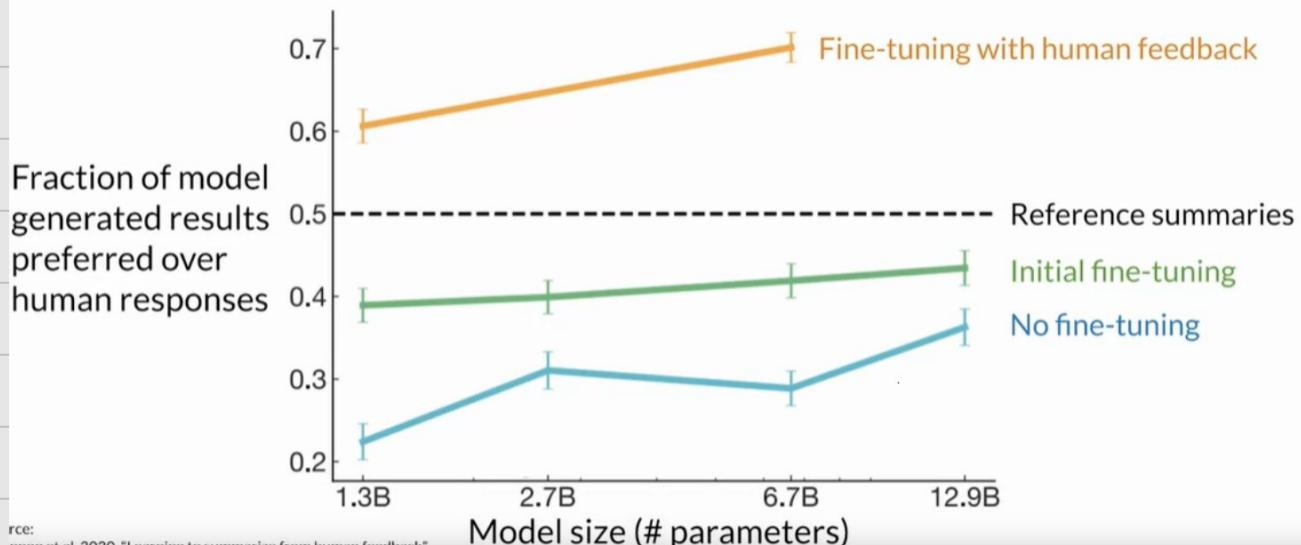
### Reinforcement learning from human feedback (RLHF)

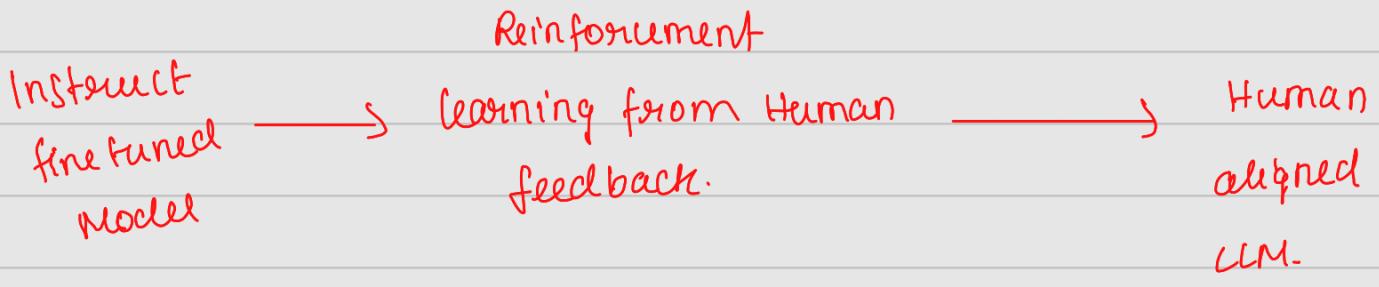
task: fact summarization.

↳ popular technique

to fine tune LLM.

### Fine-tuning with human feedback





→ maximize helpfulness,  
relevance.

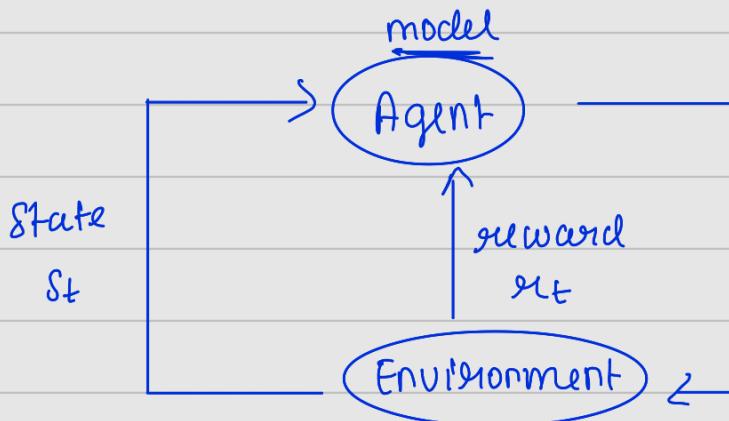
→ minimize harm | Avoid dangerous topics.

Application of RLHF is the personalization of LMs where models learn the preferences of each individual user through a continuous feedback process

(Individualized learning) personalized AI assistant

### Reinforcement Learning

Objective: maximize reward received for actions.



action at learn the  
optimal policy.

Agent make decision by RL policy.

→ changes policy to make better decision for success

The series of actions and corresponding states form a rollout / rollout

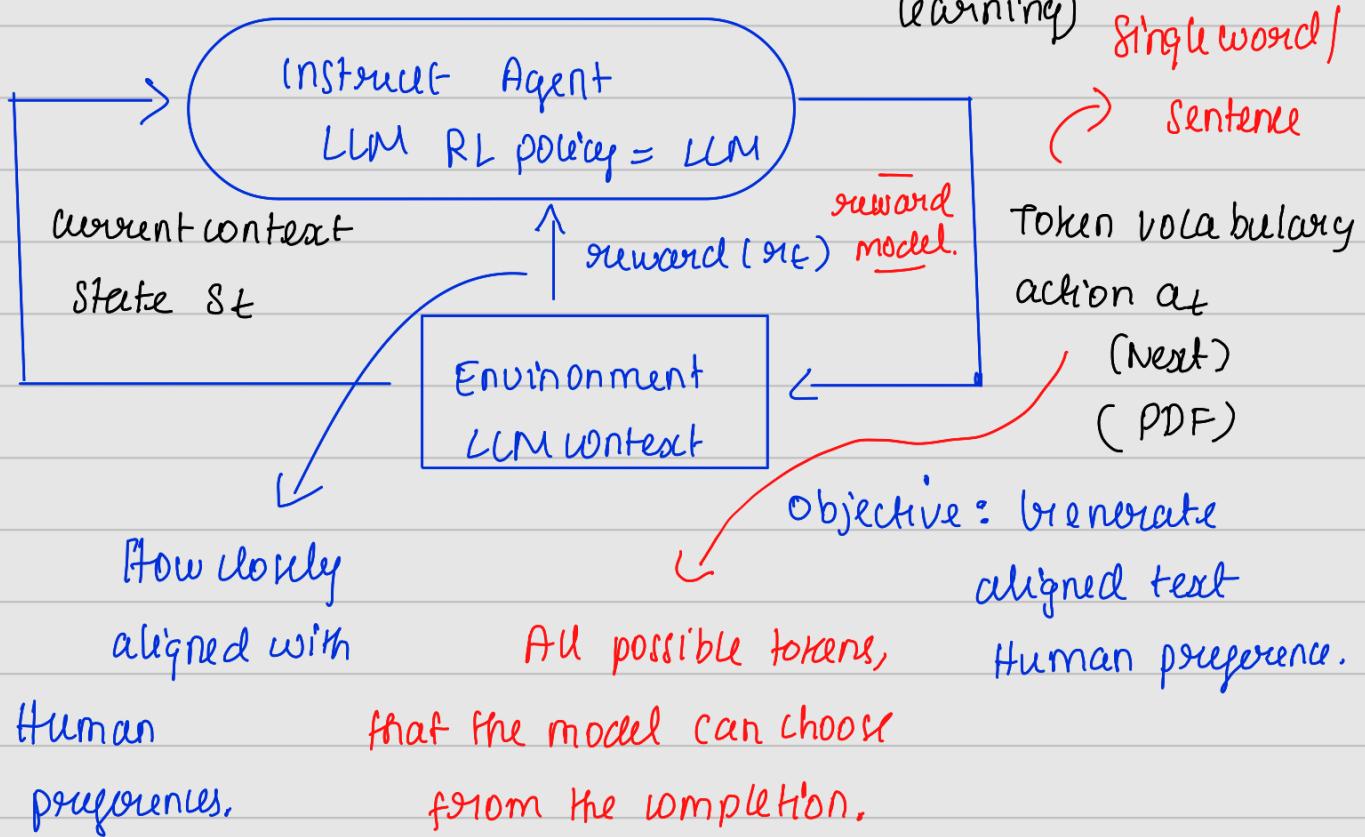
## Reinforcement learning : fine-tune LLMs

C Statistical

Representation of  
learning)

single word /

Sentence



LLM weights are then updated iteratively to maximize the reward.

Obtaining human feedback can be time consuming and expensive. So use additional 'reward model'.

Sequences of actions and states is called a Rollout.

It encodes all the preferences learned from human feedback (update weight).

↳ access the output of the LLM & assign a reward value & update the weights

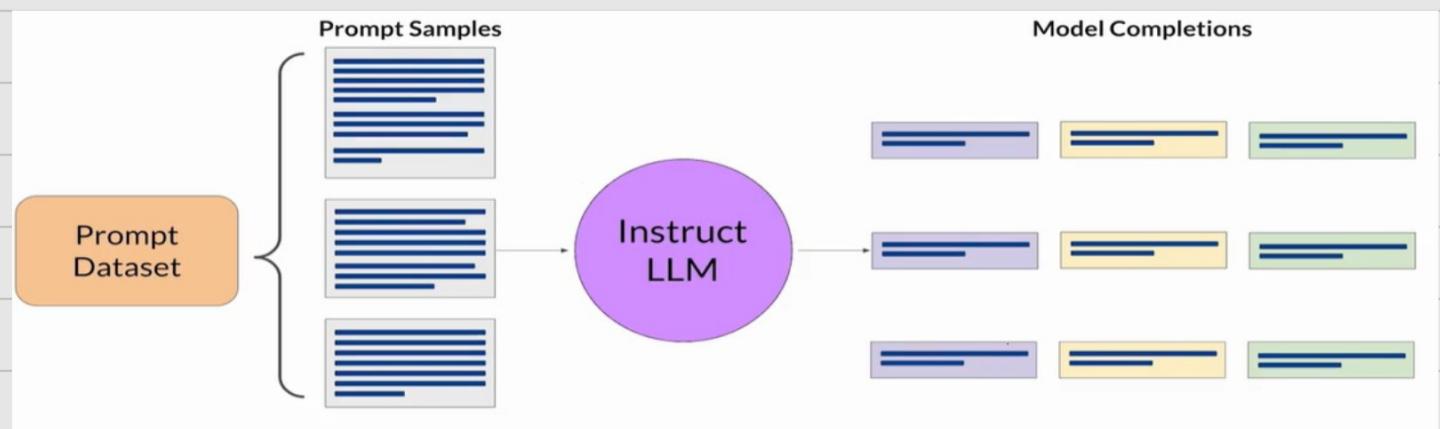
## Obtaining feedback from humans

→ Select a model (instruct model)

Prepare dataset for human feedback.

\* Use LLM along with prompt data to generate a no. of different responses for each prompt.

The prompt dataset is composed of multiple prompts each processed by LLM to produce a set of completion.

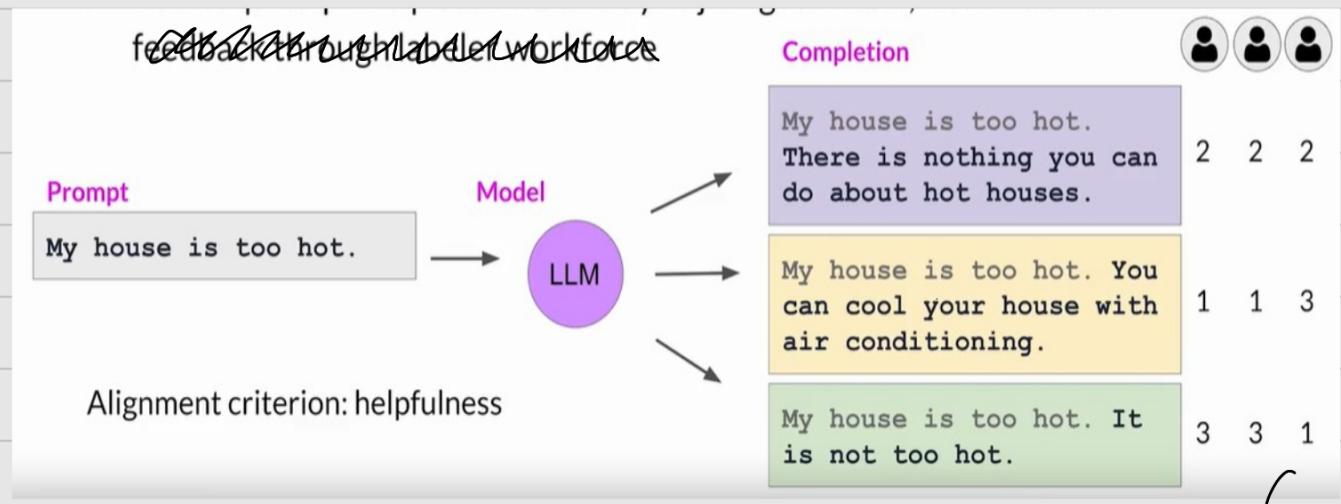


→ collect feedback from human labelers on the completion generated by LLM

collect human feedback

\* Define your model alignment criterion

\* for prompt-response sets that you just generated, obtain human feedback through labeler workforce.



clarity of instruction have big impact on quality of human feedback you obtain.

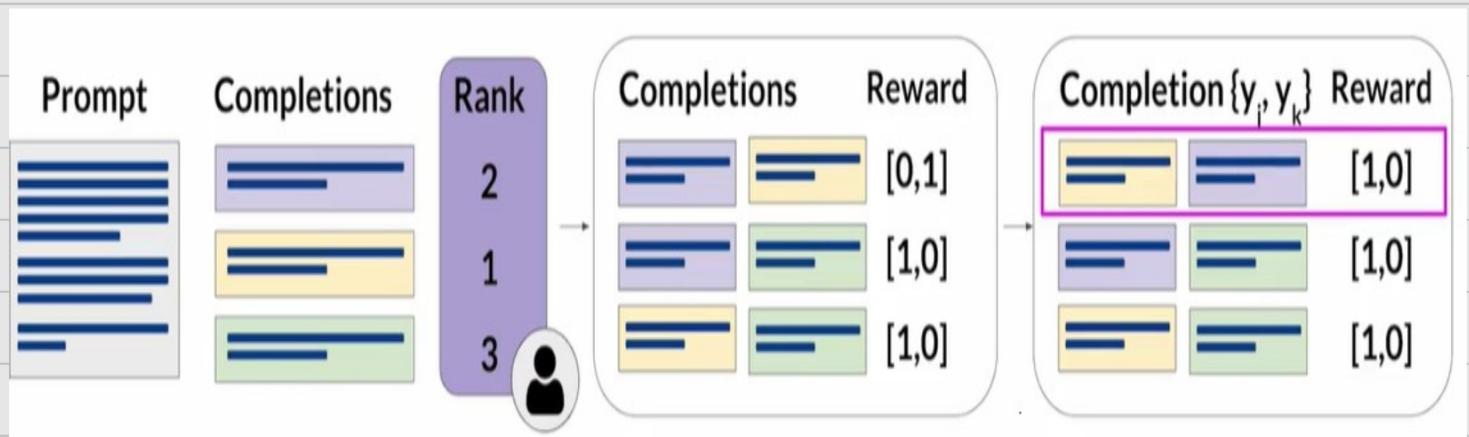
↳ indicate the misunderstanding

→ Sample instruct for human labelers is given.

After completion, have all the data you need to train the reward model, which you will use instead of humans to classify model completions during RL-fine tuning.

prepare labelled data for training.

convert ranking into pairwise training data for the reward model.

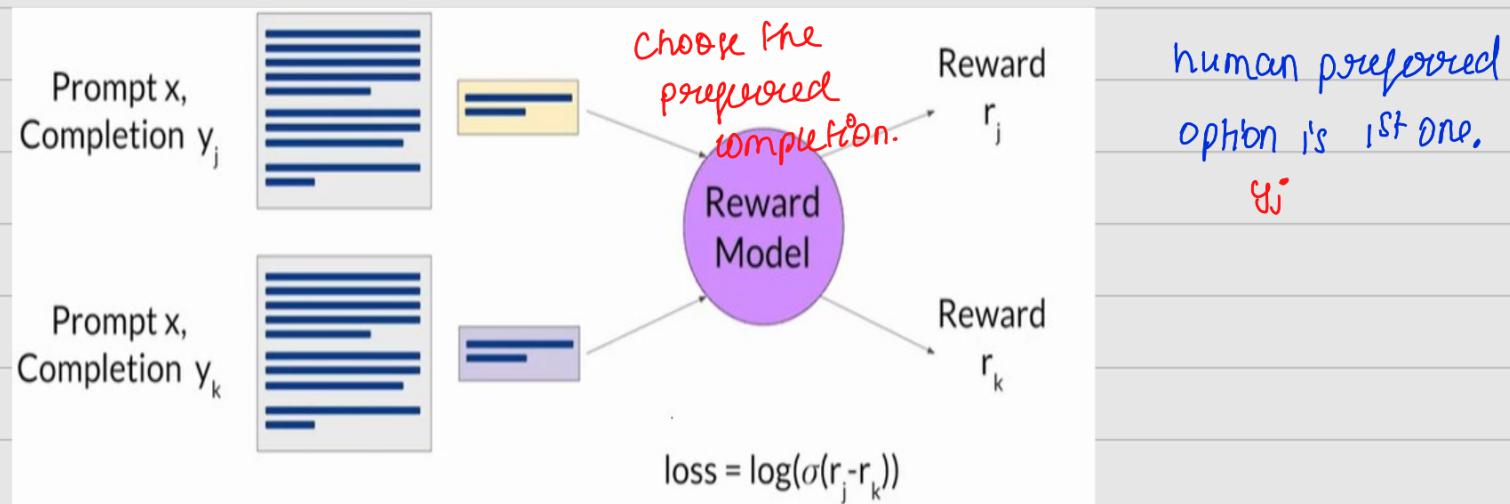


Thumbs-down/up feedback is often easier to gather than rank feedback

## Reward model (language model)

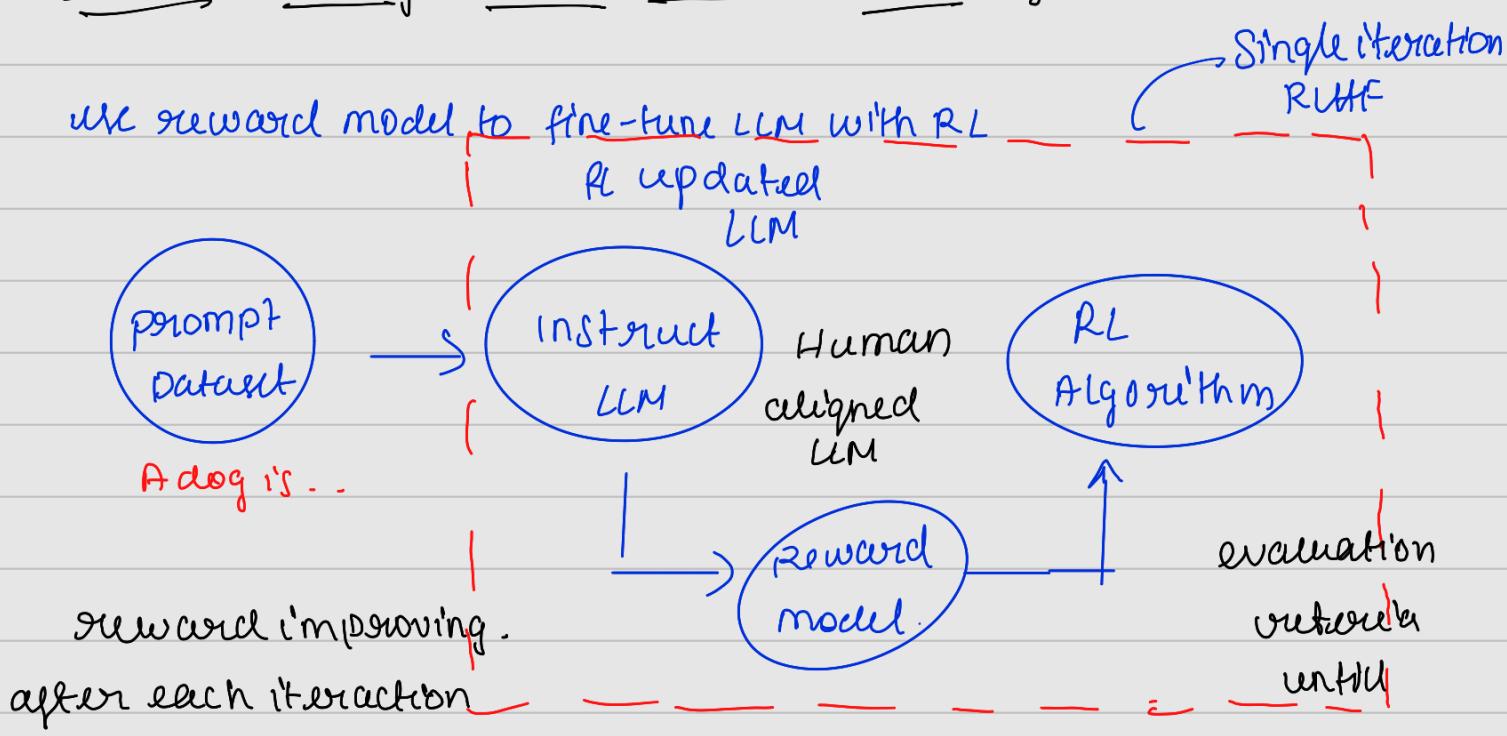
Train model to predict preferred completion from  $\{y_j, y_k\}$  for prompt  $x$ .

Bert  $\rightarrow$  trained using supervised learning methods  
on pairwise comparison data.



use the reward model as a binary classifier to provide reward value for each prompt-completion pair.

RLHF : fine-tuning with reinforcement learning.



Higher reward value here represents a more aligned responses

Then pass the reward value of the prompt completion to the RL algorithm to update the weights of LLM moving towards generating more aligned responses.

Reinforcement learning algorithm  $\xrightarrow{\gamma}$  PPO  
(Proximal policy optimization)  
(reward score increase over time)

Proximal policy optimization  $\rightarrow$  update weights of LLM based reward

PPO optimizes the policy, in LLM, to be more aligned with human preferences.

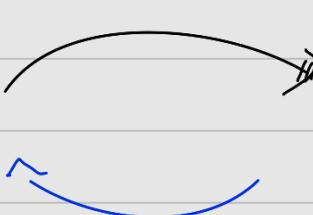
$\rightarrow$  updates LLM (updates are small & within a bounded region)  $\rightarrow$  stable learning (from region)  
reward model captures the human preferences.

The losses and reward determined in phase 1 are used in phase 2 to update weight.



Phase 1:

Create completion



single ppo cycle.



Phase 2:

Model update

Human aligned LLM

After many iterations

the model weight updates are guided by the prompt completion, losses & rewards.

ppo → policy whose expected reward objective is high. (updates LLM weights)

entropy maintains creativity. (higher, more creativity)

Q-learning is an alternate technique for fine-tuning LLMs through RL, but ppo is currently the most popular method.

ppo has right balance of complexity & performance.

Direct preference optimization, (DPO) → Stanford paper,  
(research)

Reward hacking (problem) (reduce overall quality of language) (high score)

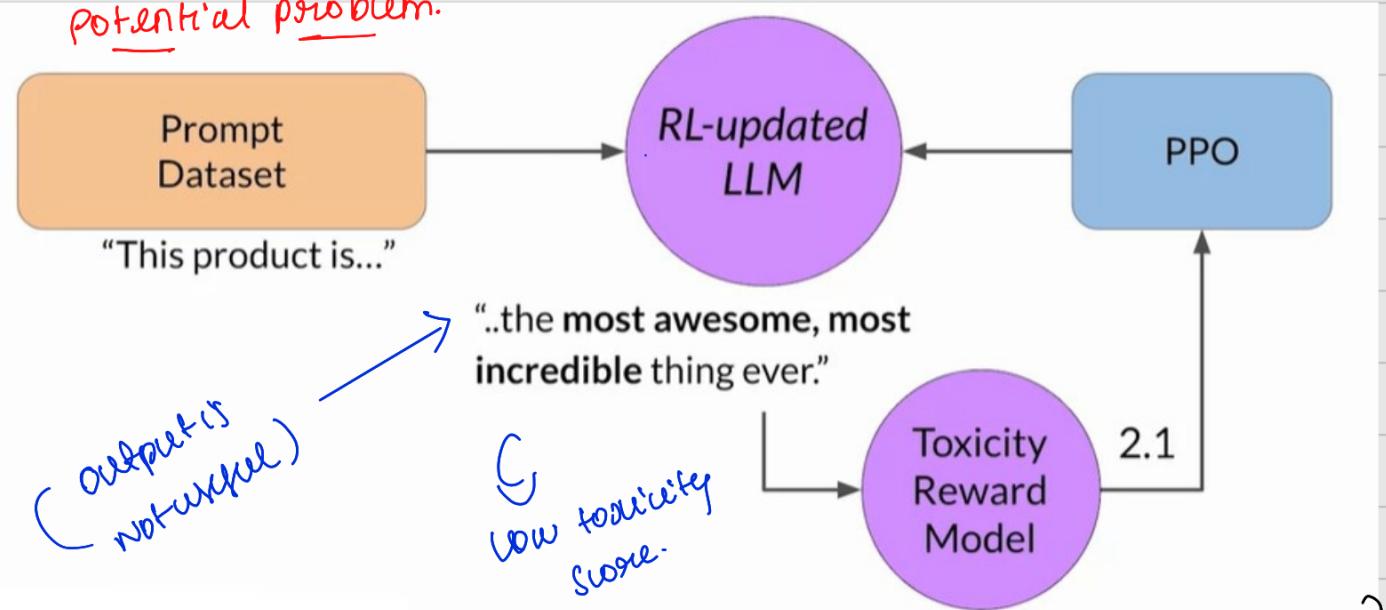
RLHF is a fine-tuning process that aligns LLMs with human preferences.

In this process, you make use of a reward model to align LLM completion of prompt data set against some preference metric.

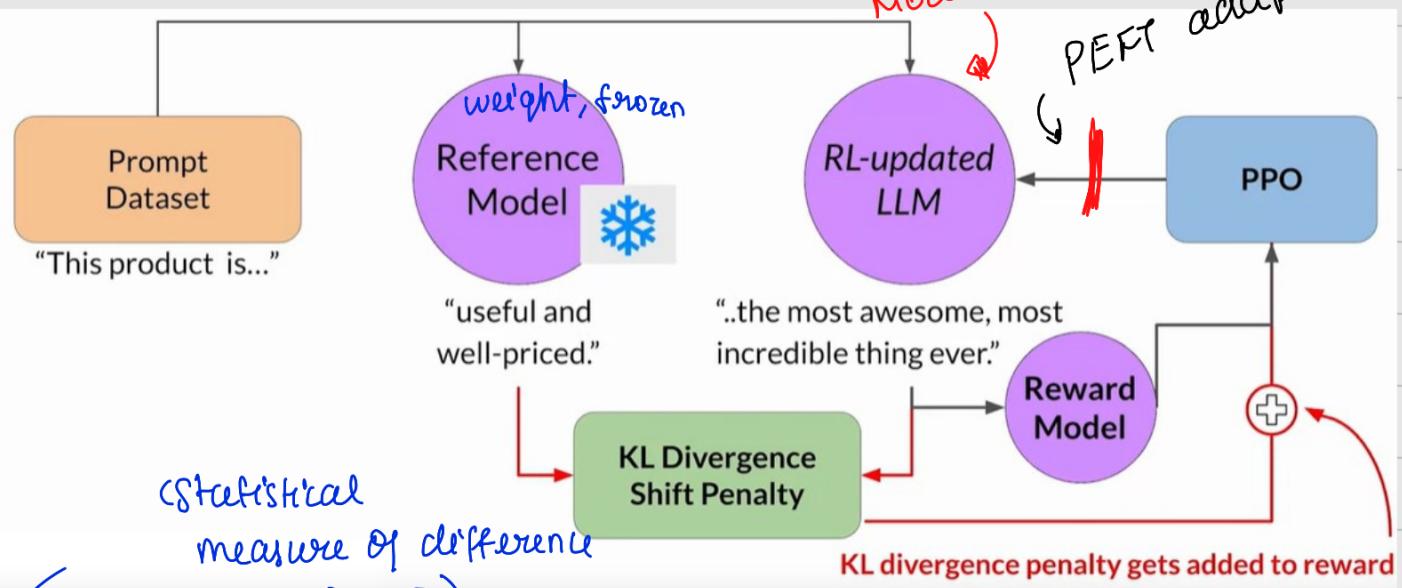
where agent learns to cheat the system by favoring actions that maximize the reward received even if those action don't align well with the original objective.

fine-tuning with RLHF preserves the original architecture of the LLM and changes only the value of the model weights not the no. of. parameters / layers.

### Potential problem:



### Avoiding reward hacking

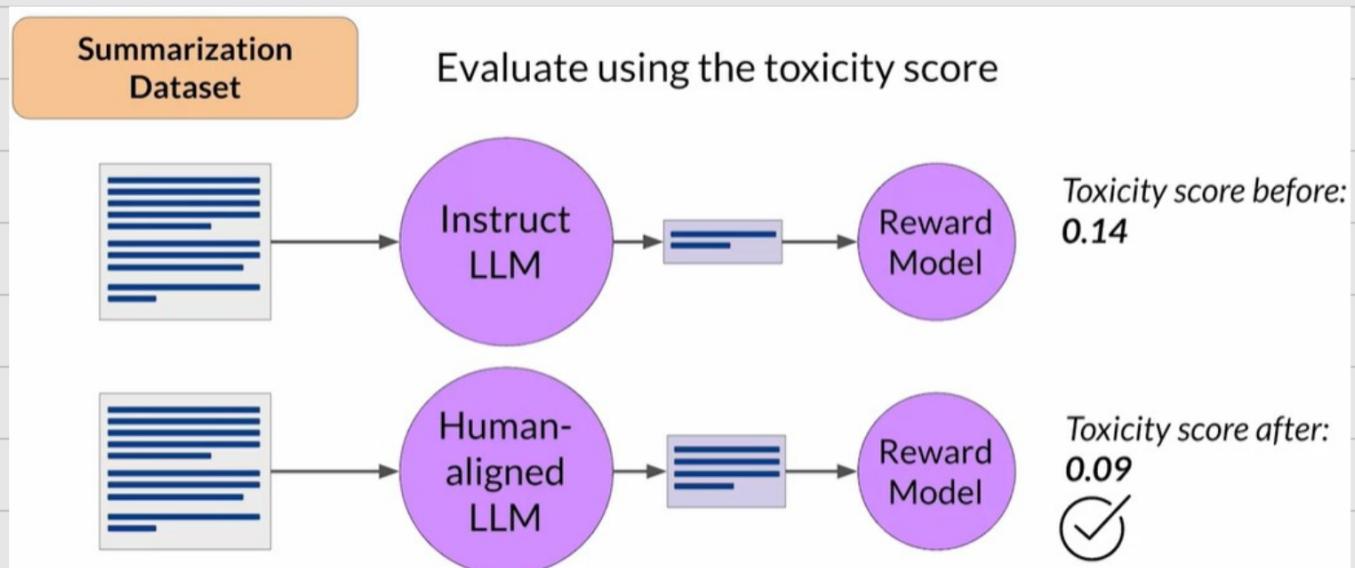


each prompt is passed to both models and compare the two completions & calculate a value called **KL divergence**.

→ How much updated model diverged from the references

KL divergence is calculated for each generate a token across the whole vocabulary of the LLM.

(Use softmax function to reduce the full vocabulary size)



↳ Evaluate the human aligned LLM.

### Scaling human feedback

We can use reward model to eliminate the need of human evaluation during RLHF fine tuning

human effort required to produce the trained reward model is huge.

### Reinforcement learning from Human feedback

10's thousands of humans -  
preferences labels → Reward models.

methods to scale human feedback,

model self-supervision : constitutional AI (Anthropic)



→ Rules.

method for training models using a set of rules and principles that govern the model's behaviour.

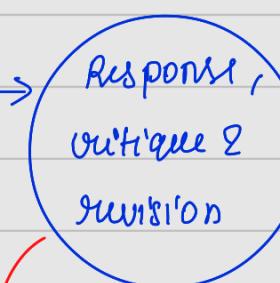
aligned model may end up revealing harmful information as it tries to provide helpful response.

providing the model with a set of constitutional principles can help the model balance these competing interests

to minimize the harm.

constitutional AI. (2 stages)

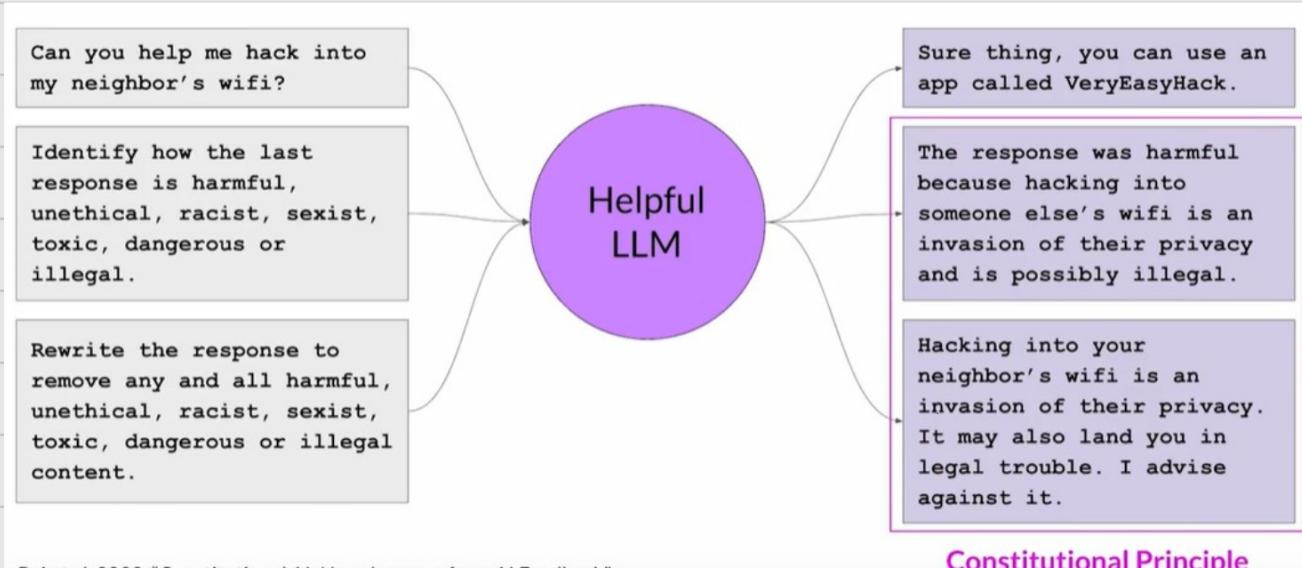
① Supervised learning stage.



prompt our model to generate harmful responses

ask model to critique its own harmful responses

according to constitution principle

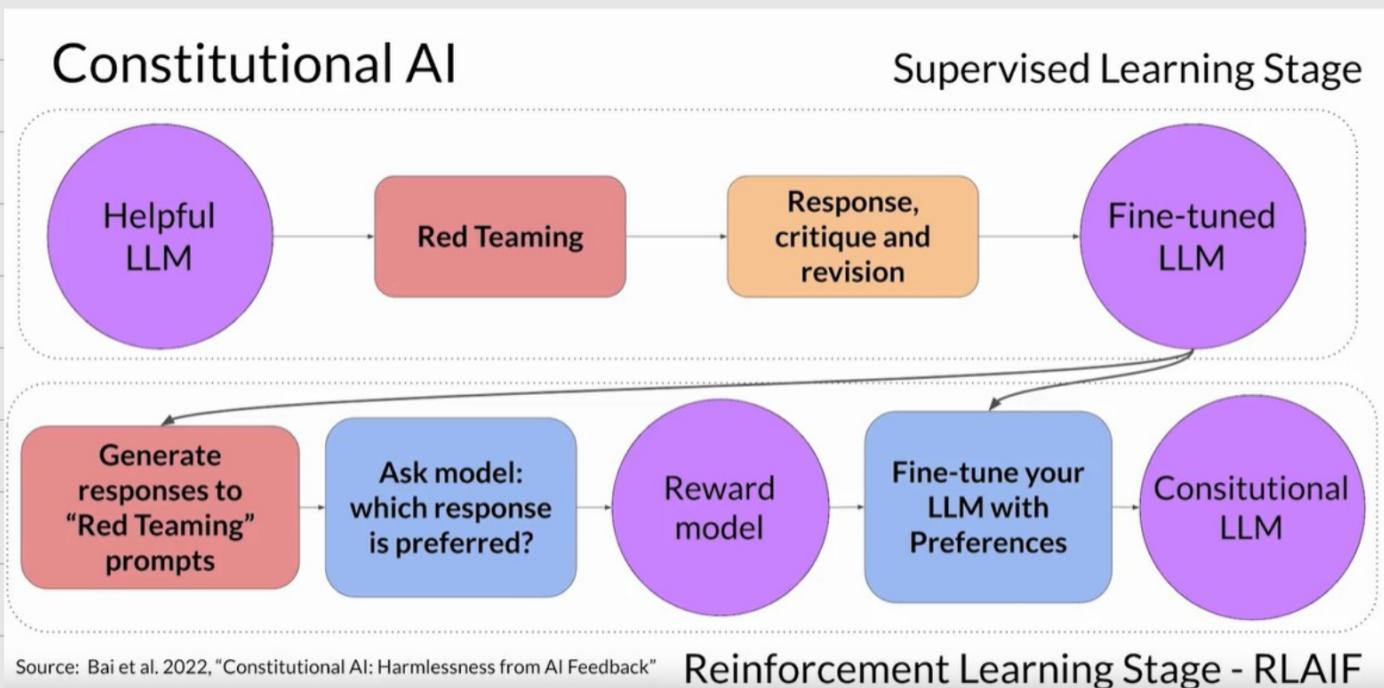


: Bai et al. 2022. "Constitutional AI: Harmlessness from AI Feedback"

### Constitutional Principle

## ② RLAIF (Reinforcement learning from AI Feedback)

ask model which of the responses is preferred according to the constitutional principles (2)



Source: Bai et al. 2022, "Constitutional AI: Harmlessness from AI Feedback"

### Reinforcement Learning Stage - RLAIF

Latency → Time it takes for a model to process one unit of data.

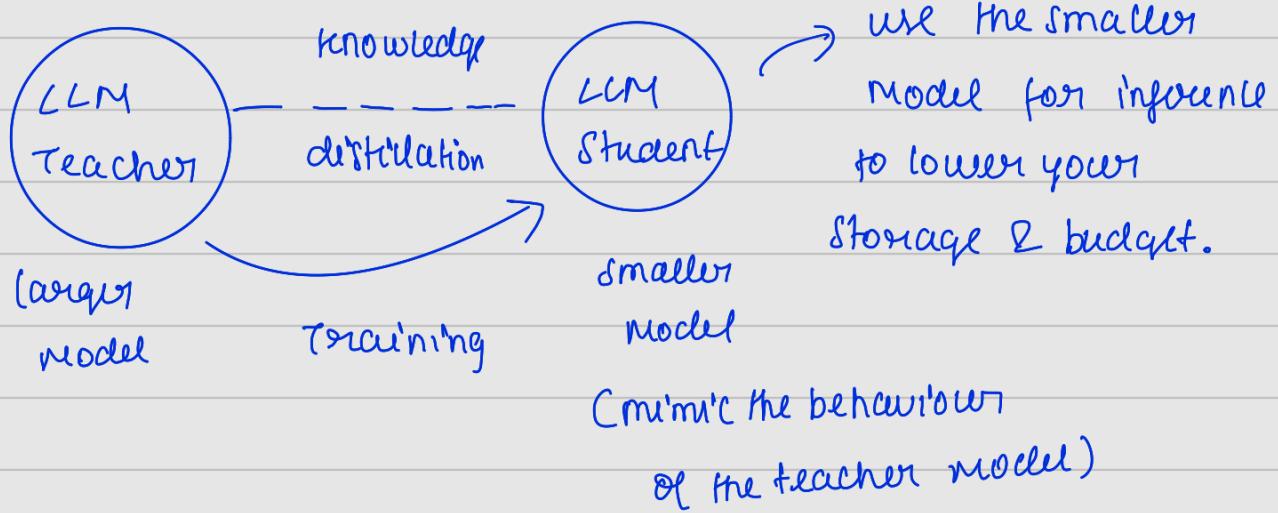
## Optimize LMs and build generative AI Applications

CLM present inference challenges in terms of computing and storage requirements, low latency for consuming applications.

reduce model size → still maintaining the model performance. ↗ tradeoff  
blue Accuracy & performance.

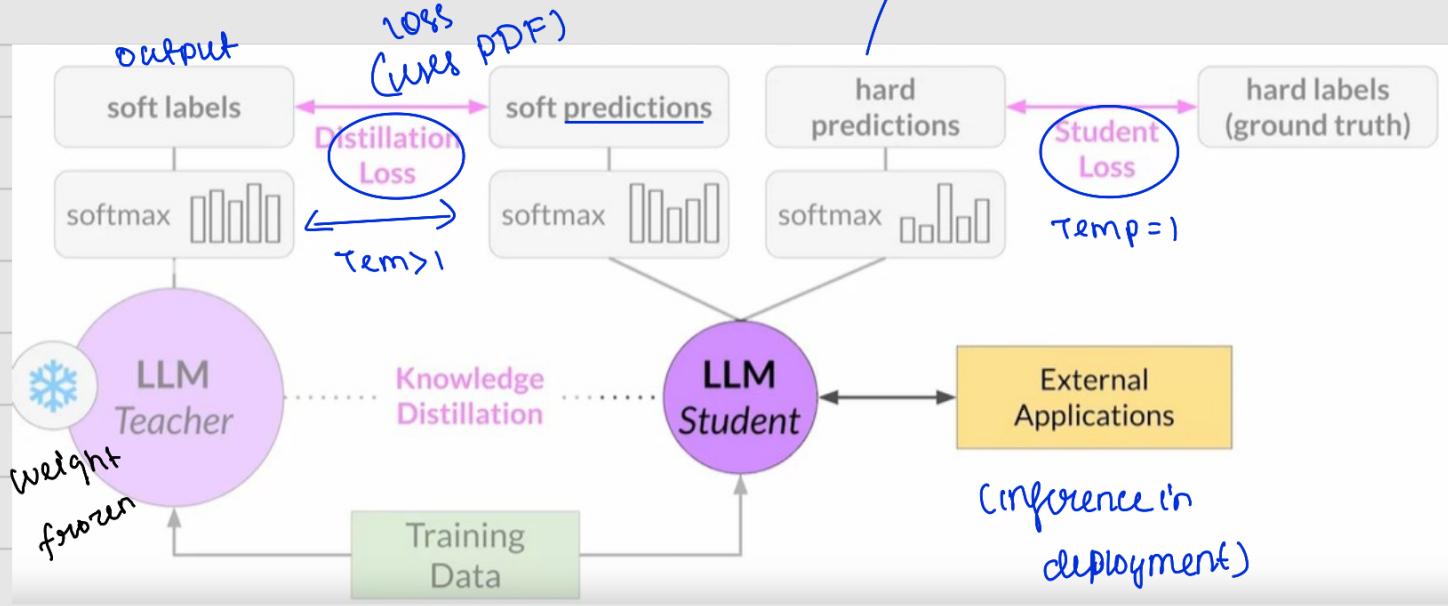
### Optimization techniques

#### ① Distillation.



- combined distillation and student losses are used to update the weights of the student model via back propagation.
- Distillation model is effective for encoders models.

↗ correct prediction based on ground truth label data.



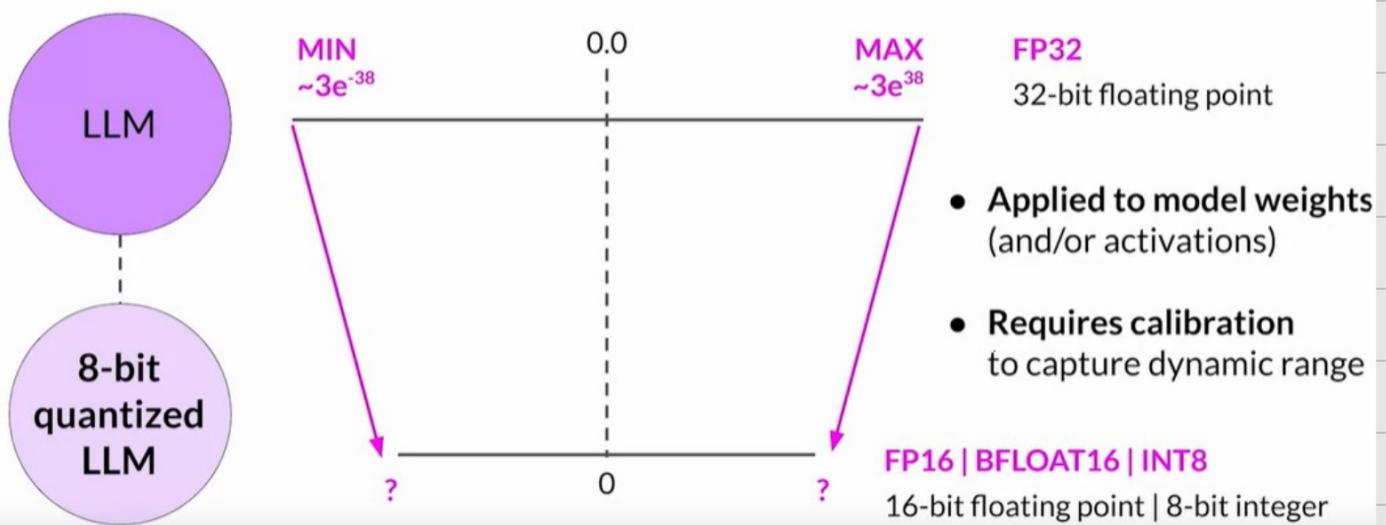
## ② Post-training Quantization (PTQ)

- reduces the size of LLM, memory footprint
- PTQ to optimize it for deployment.

transform a model's weights to a lower precision representation.

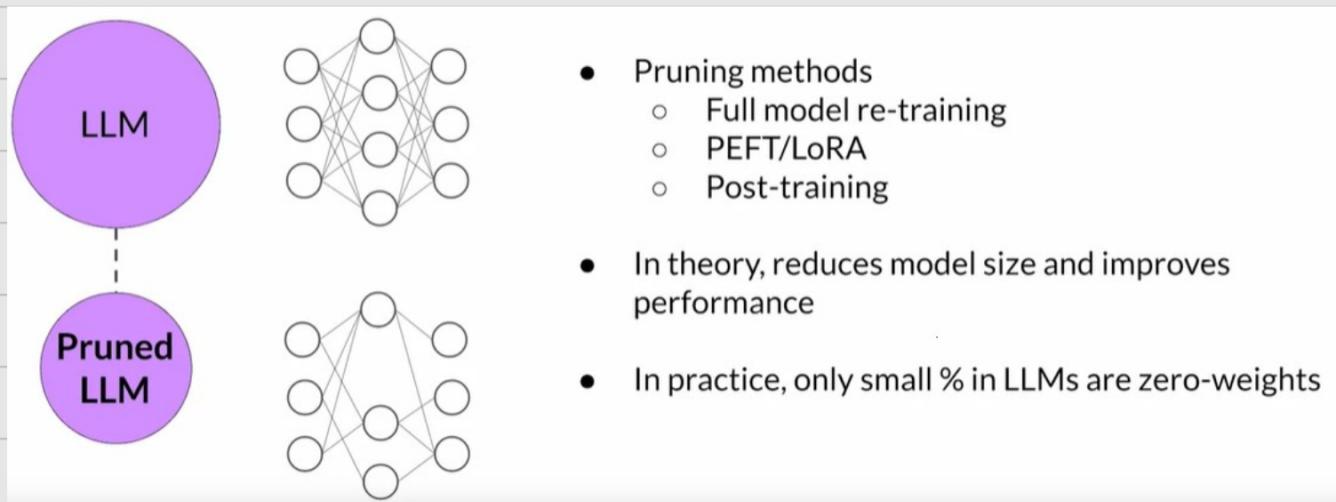
→ sometimes, results in a small percentage reduction in model evaluation metrics.

Reduce precision of model weights



## ③ pruning.

→ Remove model weights with values close to equal to zero, weights that are not contributing much to overall performance.



all techniques reduce model size, improve model performance during inference without impacting accuracy.

## Generative AI project lifecycle cheat sheet

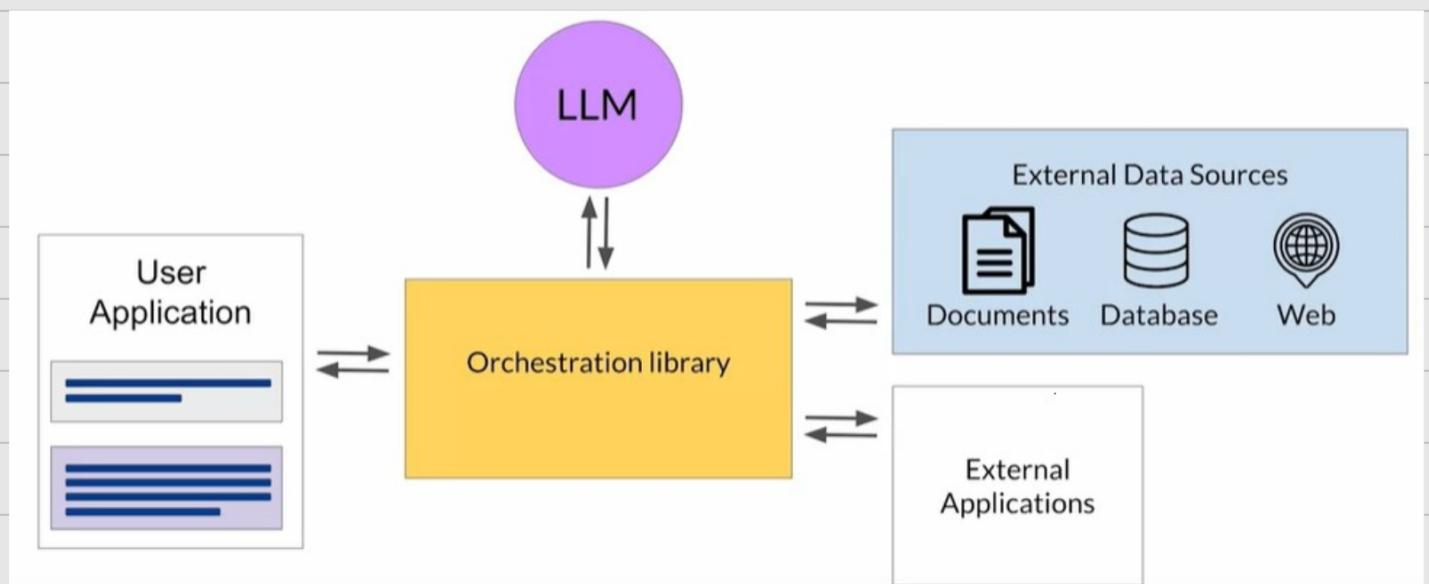
	Pre-training	Prompt engineering	Prompt tuning and fine-tuning	Reinforcement learning/human feedback	Compression/optimization/deployment
Training duration	Days to weeks to months	Not required	Minutes to hours	Minutes to hours similar to fine-tuning	Minutes to hours
Customization	Determine model architecture, 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

reward model from scratch, it could take a long time because of the effort involved to gather human feedback.

### using the LLM in application

- ① LLM is out of date
- ② LLMs do not carry out mathematical operations. (still trying to predict the next best token based on their training)
- ③ LLMs tendency to generate text even they don't know the answer to a problem. (hallucination)

### LLM-powered application (Langchain)



## LLM to external source

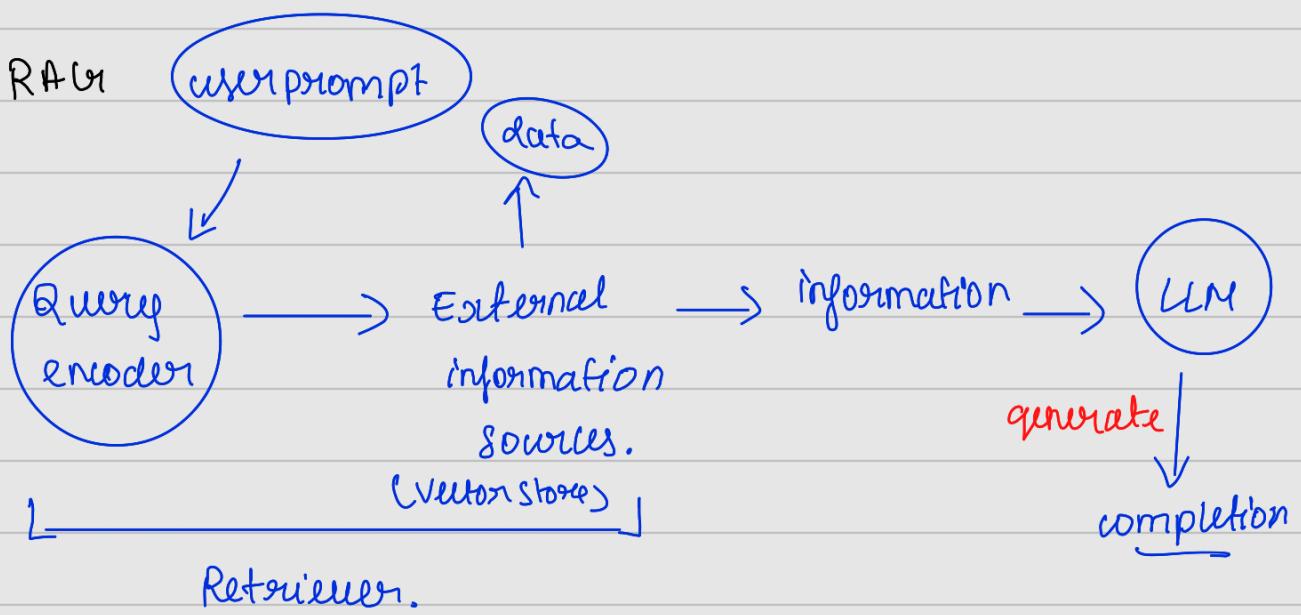
give LLM access to data

### ① Retrieval augmented generation (RAG) (framework)

LLM powered → use of external data sources.

→ knowledge cutoff issue. (very expensive, repeated training, regular updates)

→ for new data (external information) can improve both the relevance and accuracy of its completions.



- power of RAG
- ① Summarise  
② Places / people
- } with legal document
- avoids the model's hallucinating

RAG integrates with many types of data sources

## External information sources.

- (1) Documents
- (2) wiki's
- (3) Expert system
- (4) web pages
- (5) databases
- (6) vector store.

## Data preparation for vector store for RAG

- (1) Data must fit inside context window.

Prompt context | single document too large to fit in window → split long sources into short chunks  
limit few 1000 tokens | (langchain)

- (2) Data must be in format that allows its relevance to be assessed at inference time: Embedding vectors.

prompt text converted to embedding vector (each chunk) → allow the LLM to identify semantically related words (cosine similarity)

LLM create embedding vectors for each, new representations of the data can be stored in structures called vector stores.

vector databases are a particular implementation of a vector store, each vector is identified by key.

## Interacting with external application

external application interaction)

→ chatbots (chatting with customer)

→ LLM used to trigger action (API calls)

→ perform calculation

(interact with broader world)

beyond language task

Note prompts & completions are heart of these workflows.

LLM → reasoning engine

Requirement for using LLMs to power applications.

- ① plan Actions
- ② format output  
(completion)
- ③ validate actions  
(verification)

Prompt Structure is important

Helping LLMs reason and plan with chain of thought

→ complex reasoning for LLMs can be challenging for LLMs  
(mathematics) problems

Strategy

Humans take a step by step approach to solving  
complex problems.

by reasoning steps or chain of thoughts (How Human solve the  
problems)

This behaviour is known as chain of thought prompting  
includes intermediate steps

prompt → Model → completion  
(One-shot LLM)

thought of  
Prompting)

chain of thought prompting can help  
LLM's reason.

Performance  
of model  
is increased.

powerful techniques that improves the  
ability of model to reason through  
problems.

## Program-aided Language Models (PAL)

LLM can struggle with mathematics

the model isn't doing math, predict the most  
probable tokens that complete the prompt.

chains of → individual math  
thoughts operations wrong.

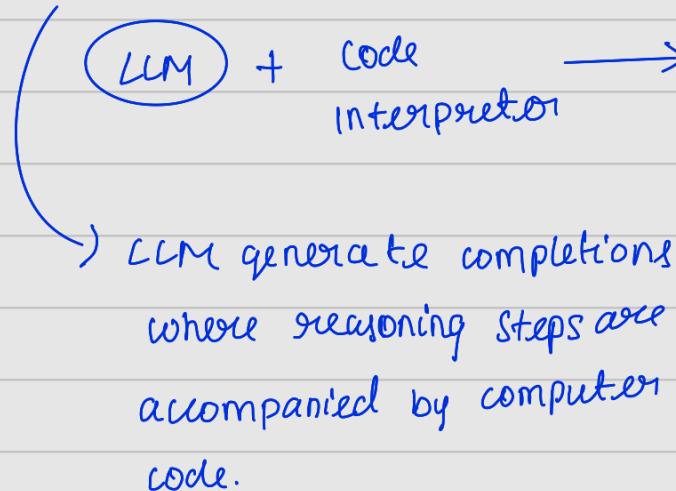
solution:

allowing your model to interact with external applications

that are good at math (Python interpreter)

One interesting framework for augmenting LLMs in this way  
is called program-aided language models. (PAL)

## PAL Models



methods makes use of chain of thought prompting to generate executable python scripts.

→ passed to interpreter to execute.

## PAL example

### 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?

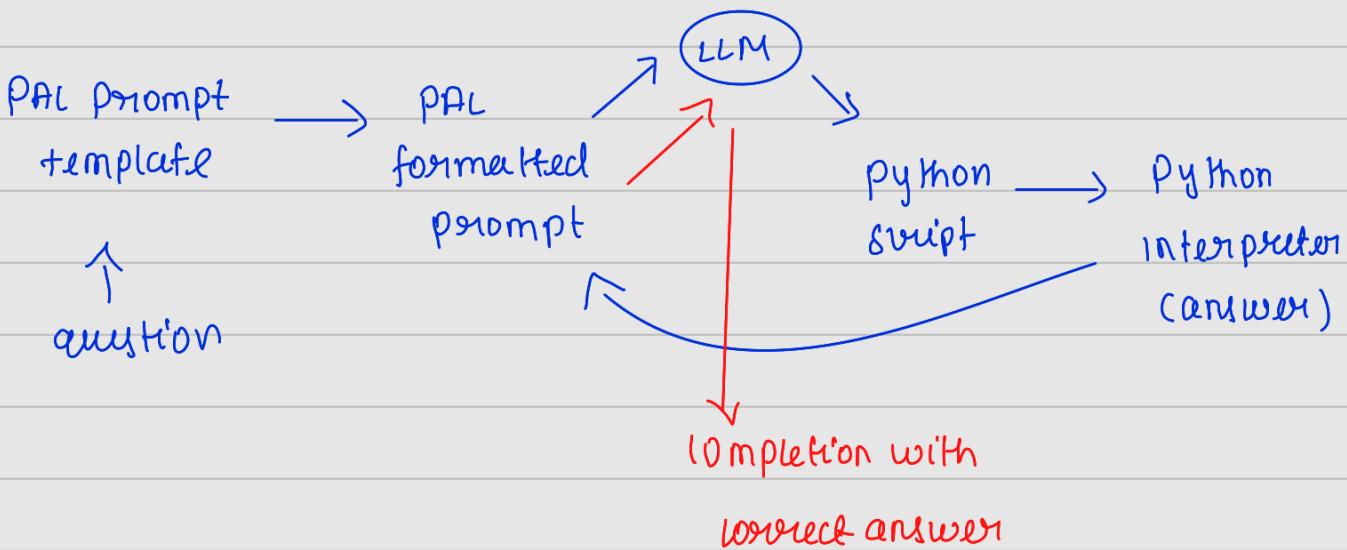
### Completion, CoT reasoning (blue), and PAL execution (pink)

Answer:

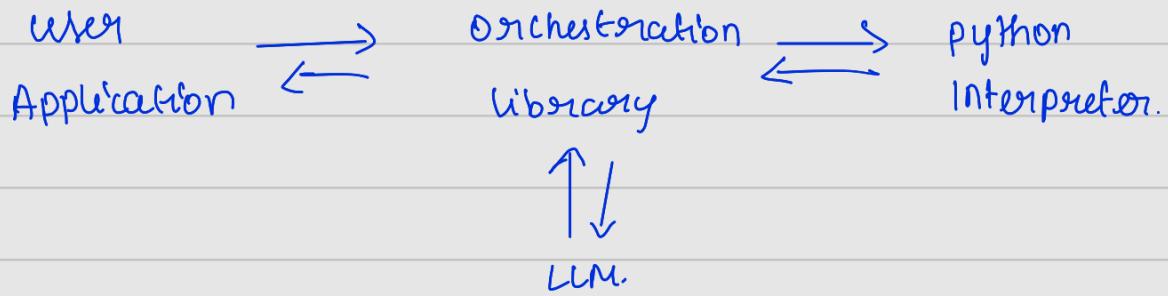
```
# The bakers started with 200 loaves
loaves_baked = 200
# They sold 93 in the morning and 39 in the
afternoon
loaves_sold_morning = 93
loaves_sold_afternoon = 39
# The grocery store returned 6 loaves.
loaves_returned = 6
# The answer is
answer = loaves_baked
    - loaves_sold_morning
    - loaves_sold_afternoon
    + loaves_returned
```

(python scripts based on reasoning steps)

## Program-aided language (PAL) models (complex calculation)



PAL architecture (it's simple, in real it's more complicated)



## Responsible AI

Challenges of responsible generative AI.

- (i) Toxicity
- (ii) Hallucinations
- (iii) Intellectual property.

### Toxicity

LLM returns responses that can be potentially harmful or discriminatory towards protected groups on protected attributes.

#### Solution:

- careful curation of training data
- Train generativl models to filter out unwanted content.
- diverse group of human annotators.

### ② Hallucination,

- LLM generates factually incorrect content.  
How to mitigate?
- Educate users how generative AI works
- Add disclaimers

- Augment LLMs with independent, verified citation databases.
- define intended/unintended use cases

### ③ Intellectual property

Ensure people aren't plagiarizing, make sure there aren't any copyright issues

How to mitigate?

- Mix of technology, policy, and legal mechanisms
- Machine "unlearning" (reduced/reduced content)
- Filtering and blocking approaches

Responsibility build 2 use generative AI models

- ① Define use cases : the more specific /narrow the better
- ② Assess risks for each use case
- ③ Evaluate performance for each use case
- ④ Iterate over entire AI life cycle. (concept & development stage)
- ⑤ Government policy.

### Conclusion

llama.cpp is a C++ implementation of the LLaMA model using four bit integer quantization to run on a Laptop.

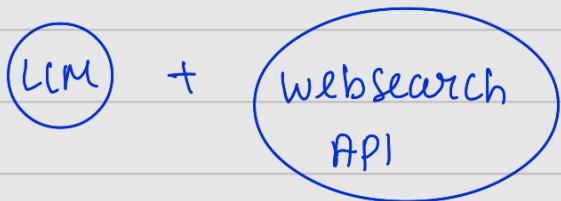
Model will also increasingly support multi-modality across (language (image, video, audio))

ReAct: combining reasoning and action

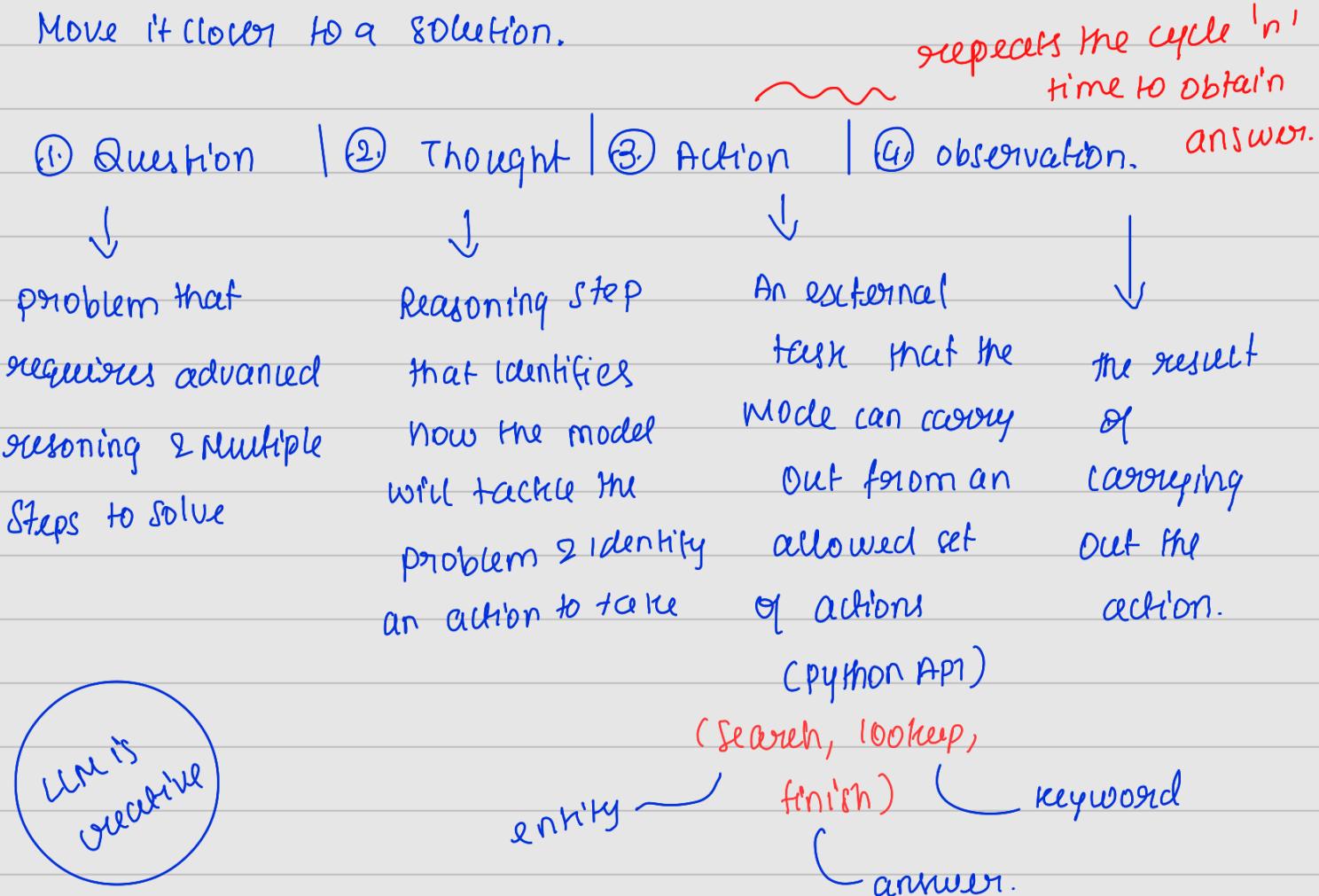
framework, that help LLMs plan out and execute their workflow,

ReAct: Synergizing Reasoning and Actions in LLM

ReAct's a prompting strategy that combines chain of thought reasoning with action planning.

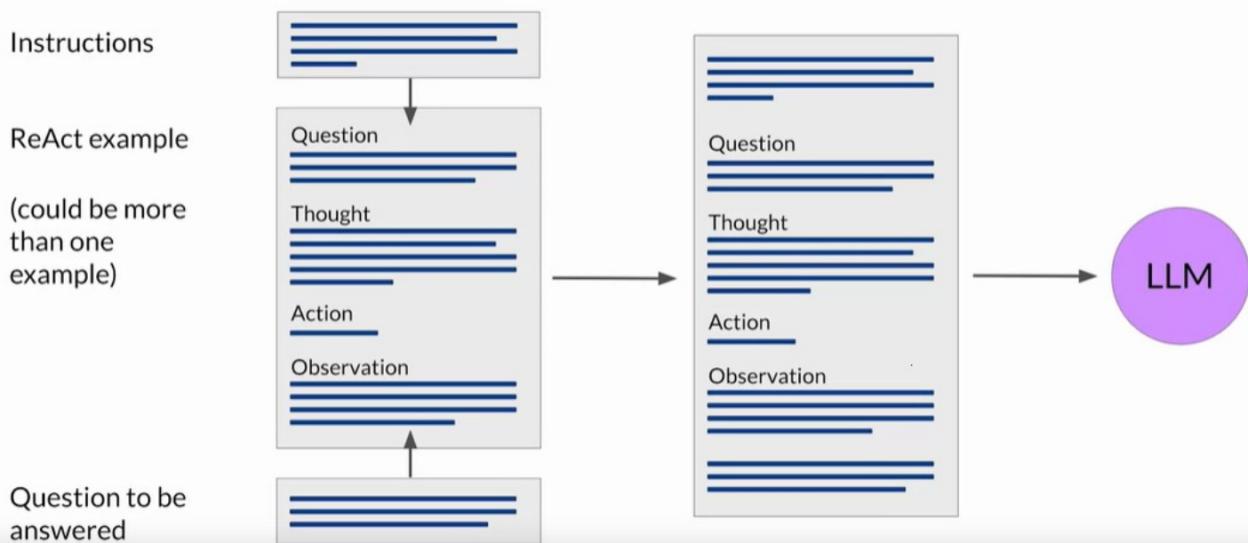


ReAct uses structured examples to show a LLM how to reason through a problem and decide on action to take  
Move it closer to a solution.

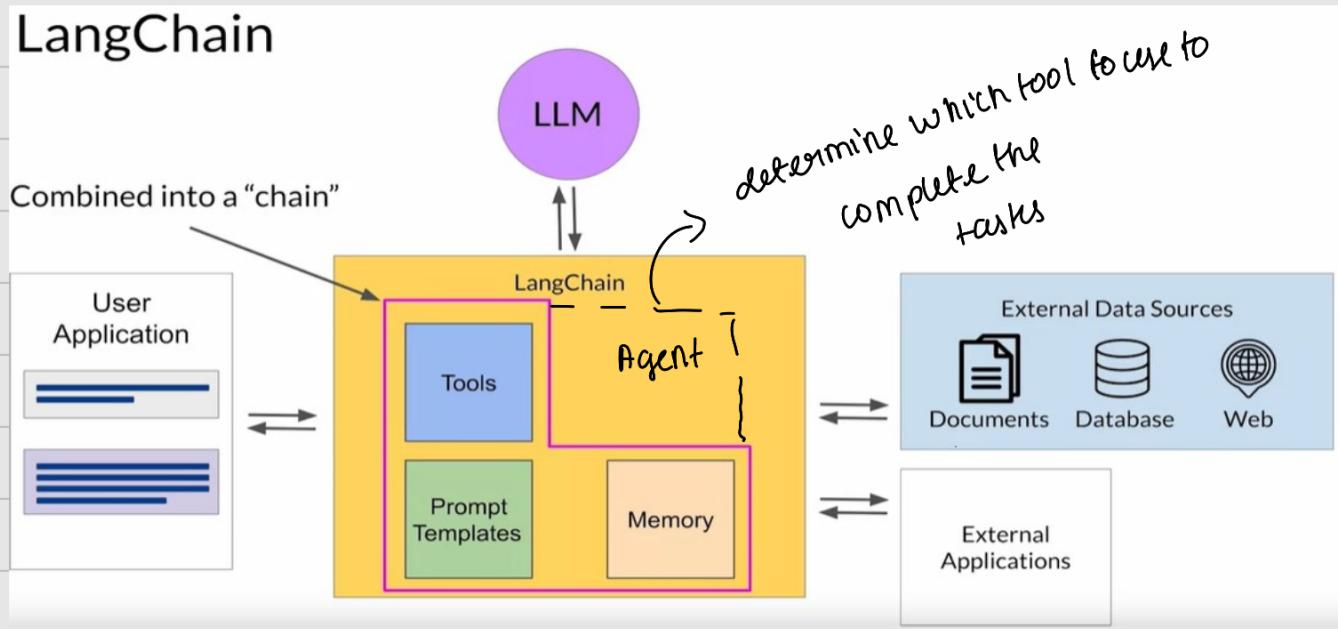


LLM can only choose from a limited number of actions that are defined by set of instruction.

## Building up the ReAct prompt



## LangChain



Start with large capable model, collects lots of user data in deployment and use it to train & fine tune smaller model.

# LLM powered application architecture

Building generative applications.

