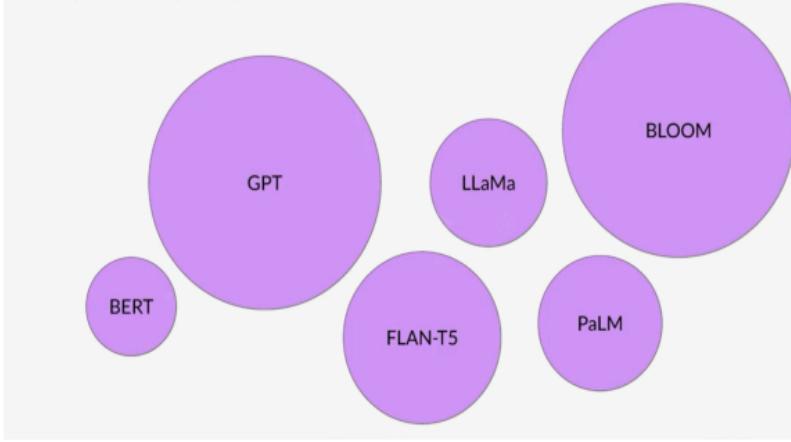


LARGE LANGUAGE MODEL

Notes by: Nikhil

Generative AI is a subset of traditional machine learning.

Large Language Models



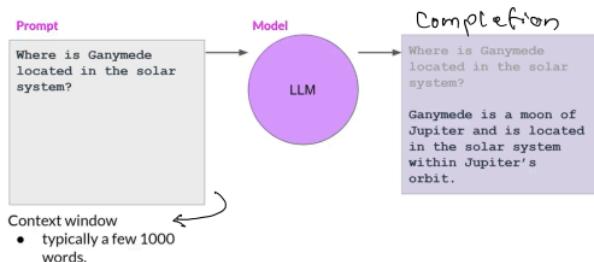
→ The models shown in the picture are the foundation models also called base models.

→ The size of the circle shows how big a model is.

→ The size of the model is based on parameters.

- The more parameter it has the more memory it has and can perform more sophisticated tasks.
- The text or instruction that we pass to the LLM is basically prompt.
- The text or instruction that we pass to the LLM is basically prompt.

Prompts and completions



→ The output of the model is called completion.

→ The act of using a model to generate output is called inference.

- LLM is not a very new concept the earlier way to predict or generate next words or sentences was achieved through RNN Recurrent Neural Network.
- But the limitation the RNN was with the compute power and only few previous words can't predict next word accurately.
- This limitation was overcome by **Transformer** introduced by google.

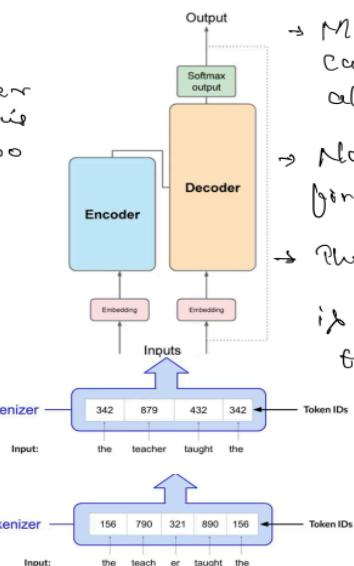
→ Now let's understand how transformer works :

Transformers

The transformer architecture is split into two distinct parts

- ↳ Encoder
- ↳ Decoder

Type - I



Type - II

→ Machine learning is basically a massive calculator and it understands numbers also called tokens but not words.

→ Now to work with this we need to first convert words into tokens

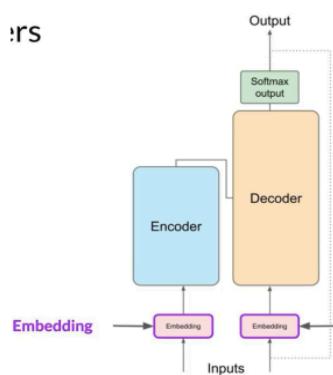
→ There are two ways to create tokens

if you can either use complete words to create token ID's

if you can use parts of words to create token ID's

→ Once the input is selected with the anyone of the token id's then it is going to be used throughout

!rs



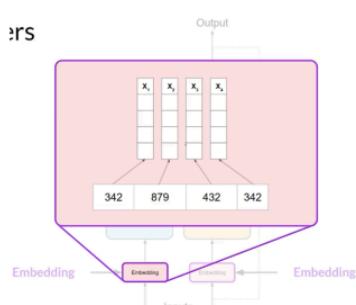
→ Once the input is converted into tokens then you can pass it into embedding layer.

→ This layer is a high trainable vector space and in this layer each token is represented as a vector and occupies a unique location in a space.

→ These vectors learn to understand the meaning and context of individual tokens in the sequence of input text.

→ Traditional NLP uses this concept like word2vec

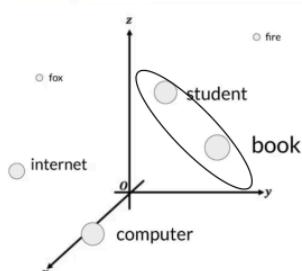
!rs

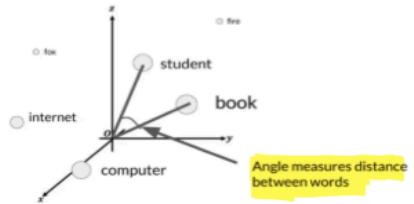


→ In this case you can see that each words is matched to a token ID and each token then mapped into vector.

→ In the original transformer paper the vector size is mentioned as 512.

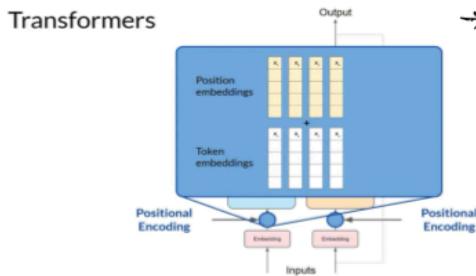
→ Now let's take an example if the vector size is 3 then we can plot the words into 3 dimensional space and understand the relationship based on how closely the words are placed into vector space.



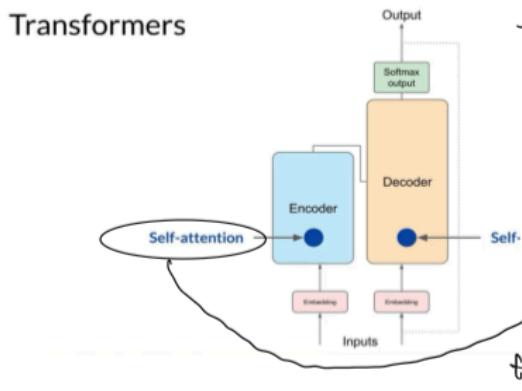


→ To calculate the distance between words by calculating the angle.

→ The model also adds positional encoding.



→ The model processes all the tokens parallelly and the use of position encoding helps the model to remember the position in the sentence.

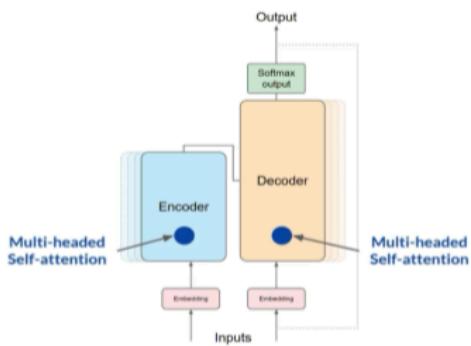


→ Once we combine the input token and the positional encoding then it passes to the self attention layer.

→ Here the model analyzes the relationship of input token into your input sentence.

→ The attention weights are learned during the training and stored into this layer reflects the importance of words into input sequence with all the words into the sentence.

→ The attention weights are learned during the training and stored into this layer reflects the importance of words into input sequence with all the words into the sentence.

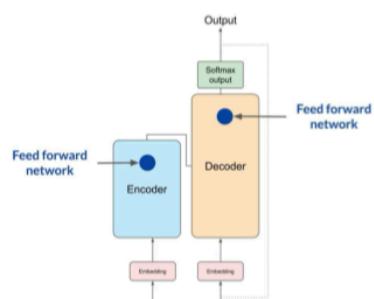


→ The learning of attention weights are not done only once but it creates multi-headed attention layer to learn the different weights parallelly and independently of each other.

→ And multi-headed self attention depends on model to model but ranging from 18-100

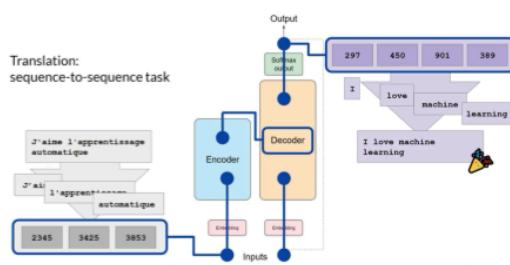
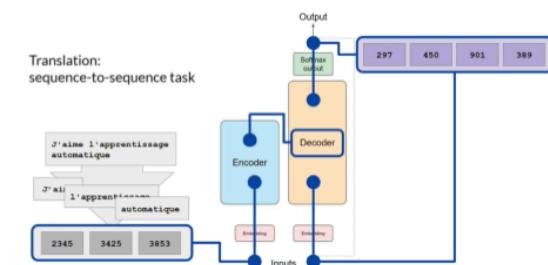
→ The idea of multi-headed is to learn every aspect of sentence like one layer learn about entity other layer learn about behaviour etc.

→ The weights of the layer assigned randomly and learn during training.



→ Now that we have weights learned during training then the input is processed through a fully connected feed forward network

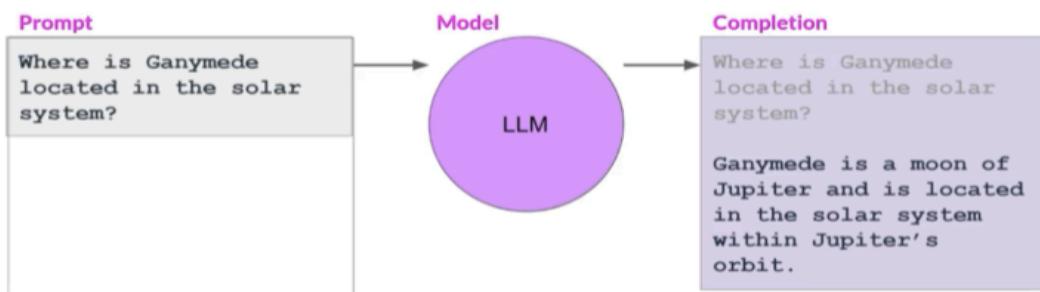
→ The output of this layer is a vector of logit proportional to probability of each word and then this logit is passed to the softmax layer.



- So far we have seen architecture of encoder and decoder model however we can also use encoder only model or decoder only model.
- Encoder only model is an example of seq-to-sequence model.
- BERT is an example of encoder only model, also used for sentiment analysis.

- The text that we feed into the model is called prompt
- The act of generating text is called inference and the output of the model is called completion
- The memory available to use the prompt is called context window.

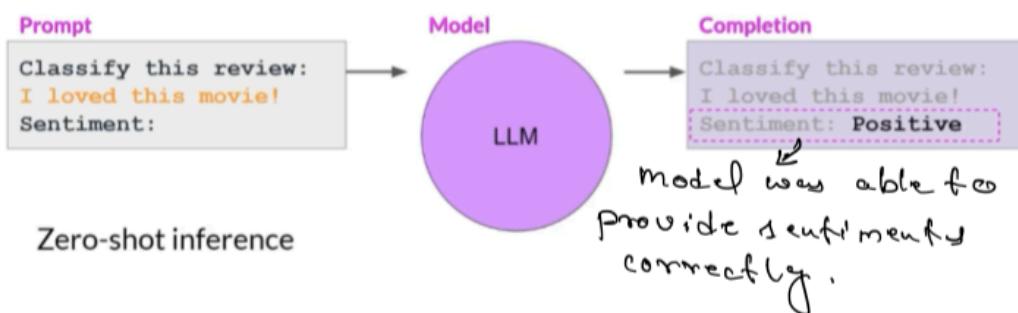
Prompting and prompt engineering



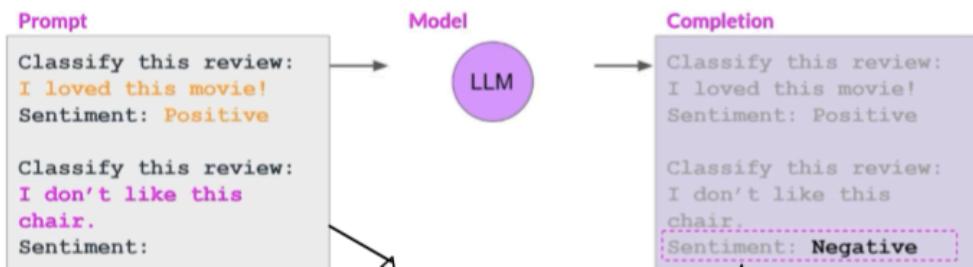
Context window: typically a few thousand words

- It's not necessary that the model would behave as per you on the first try but if you want this to happen you must provide prompt multiple times to get the desired output. This work to improve and develop model is known as prompt engineering.
- Providing example into the context window is called in-context learning.

In-context learning (ICL) - zero shot inference



In-context learning (ICL) - one shot inference



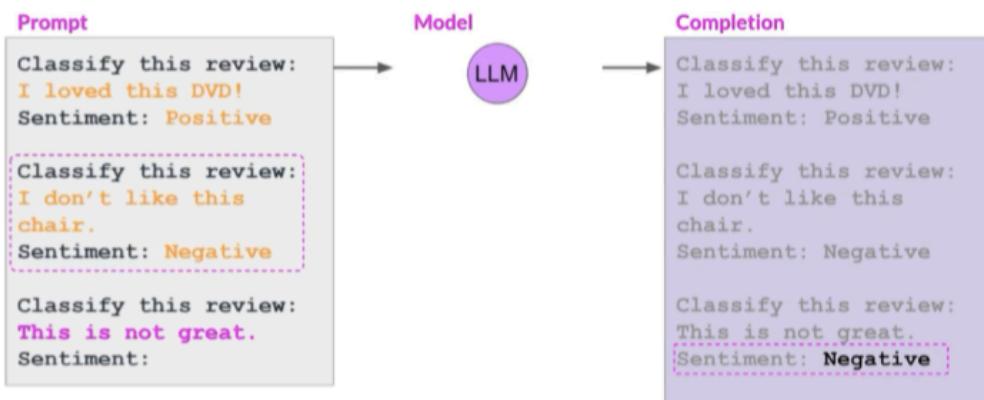
One-shot inference

In this example we have provided few example for the model to understand it more accurately.

The inclusion of single example into context window is known as One-shot inference.

- Sometimes single example won't be sufficient for the model to understand so in that case we can include multiple example and this is called few-shot inference.
- Larger model won't require one-shot or few-shot to produce better output however smaller model require one or few shot to produce completion accurately.

In-context learning (ICL) - few shot inference



Summary of in-context learning (ICL)

Prompt // Zero Shot	Prompt // One Shot	Prompt // Few Shot
Classify this review: I loved this movie! Sentiment:	Classify this review: I loved this movie! Sentiment: Positive Classify this review: I don't like this chair. Sentiment:	Classify this review: I loved this movie! Sentiment: Positive Classify this review: I don't like this chair. Sentiment: Negative Classify this review: Who would use this product? Sentiment:

→ The limitation of in-context learning is the memory size of context window if you realize that even after providing lot of example if your model is not providing better output in that case you have to fine-tune your model.

→ BERT was trained on 110M and BLOOM was trained on 176B parameters.

Generative configuration - inference parameters

But these inference parameters are different than the training parameters which was learnt during training,

Enter your prompt here...

Max new tokens: 200
Sample top K: 25
Sample top P: 1
Temperature: 0.8
Submit

Inference parameters are influenced during inference to provide better output.
Inference configuration parameters

Generative configuration - inference parameters

Enter your prompt here...

Inference configuration parameters

Max new tokens	200
Sample top K	25
Sample top P	1
Temperature	0.8

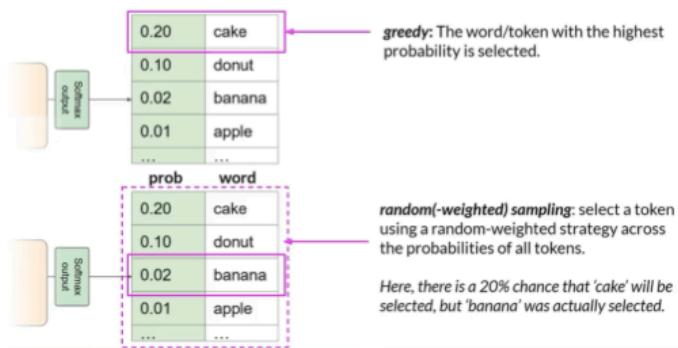
Submit

* max new tokens
basically gives you maximum words/tokens in the completion.

* This can be used

to limit the tokens in output or you can think this as putting a cap on it.

Generative config - greedy vs. random sampling



→ The output of softmax layer is the probability distribution across the entire dictionary of words that the model chose.

→ Most LLM operates by default on greedy sampling in order to predict new words this is the simplest technique where the model chooses the word with highest probability.

→ The second approach to introduce the variability is random sampling to select more probable words with no repetition of words.

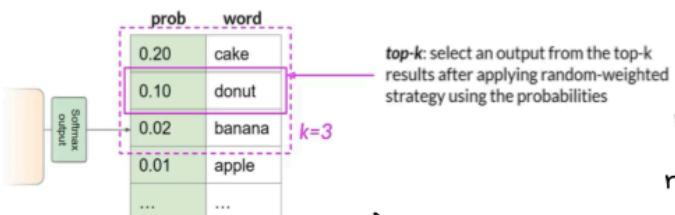
→ Based on the requirement either we need to disable greedy and enable sampling.

Generative configuration - inference parameters



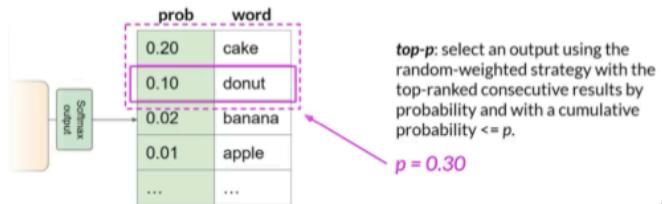
Pop k and top P are two sampling technique that help limit the random sampling and produce more sensible.

Generative config - top-k sampling



→ In this example you are setting $k=3$ which means you are restricting model to select between the top 3 words and the probable selection is "donut" here. By this method model has some randomness while avoiding to choose from more improper selection.

Generative config - top-p sampling

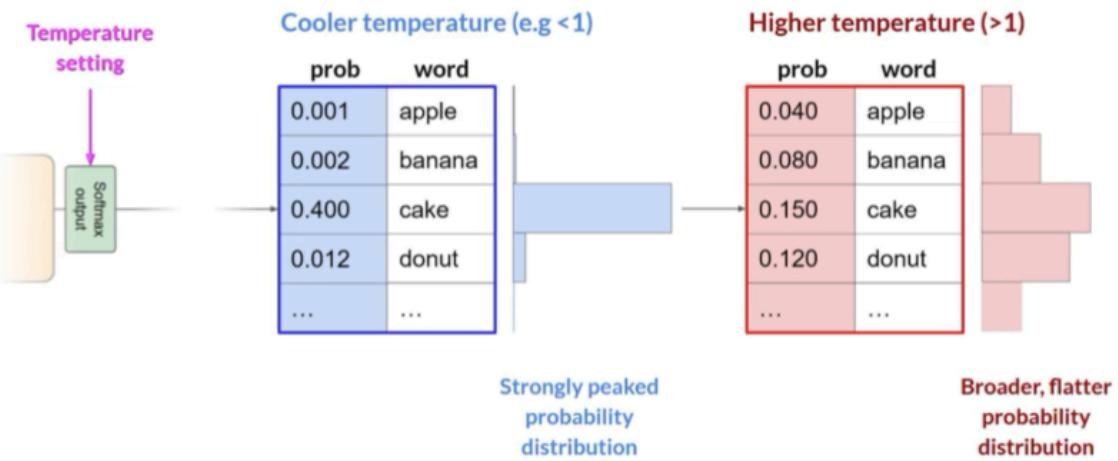


→ The other method is to set the value of p and based on that the words are randomly selected.

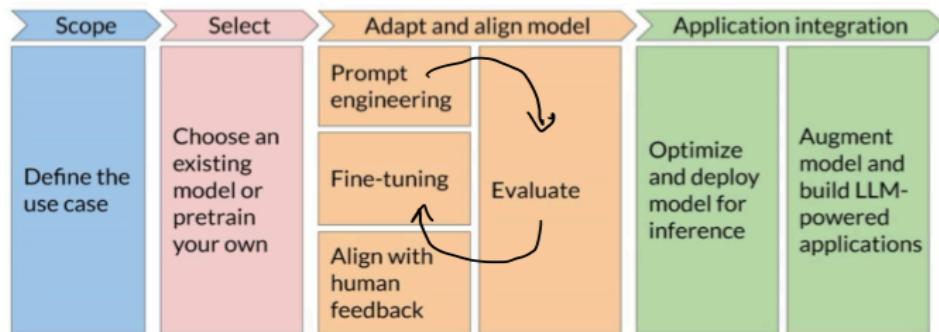
→ In this case the probability is set to 0.30 so the model has selected two words that fall within the probability and then it perform random selection.

→ One more parameter is **temperature** which control the randomness means higher the temperature higher the randomness and lower the temperature lower the randomness.

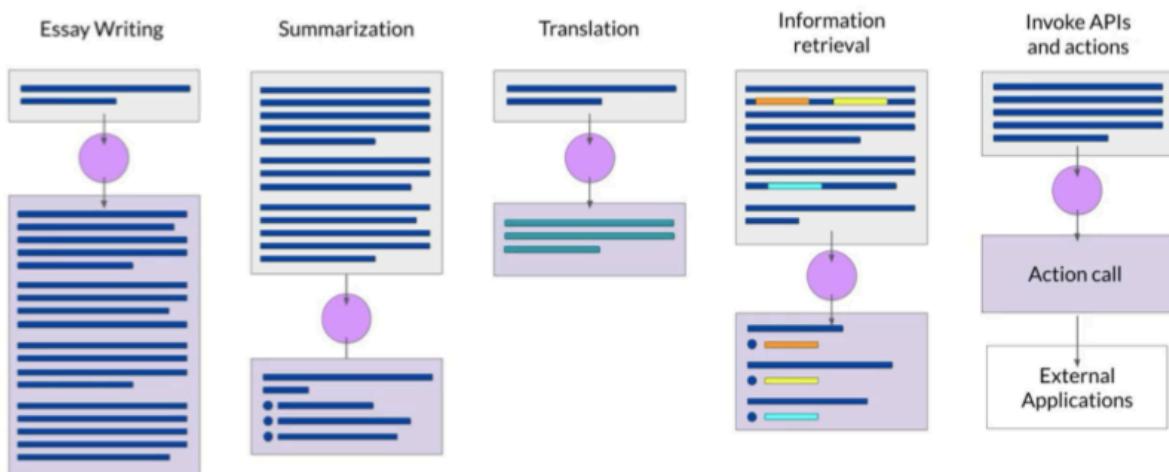
Generative config - temperature



Generative AI project lifecycle



→ The first part is to define or decide on the use case means what is the problem you want to solve. Below are the examples of multiple use cases :-



- The second part is to check if we want to train our own model from scratch or we can use the existing model. In general we use existing model.
- 3rd part is to do the assessment of the model. So that we can fine-tune the model like some time prompt engineering is sufficient enough to get sufficient result so in that case we can learn more about in-context learning.
- If not improved then we start fine-tuning the model.

* Autoencoding models.

→ Good cases.

↳ Sentiment analysis,

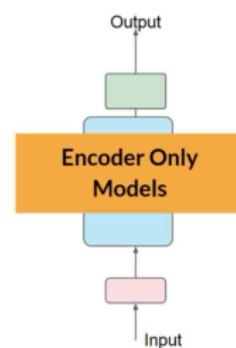
↳ Named Entity recognition,

↳ Word classification,

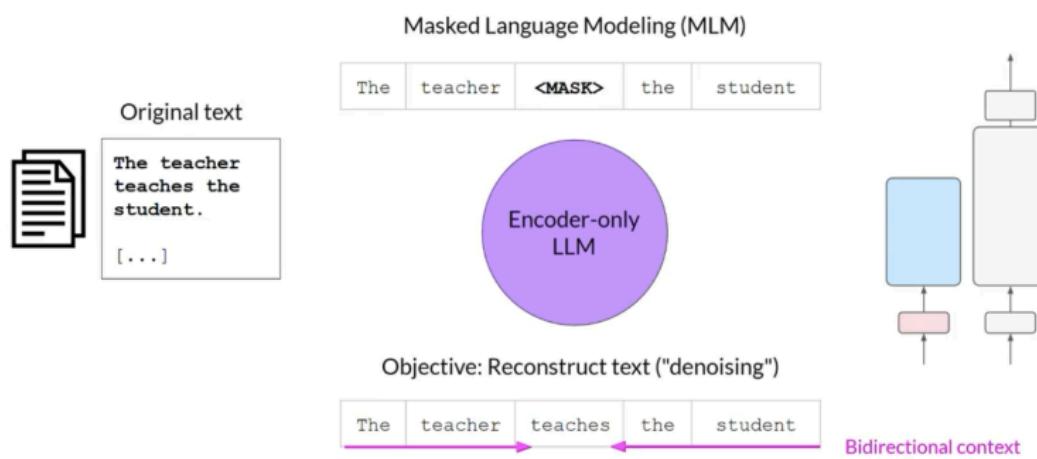
Example models ↳

↳ BERT

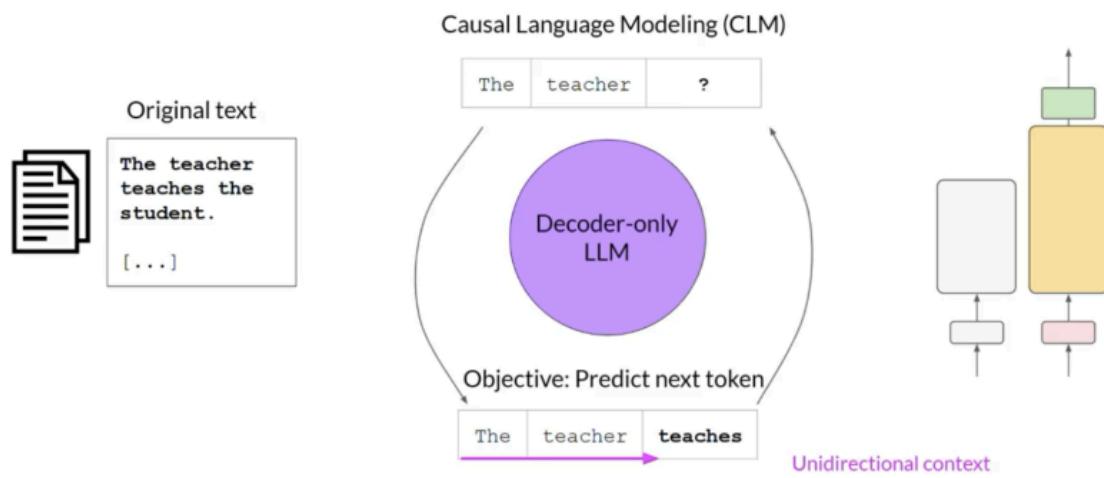
↳ ROBERTA.



Autoencoding models

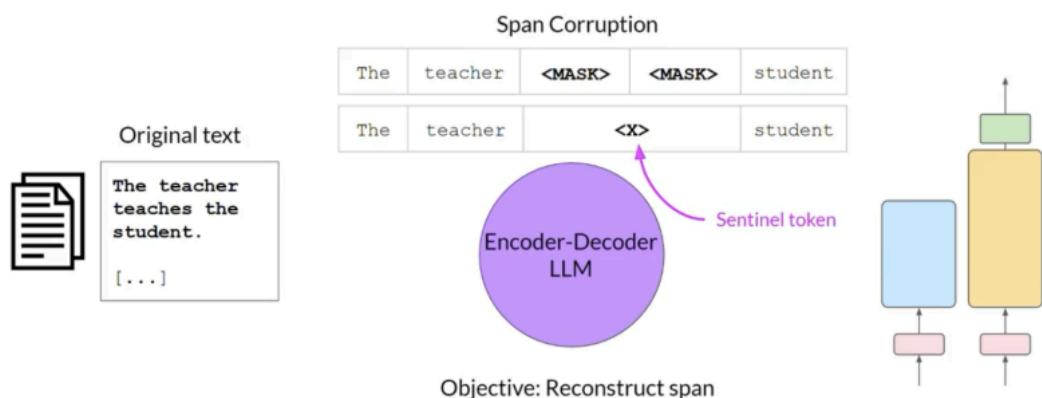


Autoregressive models



These models are good for text generation and few models are GPT and BLOOM.

Sequence-to-sequence models

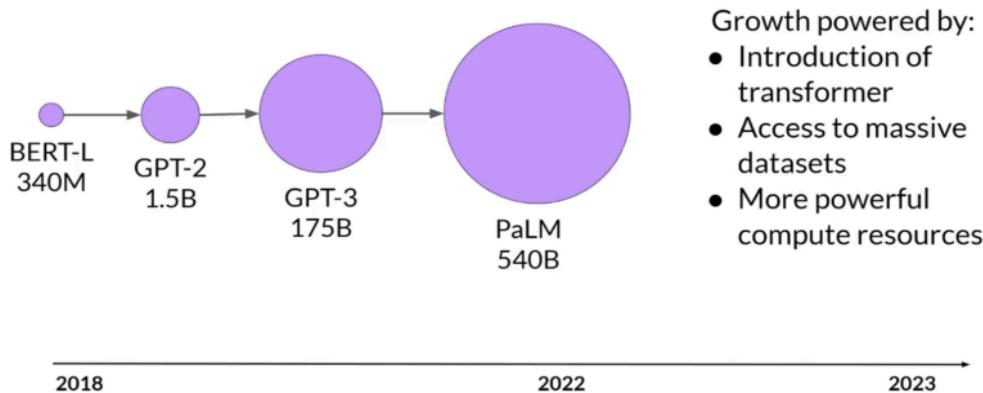


- The final is sequence-to-sequence that uses both encoder and decoder model.
- Good use cases :-
 - ↳ Translation
 - ↳ Text summarization,
 - ↳ Question answering,

Example in T5 and BART.

→ Larger the model it is more likely to work best with use case without any context learning and further fine-tuning.

Model size vs. time



→ The common challenge that anyone can face while training a large language model is running out of time.

Out Of Memory Error : 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



Approximate GPU RAM needed to train 1B-params

Memory needed to store model



4GB @ 32-bit
full precision

Memory needed to train model

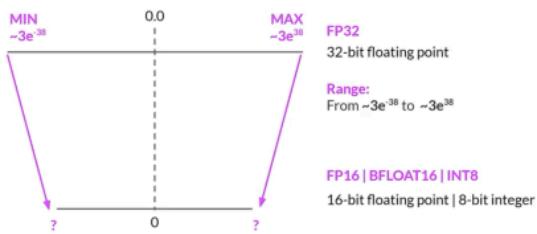


→ What are the options to reduce the memory required to train the model?

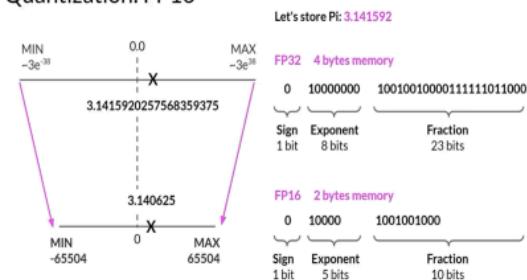
↳ Quantization.

→ By reducing their precision from 32-bit to 16-bit floating point.

Quantization



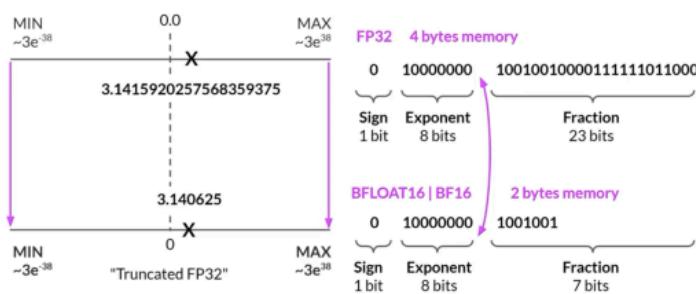
Quantization: FP16



→ It will reduce the memory requirement by half.

→ One alternative recently developed to replace FP16 is BFLOAT16 which is Brain Float16 developed by Google and is used in PyTorch training.

Quantization: BFLOAT16

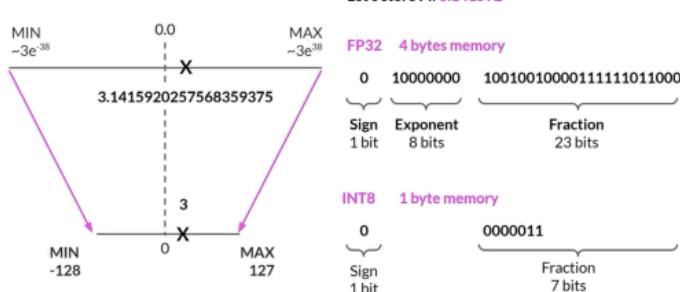


→ It uses 8-bit exponent but truncates the fraction to 7-bits which is why it is called truncated FP32.

→ Downside it can't be used for integer types.

→ To deal with that we have DNPB but it would compromise with precision to reduce the memory requirement to train the model.

Quantization: INT8



Quantization: Summary

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

FLAN
T5

- 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

Approximate GPU RAM needed to store 1B parameters Approximate GPU RAM needed to train 1B-params



Pre-training for domain adaptation

Legal language

The prosecutor had difficulty proving mens rea, as the defendant seemed unaware that his actions were illegal.

The judge dismissed the case, citing the principle of res judicata as the issue had already been decided in a previous trial.

Despite the signed agreement, the contract was invalid as there was no consideration exchanged between the parties.

Medical language

After a strenuous workout, the patient experienced severe myalgia that lasted for several days.

After the biopsy, the doctor confirmed that the tumor was malignant and recommended immediate treatment.

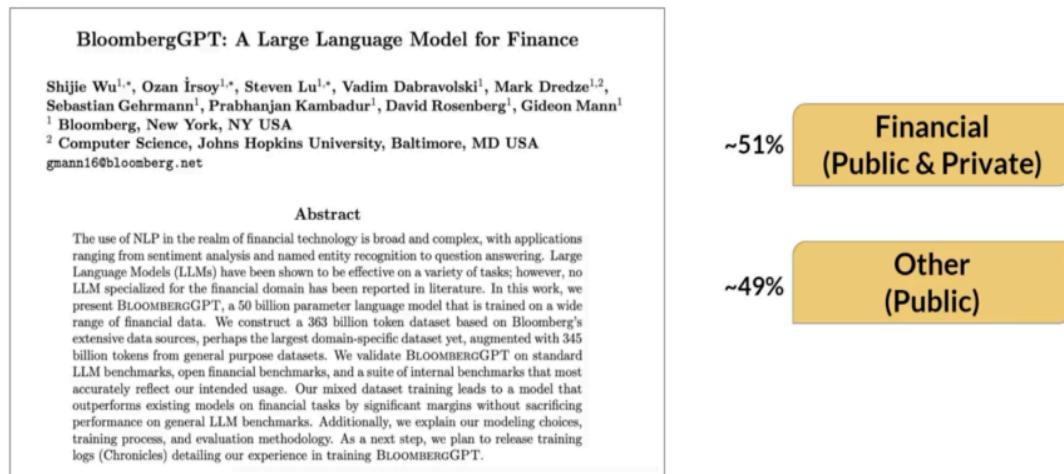
Sig: 1 tab po qid pc & hs



Take one tablet by mouth four times a day, after meals, and at bedtime.

- Pre training a model for domain specific would provide you more accurate result. like law, medicine, finance etc.
- Bloomberg GPT is specially trained for finance.

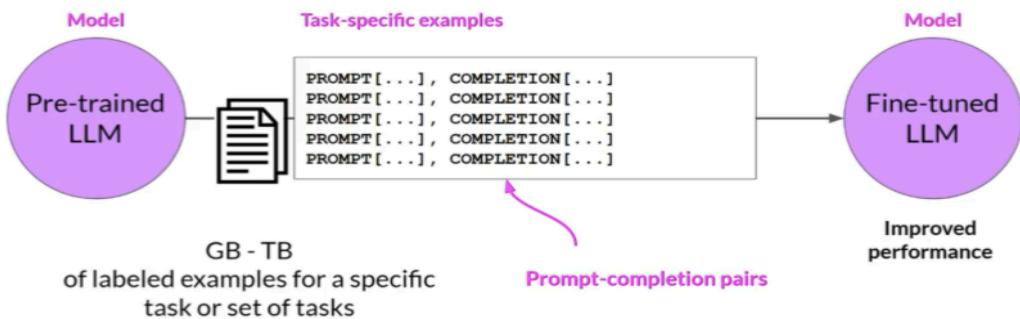
BloombergGPT: domain adaptation for finance



- * **Fine tune LLM with Instruction prompt.**
 - Fine tune is important as in-context learning has some limitation.
 - ↳ In context learning doesn't work with smaller model even if we include multiple example in-context.
 - ↳ The second is it takes important space from context window which limit us to include important information.
 - Fine tuning is a supervised process where we use labeled data to update the weight of LLM.
 - The labeled examples are prompt completion pairs.

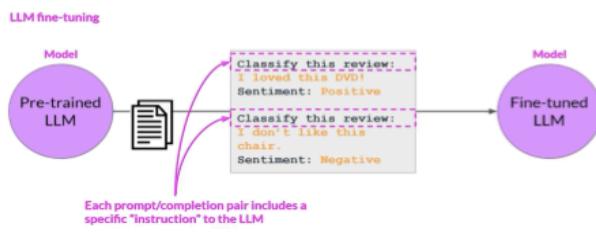
LLM fine-tuning at a high level

LLM fine-tuning



- Fine tuning process extends model to provide improved performance. Our strategy is to use instruction fine tuning.

Using prompts to fine-tune LLMs with instruction



→ The instruction - fine tuning needs a dataset that contains prompt and completion pairs to fine tune a model.

- Instruction fine tuning where all weight of LLM gets updated is known as full-fine tuning. This process result in new version of model with the updated weights.

- Just like pre-training, instruction fine tuning required enough memory to fine-tune gradients, optimizer and other components of the model.

- * Steps in instruction fine tuning,

Step 1: Prepare your training data. But this we need to create a dataset of prompt instruction temp. For that developers have created prompt template libraries that will take huge data and convert into prompt instruction template.

Sample prompt instruction templates

Classification / sentiment analysis

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

In this we would pass original review to the template.

→ In this we would pass original review to the template, → This is true.

Text generation

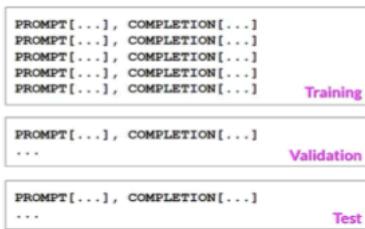
```
jinja: Generate a {{star_rating}}-star review (1 being lowest and 5 being highest)  
about this product {{product_title}}. ||| {{review_body}}
```

Text summarization

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

Prepared instruction dataset

Training splits



→ Once the dataset is ready then we divide the dataset into train, test and validation set.

LLM fine-tuning

Prepared instruction dataset



Model



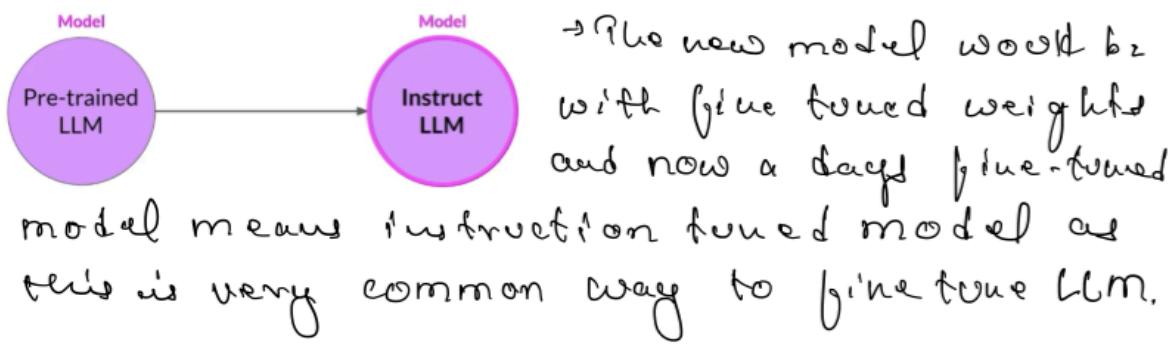
→ During fine-tuning we take prompt from training dataset and pass it through LLM which generates output called completion.

→ Then we compare the response with the training data to check the accuracy. We can see the response is "Neutral" however the label from training is "Positive". Not a good outcome.

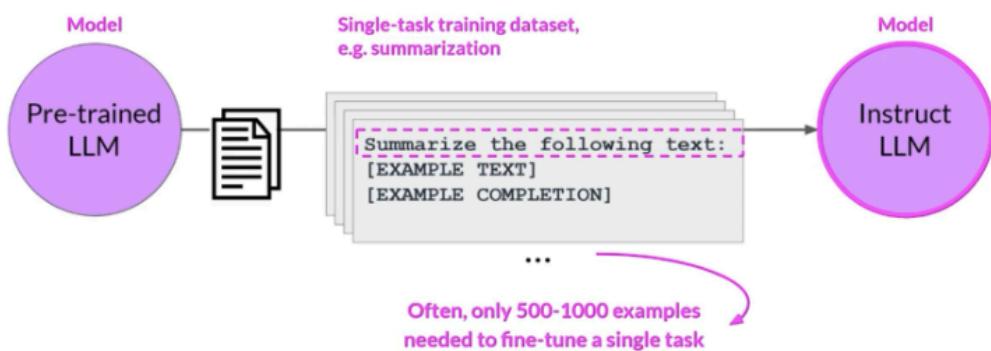
→ The output of a CCM is a probability distribution across tokens.

→ We would use cross-entropy method to calculate the loss between actual and predicted. Then we would calculate loss to update the weights in back propagation, for many batches and over several epochs to improve the performance.

→ Then check with hold-out validation dataset and then finally with test dataset.



Fine-tuning on a single task



→ The potential downside to fine-tune on a single task. The process can lead to the phenomenon called catastrophic forgetting.

→ The catastrophic forgetting happens because the full fine tuning modifying the weights of original LLM model.

→ This would improve performance for that particular task but degrade the performance on other tasks.

Before fine-tuning

Prompt

Classify this review:
I loved this DVD!
Sentiment:

Model



Completion

Classify this review:
I loved this DVD!
Sentiment: loved a
very nice book review

Before fine-tuning

Prompt

What is the name of
the cat?
Charlie the cat roamed
the garden at night.

Model



Completion

What is the name of
the cat?
Charlie the cat roamed
the garden at night.
Charlie

- Before fine-tuning the model is easily carried out named-entity-recognition like predicting "Charlie"

After fine-tuning

Prompt

What is the name of
the cat?
Charlie the cat roamed
the garden at night.

Model



Completion

What is the name of
the cat?
Charlie the cat roamed
the garden at night.
The garden was
positive.

- Post fine tuning the model is not able to perform well and confuse with the new task "The garden was positive"

* How to avoid Catastrophic forgetting.

Step 1:

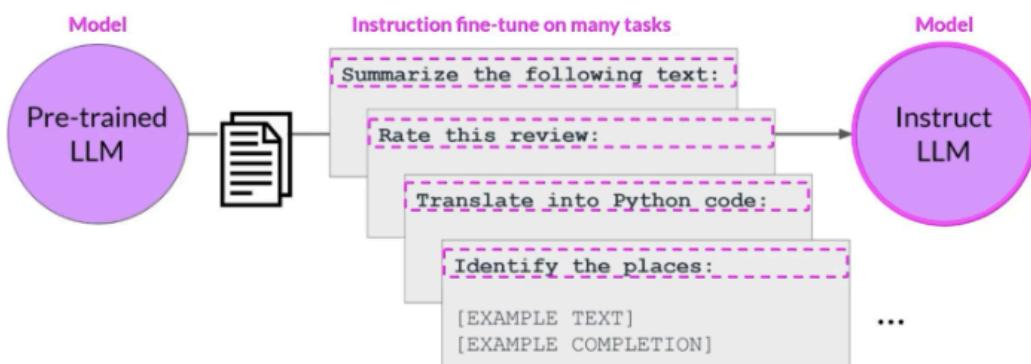
- It is not necessary to avoid as first check if your model actually not performing well on other's task. If not then it is not required.

- If your model is not able to perform well on other's task then to avoid one should fine tune the model on multiple tasks rather than one.

- It would take 800 - 1,000,000 data records all individual case to fine tune model on general task at the same time.
- Step 2:
- Consider Parameter Efficient fine-tuning (PEFT)
 - PEFT preserves the original weight of LLM and fine tune on adaptive layer for smaller task specific.

* Multi-task instruction fine-tuning.

Multi-task, instruction fine-tuning

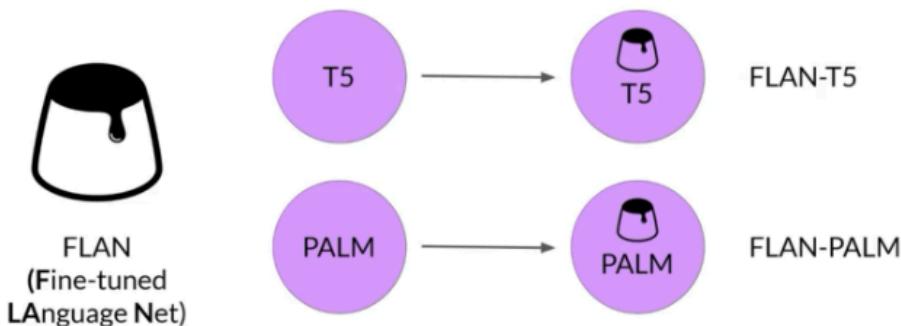


- Rather than fine tuning on one task in this we fine-tune the model on multiple task simultaneously as show above to avoid catastrophic forgetting.

→ Drawback is we need a lot of data to fine-tune a model 500 - 100000 datasets.

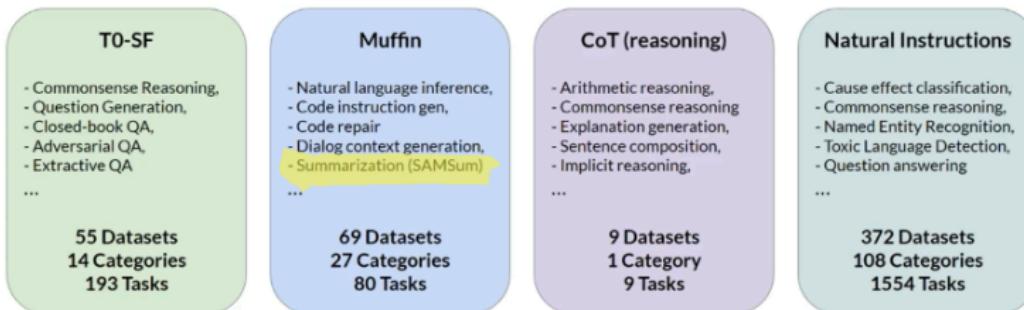
Instruction fine-tuning with FLAN

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



FLAN-T5: Fine-tuned version of pre-trained T5 model

- FLAN-T5 is a great, general purpose, instruct model



→ One example of prompt dataset used for summarization task in **FLAN-T5** is **SAMSUM**.

→ It's a part of muffin task and dataset used for summarization dialogue task.

SAMSum: A dialogue dataset

Sample prompt training dataset (**samsum**) to fine-tune FLAN-T5 from pretrained T5

Datasets: samsum	Tasks:	Summarization	Languages:	English
dialogue (string)	summary (string)			
"Amanda: I baked cookies. Do you want some? Jerry: Sure! Amanda: I'll bring you tomorrow :-)"	"Amanda baked cookies and will bring Jerry some tomorrow."			
"Olivia: Who are you voting for in this election? Oliver: Liberals as always. Olivia: Me too!! Oliver: Great"	"Olivia and Olivier are voting for liberals in this election."			
"Tim: Hi, what's up? Kim: Bad mood tbh, I was going to do lots of stuff but ended up procrastinating Tim: What did...	"Kim may try the pomodoro technique recommended by Tim to get more stuff done."			

Sample FLAN-T5 prompt templates

```
"samsum": [  
    ("{dialogue}\n\nBriefly summarize that dialogue.", "{summary}"),  
    ("Here is a dialogue:\n{dialogue}\n\nWrite a short summary!",  
     "{summary}"),  
    ("Dialogue:\n{dialogue}\n\nWhat is a summary of this dialogue?",  
     "{summary}"),  
    ("{dialogue}\n\nWhat was that dialogue about, in two sentences or less?",  
     "{summary}"),  
    ("Here is a dialogue:\n{dialogue}\n\nWhat were they talking about?",  
     "{summary}"),  
    ("Dialogue:\n{dialogue}\n\nWhat were the main points in that "  
     "conversation?", "{summary}"),  
    ("Dialogue:\n{dialogue}\n\nWhat was going on in that conversation?",  
     "{summary}")]
```

→ Same instruction told in different ways helps the model generalize better.

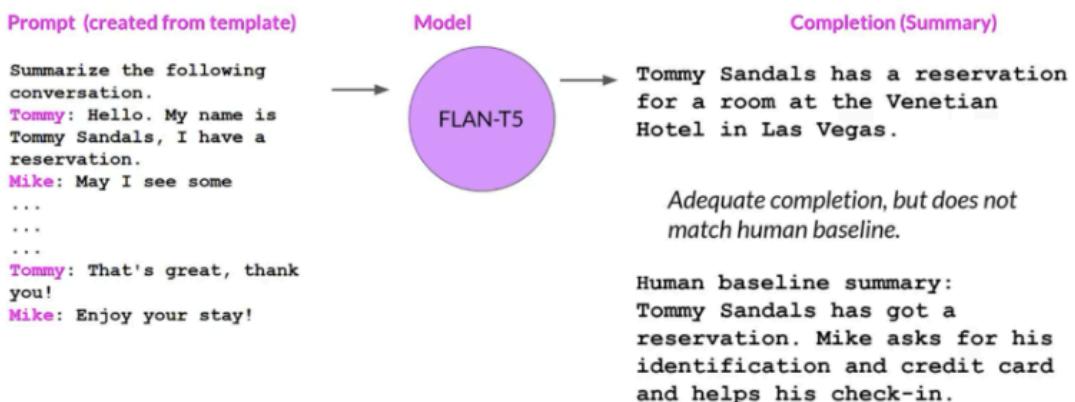
Instruction

Prompt

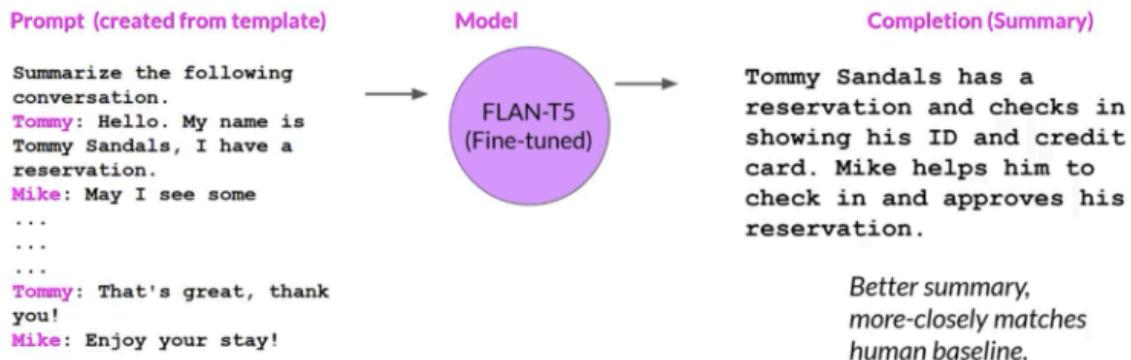
Label

→ FLAN-T5 is very popular generalize model however there are rooms of improvement.

Summary before fine-tuning FLAN-T5 with our dataset



Summary after fine-tuning FLAN-T5 with our dataset



→ FLAN-T5 model performed more accurate post fine tuning.

Model evaluation metrics

→ In traditional ML calculating or evaluating the much more realistic because we have the known output with which we can calculate, known formulae.

$$\text{Accuracy} = \frac{\text{Correct Predictions}}{\text{Total Predictions}}$$

→ However evaluating LLM models have challenges as the generated outcome are unknown or non-deterministic.

"Mike really loves drinking tea."



=



"Mike adores sipping tea."

"Mike does not drink coffee."



≠



"Mike does drink coffee."

- First example where both the sentences have same meaning.
- 2nd example only one word is different but both are not same.

LLM Evaluation - Metrics



- Used for text summarization
- Compares a summary to one or more reference summaries
- Used for text translation
- Compares to human-generated translations

- These are the two LLM evaluation metrics

ROUGE → Recall Oriented Understudy for getting evaluation compare the quality of automatically generated text with human generated.

- **BLEU** → Bilingual evaluation understudy is an algorithm to check the quality of machine translated text with human translated text.

LLM Evaluation - Metrics - ROUGE-1

Reference (human):

It is cold outside.

Generated output:

It is very cold outside.

$$\text{ROUGE-1} = \frac{\text{unigram matches}}{\text{unigrams in reference}} = \frac{4}{4} = 1.0$$

$$\text{Precision: } \text{ROUGE-1} = \frac{\text{unigram matches}}{\text{unigrams in output}} = \frac{4}{5} = 0.8$$

$$\text{F1: } \text{ROUGE-1} = 2 \frac{\text{precision} \times \text{recall}}{\text{precision} + \text{recall}} = 2 \frac{0.8}{1.8} = 0.89$$

- Let's talk about first Evaluation Metrics ROUGE 1
- Just like typical machine learning we can calculate Recall, Precision and F1 score to get the performance of model.
- For example we have human as (Reference word) and outcome from model as generated one.
- We can see as all the 4 words from reference is presented into generated resulting recall = 1
- Similarly we have 4 out of 5 is matching resulting precision = 0.8.
- But imagine a case instead of "very" the generated one would have "not" into output. The recall and precision would be same however the meaning is completely different.

Reference (human):
It is cold outside.

Generated output:

It is not cold outside.

$$\text{ROUGE-1} = \frac{\text{unigram matches}}{\text{unigrams in reference}} = \frac{4}{4} = 1.0$$

$$\text{Precision: } \frac{\text{unigram matches}}{\text{unigrams in output}} = \frac{4}{5} = 0.8$$

$$\text{F1: } \text{ROUGE-1} = 2 \frac{\text{precision} \times \text{recall}}{\text{precision} + \text{recall}} = 2 \frac{0.8}{1.8} = 0.89$$

- To avoid this situation a better approach would be using bigrams (means group of 2 words). By using bigrams it would be difficult to fall into condition like we saw in ROUGE-1. When we use bigrams in evaluation we take this as ROUGE-2.

LLM Evaluation - Metrics - ROUGE-2

Reference (human):

It is cold outside.

It is is cold
cold outside

Generated output:

It is very cold outside.

It is is very
very cold cold outside

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

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

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

LLM Evaluation - Metrics - ROUGE-L

Reference (human):

It is cold outside.

Generated output:

It is very cold outside.

Longest common subsequence (LCS):

It is
cold outside

2

→ Rather than using bigram, trigram or n-gram as the length increases, let's check another metric called ROUGE-L,

→ In this we calculate using LCS which is longest common subsequence, which is highlighted into red box,

→ We can now use LCS to calculate precision and recall. We would use LCS in both precision and recall.

Reference (human):

It is cold outside.

Generated output:

It is very cold outside.

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

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

LCS:

Longest common subsequence

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

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

The BLEU score quantifies the quality of a translation by checking how many n-grams in the machine-generated translation match those in the reference translation. To calculate the score, you average precision across a range of different n-gram sizes. If you were to calculate this by hand, you would carry out multiple calculations and then average all of the results to find the BLEU score.

Parameter efficient fine-tuning (PEFT)



→ PEFT doesn't change the weight of entire Lm.

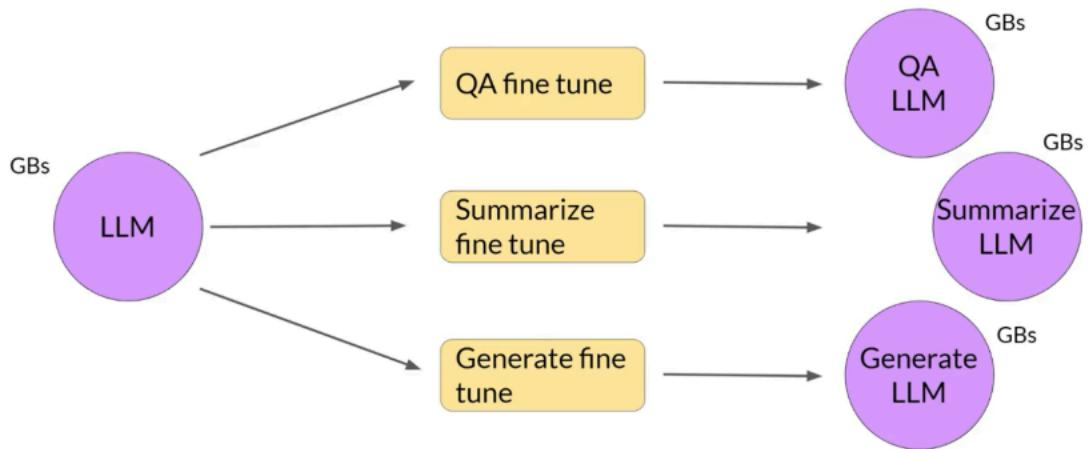
→ Some PEFT technique freeze the most layer of Lm and update the weights of small layer or Parameters.

→ However some PEFT technique doesn't even touch the existing layers rather than it add new layer to update weights of Lm.

→ PEFT can often run on single GPU.

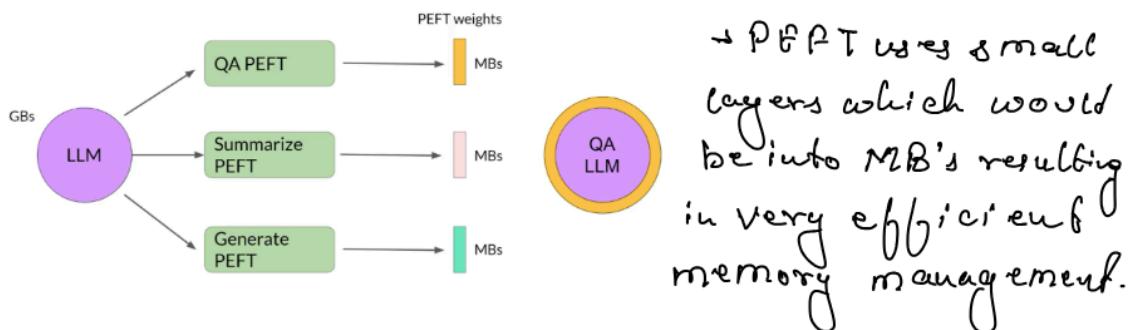
→ PEFT also eliminates catastrophic forgetting which is a common phenomenon of full-fine tuning.

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



- Full fine-tuning creates multiple copy of original LLM resulting in difficult for memory management.
- Let's see how we can use PEFT to avoid situation of memory management.

PEFT fine-tuning saves space and is flexible



PEFT Trade-offs

Parameter Efficiency

Memory Efficiency

Training Speed

Model Performance

Inference Costs

PEFT methods

Selective

Select subset of initial LLM parameters to fine-tune

Reparameterization

Reparameterize model weights using a low-rank representation

LoRA

Additive

Add trainable layers or parameters to model

Adapters

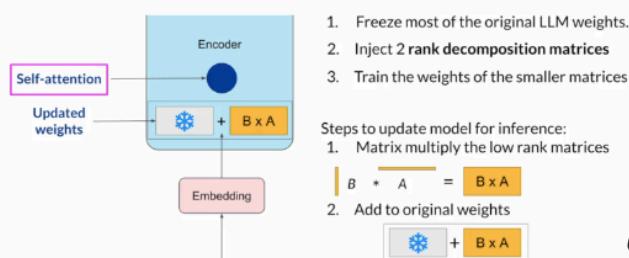
Soft Prompts

Prompt Tuning

Low-Rank Adaptation of Large Language Models (LoRA)

→ LoRA is a PEFT technique which falls into reparameterization method of PEFT.

LoRA: Low Rank Adaption of LLMs



→ Let's talk about transformers first
→ The input is converted to tokens.
→ Then the tokens are converted to embeddings.

layers which is connected to encoder and decoder of the transformer,

→ In encoder and decoder layer we have a neural network called self attention layer and feed-forward network.

→ This is a layer where the weight of the neural network updated during training.

→ In case of full fine tuning all the weight of layers get updated.

- LoRA helps this by freezing original layers of LLM and update the small part of network layer.
- Then itself smaller dimension matrix. Then it train the smaller matrix using supervised technique and then it added with the original weights kept frozen.
- The researchers has confirmed that updating LoRA only on self-attention layers helps in achieving any task. However it can also be used for other components like feed-forward layer.

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

In LoRA with rank $r = 8$:

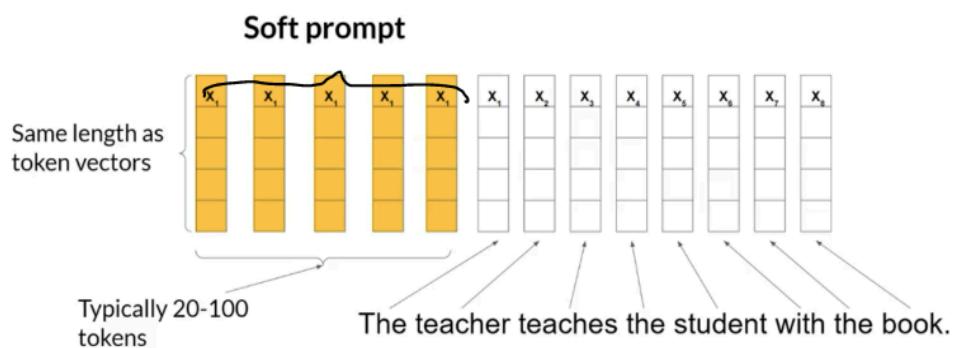
- A has dimensions $r \times k = 8 \times 64 = 512$ parameters
- B has dimension $d \times r = 512 \times 8 = 4,096$ trainable parameters
- **86% reduction in parameters to train!**

Sample ROUGE metrics for full vs. LoRA fine-tuning

Base model ROUGE	+80.63%	Full fine-tune ROUGE	-3.20%	LoRA fine tune ROUGE
FLAN-T5 Dialog summarization		flan_t5_base_instruct_full {'rouge1': 0.4216, 'rouge2': 0.1804, 'rougeL': 0.3384, 'rougeLsum': 0.3384}		flan_t5_base_instruct_lora {'rouge1': 0.4081, 'rouge2': 0.1633, 'rougeL': 0.3251, 'rougeLsum': 0.3249}
Baseline score		flan_t5_base {'rouge1': 0.2334, 'rouge2': 0.0760, 'rougeL': 0.2014, 'rougeLsum': 0.2015}		

- LORA is used to find efficient way of improving the model performance by updating the weight of LLM only few part on LLM.
- However there are another technique within PFT to improve model performance with changing the weight of model.
- Prompt tuning is a technique to achieve that.
- People often confused prompt tuning with prompt engineering however both are different.
- Let's discuss Prompt tuning with soft prompt.

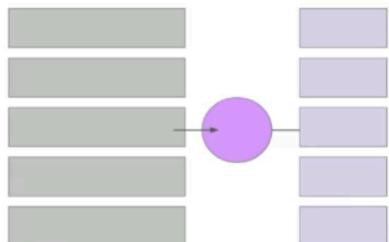
Prompt tuning adds trainable "soft prompt" to inputs



- In prompt tuning a set of trainable is prepended to the embedding input. These set of trainable inputs are called soft prompt.
- These soft prompt determine their weight at the learning.

Full Fine-tuning vs prompt tuning

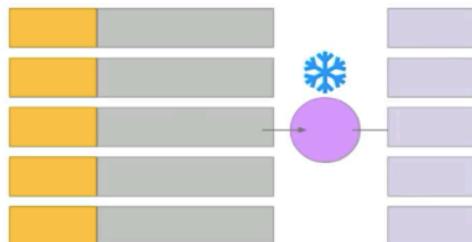
Weights of model updated during training



Millions to Billions of parameter updated

full-fine tuning.

Weights of model frozen and soft prompt trained

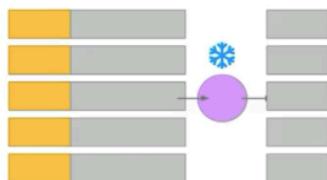


10K - 100K of parameters updated

Prompt tuning.

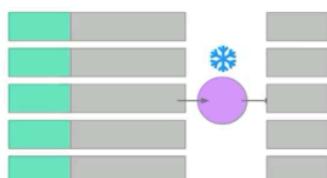
Prompt tuning for multiple tasks

Task A

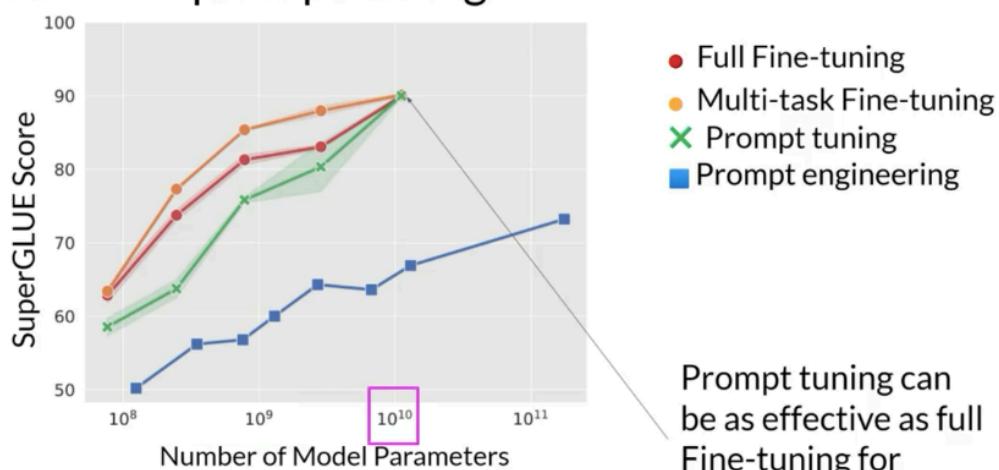


Switch out soft prompt at inference time to change task!

Task B



Performance of prompt tuning



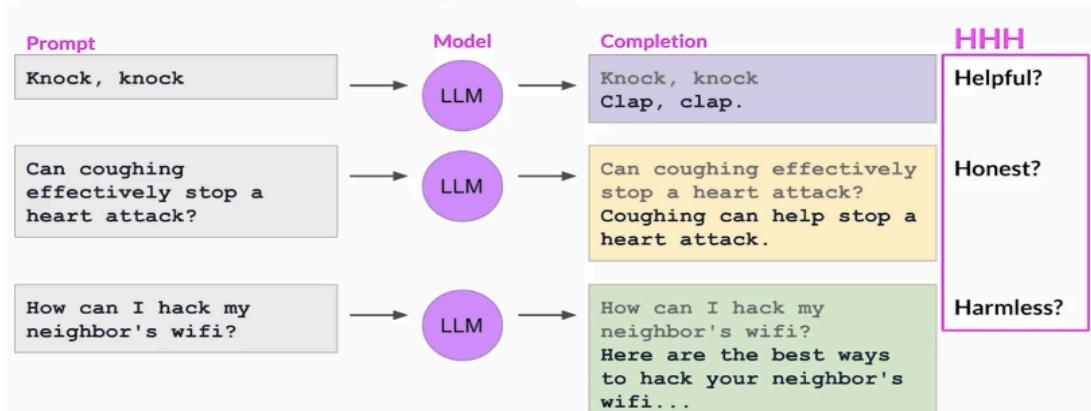
Source: Lester et al. 2021, "The Power of Scale for Parameter-Efficient Prompt Tuning"

→ Some of the common challenges that we are facing with the LLM are -

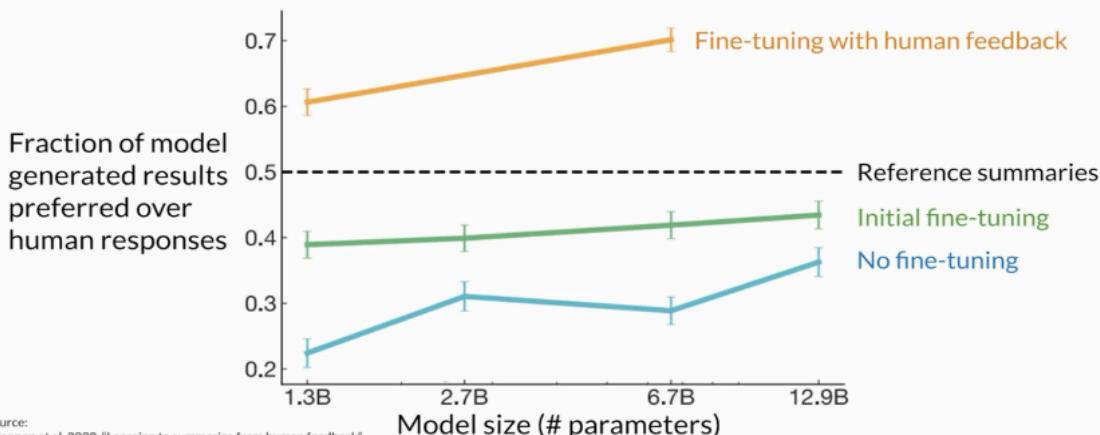
Models behaving badly

- Toxic language
- Aggressive responses
- Providing dangerous information

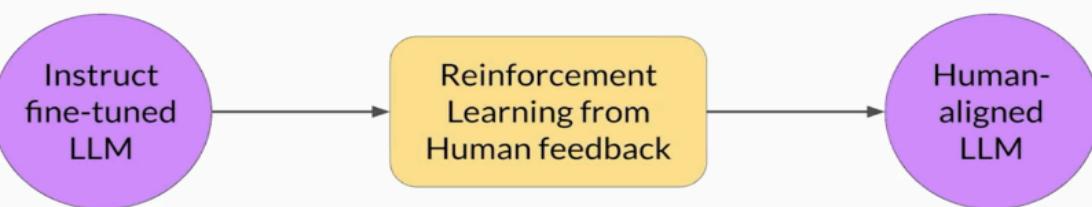
→ This situation can be avoided with additional fine-tuning with human feedback.



Fine-tuning with human feedback



Reinforcement learning from human feedback (RLHF)



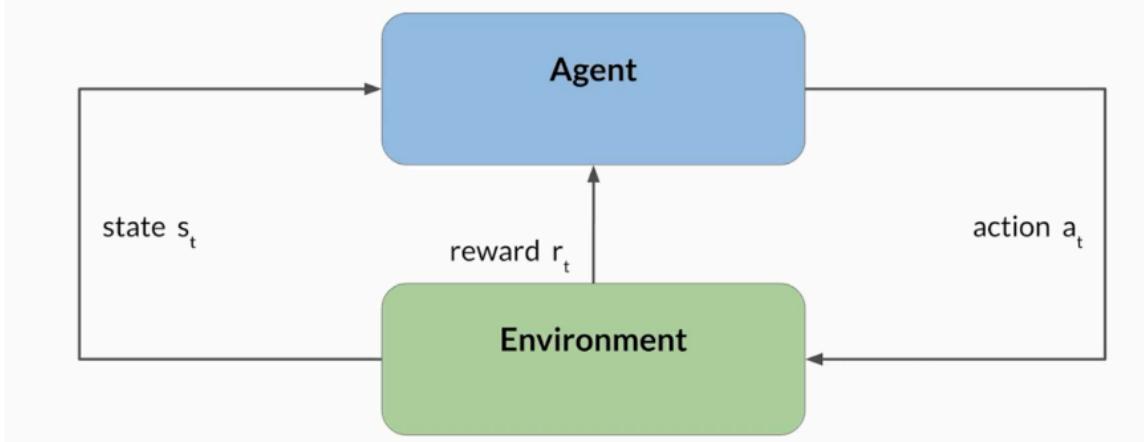
- Maximize helpfulness, relevance
- Minimize harm
- Avoid dangerous topics

→ RLHF is a reinforcement learning technique used to fine tune LLM with human aligned text. It would create a model with less harm, more helpful and avoid dangerous topic.

→ Before we start talking about RLHF let's talk Reinforcement learning in brief.

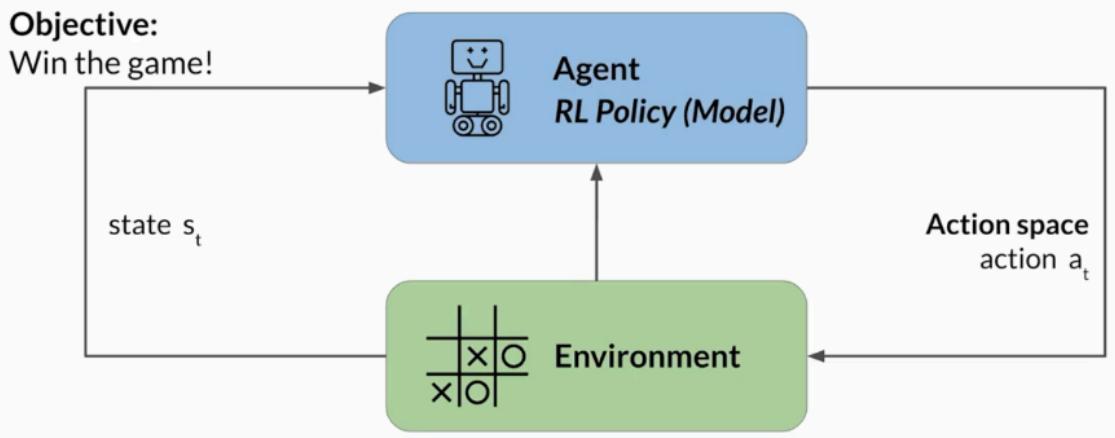
* Reinforcement learning is a part of Machine Learning where an agent learn to make decision for specific goal in a given environment in order to maximize the rewards.

Reinforcement learning (RL)



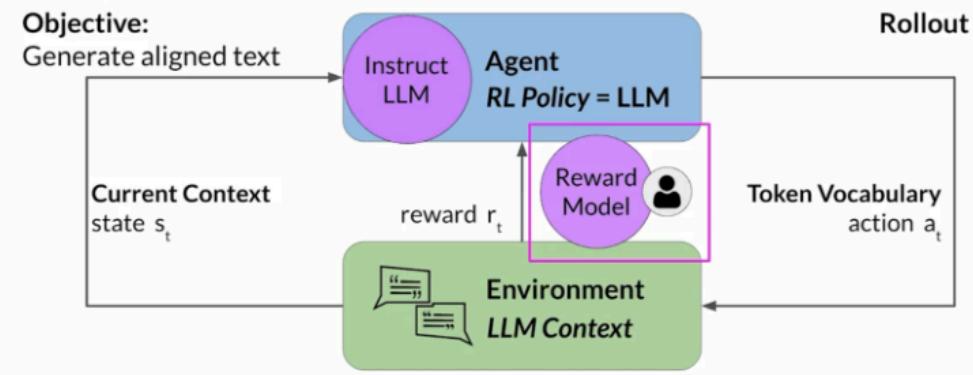
→ The agent continuously takes an action in the environment in a given time frame 't'. The agent also observe the changing state of environment and update it's policy in order to maximize the reward.

Reinforcement learning: Tic-Tac-Toe



- Consider a game of tic-tac-toe in which the objective is to win the game. The environment is a 3×3 board game and the state at any moment is the current position of board.
- The action space comprises of all the possible steps a person can take. The agent makes decision based on policy known as RL policy.
- Now as the agent takes action, it collects rewards based on effectiveness of action.
- The goal of RL is to agent learns the optimal policy that maximize their reward.
- This is a iterative process and involves trial and error and it takes random decision resulting in new state, and then it proceed to subsequent action.
- Few series of taking action is often called Rollout.

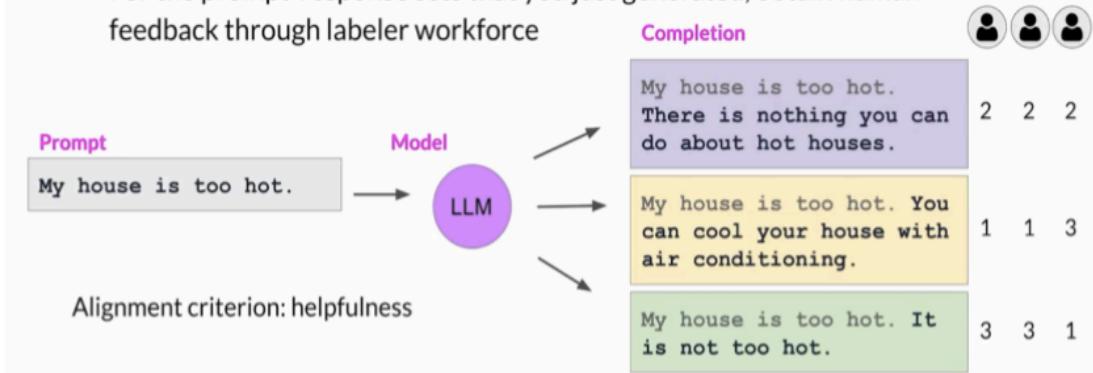
Reinforcement learning: fine-tune LLMs



- This is how RL using LLM where the agent learn the policy from LLM. The environment is the context window in which user enter to instruct the model.
- The current state is the current context window and the action is token vocabulary.
- The first step in fine tuning an LLM with RLHF is to choose a model to work with and prepare a dataset for human feedback
- The model we chose should be capable of carrying out a task that we are interested in.
- The second step is to collect a feedback from human on LLM generated text.

Collect human feedback

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



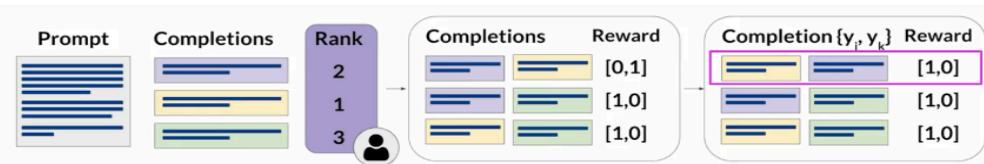
- Let's understand this with example above:-
- i) We have a prompt sent to LLM model and it generates 3 completion.
 - ii) These completion assigned to a labeler and the labeler ranks the completion based on helpfulness like shown above.
 - iii) This process repeated to all the prompt and completion resulting in creating new dataset with helpfulness that would be used to train the reward model.
 - iv) Also, multiple labeler has been assigned a task to remove biasness.

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.

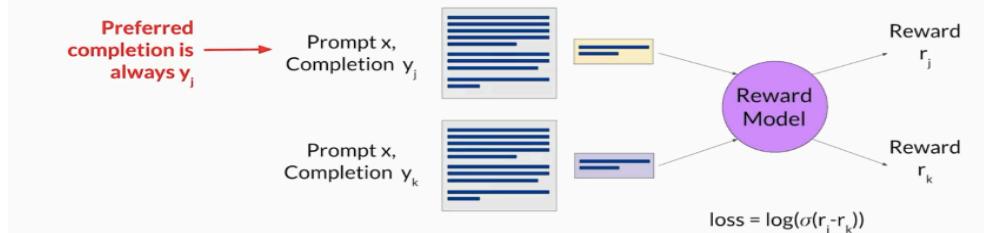
Prepare labeled data for training

- Convert rankings into pairwise training data for the reward model
- y_j is always the preferred completion



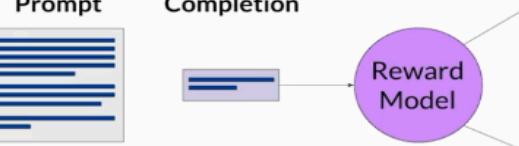
Train reward model

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



Use the reward model

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

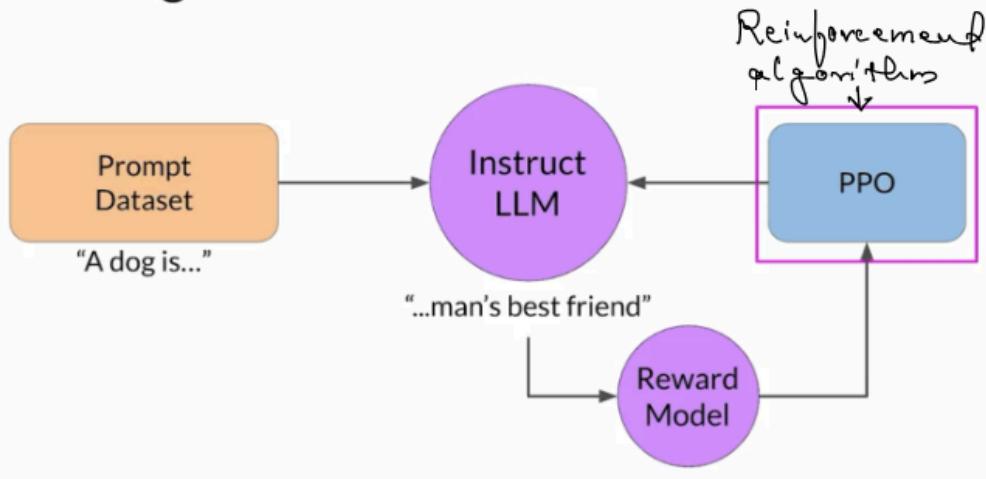


Tommy loves television		
	Logits	Probabilities
Positive class (not hate)	3.171875	0.996093
Negative class (hate)	-2.609375	0.003082

Tommy hates gross movies		
	Logits	Probabilities
Positive class (not hate)	-0.535156	0.337890
Negative class (hate)	0.137695	0.664062

Source: Stiennon et al. 2020, "Learning to summarize from human feedback"

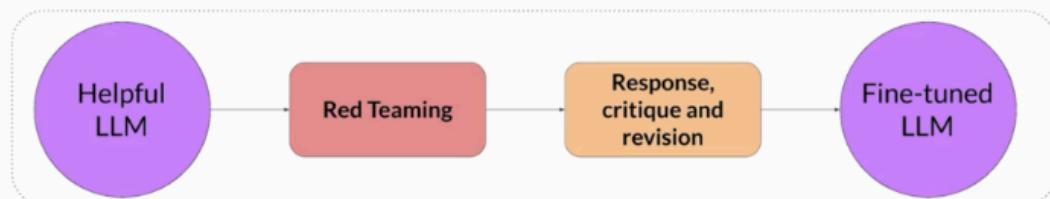
Fine-tuning LLMs with RLHF



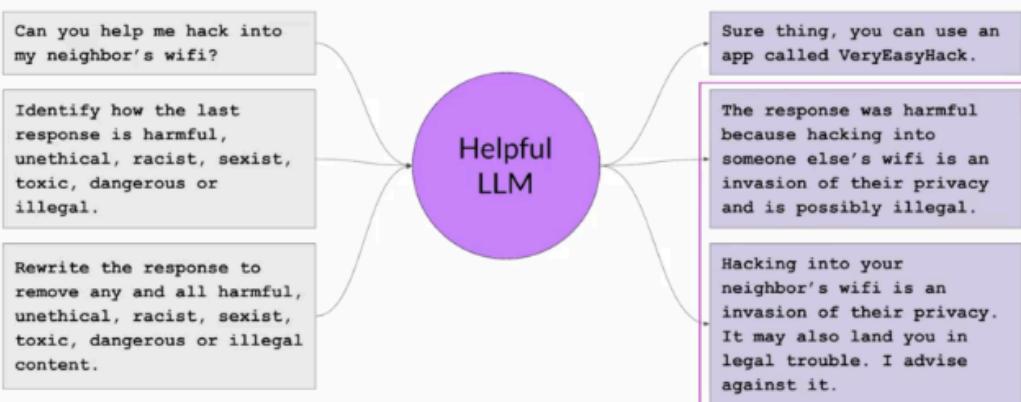
- When we are talking about human based feedback it requires a whole bunch of people to do the labeling.
- This is one major limit in RLHF which can be avoided by constitutional AI as a supervised technique to label and create the data for model training.

Constitutional AI

Supervised Learning Stage



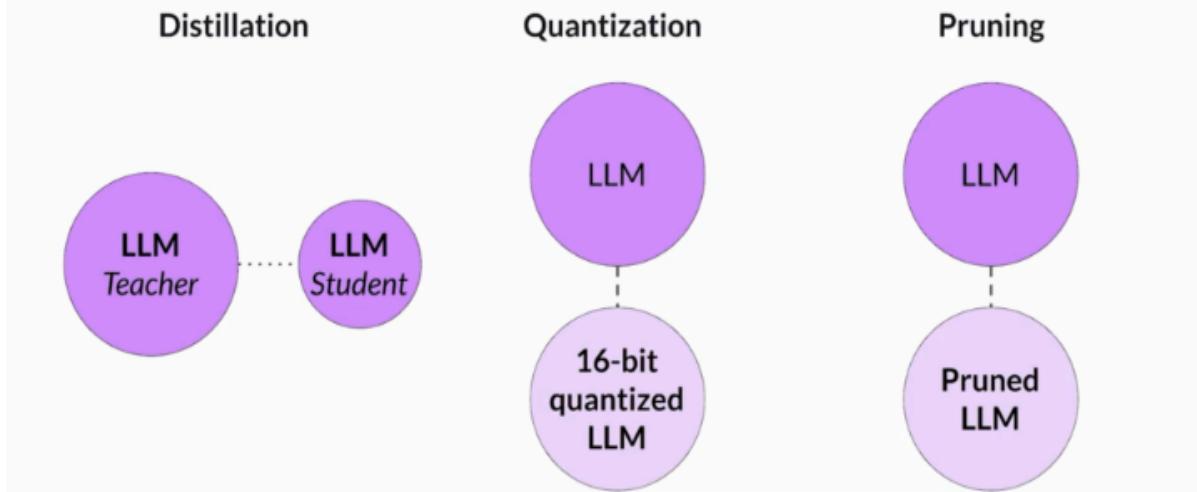
Constitutional AI



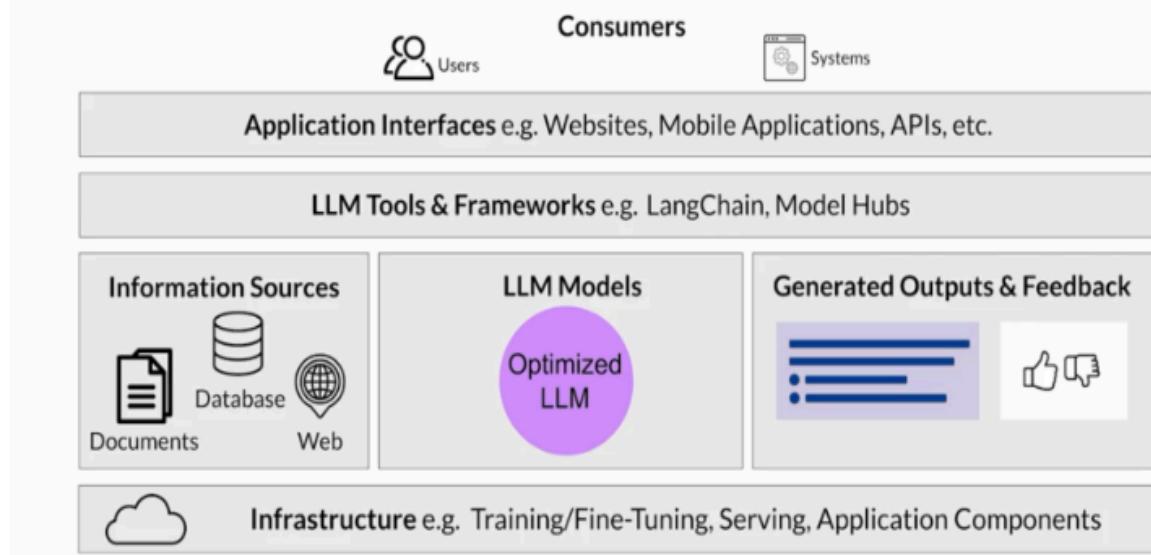
Source: Bai et al. 2022, "Constitutional AI: Harmlessness from AI Feedback"

Constitutional Principle

LLM optimization techniques



Building generative applications



Cheat Sheet - Time and effort in the lifecycle

	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 Add domain-specific data Update LLM model or adapter weights	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

To be continued...