

①

GenAI with LLM's

GenAI is subset of ML and the ML models that underpin genAI have learned these abilities by finding statistical patterns in massive datasets of contents that were originally generated by humans.

Foundation models (also called base models)

GPT, BERT, LLaMA, PaLM, BLOOM, FLAN-T5

- The text we pass to LLM is known as Prompt
- The space/memory available to prompt is called Context window
- Output of the model is called completion
- And the act of using the model to generate text is known as inference.

Transformer learns relevance & context of all words in sentence, not just its neighbour words but every word in a sentence. Then apply attention weights to those relationships so the model learns the relevance of each word to each other word no matter where they are input

#prompt engineering:- (In-context learning)

- One-shot inference: prompt with no example or label
Large LLM models are good at it.
- One shot inference: prompt with one example and label
- Few shot inference: prompt with multiple examples & label

#Generative Configuration:-

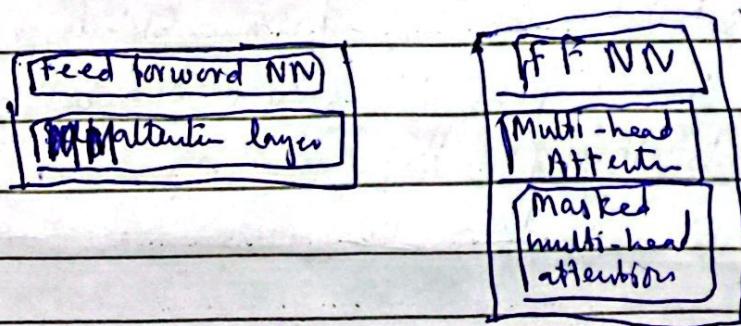
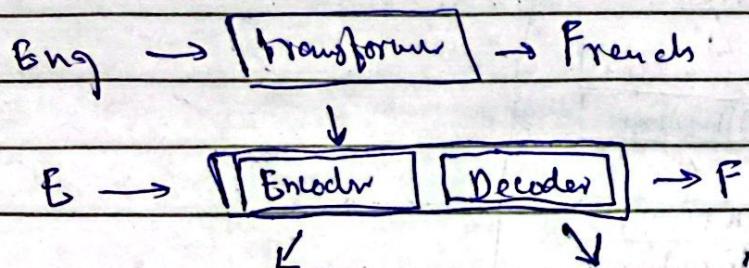
- Max new tokens: it is used to control the length of the generated text.
- Greedy Sampling: the model selects the word with the highest proba for the next token at each step.
- Random Sampling: selects word randomly
- Top-k Sampling: selects an output from the top-k results after applying random weighted strategy using proba
- Top-p Sampling: , , , ; using the random-weighted strategy with the top ranked consecutive results with proba & with a cumulative proba. $\leq p$
- Temp ↑ the ↑ the randomness & vice versa.

Attention: The ability of model to pay attention is important part of a sentence or image.

(3)

Transformers :-

- ① Input Embedding: each word is mapped to a vector. 512 length vector is used. But it is hyper parameter can be changed.
- ② Positional Encoding:

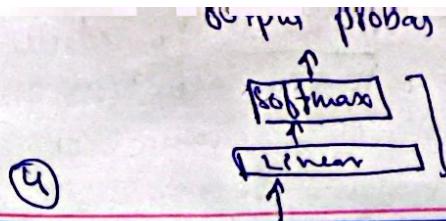


M-head Attention:-

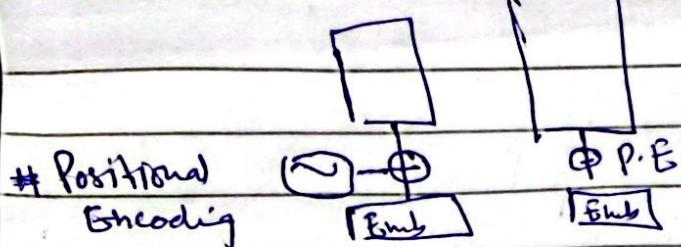
- All the words of sentences are feed together
- All the words of the s are compared to other words in sent.

FFNN:-

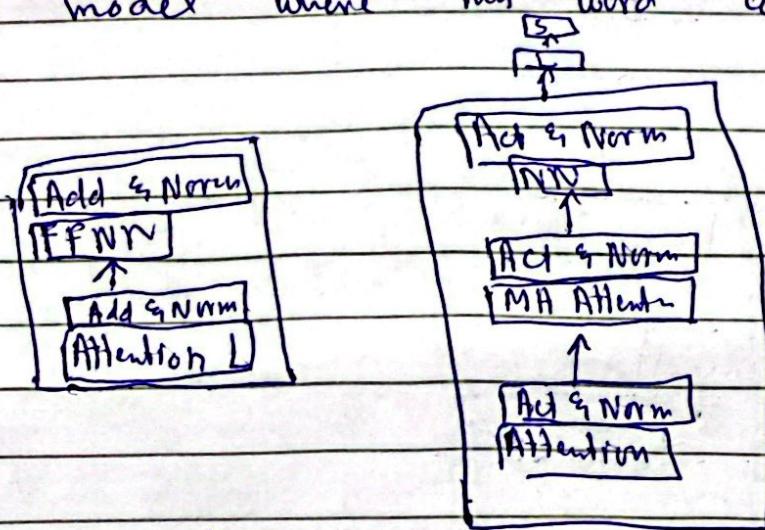
- Then they are passed through NN separately so they do not have info exchange
- NN are all different in each encoder



so that output of the decoders can be transformed into something we can understand & they basically into vector of length of amount of words we have



Annotations some info about each word which tells the model where this word comes in a sentence.

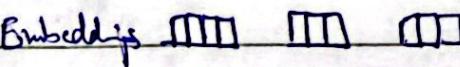


* Add \hookrightarrow Normalize layers :-

- To normalize the output coming from sub-layers.
Layer Normalization is used: an improved version of batch normalization
 - Add the info that went through the sub-layer & also just skip the sub-layer & then normalize them together
 - Multi-H Attention & Masked MHA are the same: but only difference is in a normal MHA L all the words are compared with all the other words but in MHA_L the words is compared with only the words coming before it

(5)

In beginning we have embeddings of words then
 In attention layer these ~~s~~ vector vectors are multiplied
 with some matrices. These are query, key, value matrices.
 These are values initialized randomly & are trained during
 training process to be learned, kind of like
 weights & biases. As a result we get query,
 key & value for each word, now we will keep these
 vectors to keep going with the calculations

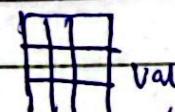
- we will calculate score of input that was fun
 each word against all the Embedds 

other words in the sentence.

- we take the dot product
 of query vector of each
 word against the key vector
 of all words

Query q_1, q_2, q_3 

Key k_1, k_2, k_3 

Value v_1, v_2, v_3 

Score $q_1 \cdot k_1, q_2 \cdot k_1, q_3 \cdot k_1$

Divide by 3 $12, 25, 5, 75$

Softmax $0, 93, 0, 07, 0, 06$

weighted $w_{v1} \square \quad w_{v2} \square \quad w_{v3} \square$

⑥

Scaling choices for Pre-training:

GOAL: maximize model performance

constraint
compute budget (GPUs, training time, cost)

- Two options to achieve

↑ Scaling choice
Dataset size
(no of tokens)



↑ Scaling choice
Model size
(no of parameters)

(7)

Types of Transformer :-

① Encoder Only: aka autoencoding models, pre-trained using Masked language model. Tokens of the input sequence are randomly masked and the training objective is to predict the mask tokens in order to reconstruct the original sentence.

- Good use cases:

- Sentiment analysis
- Named entity recognition
- Word classification
- Examples models:- BERT
- ROBERTA

② Decoder Only: aka Autoregressive models, pretrained on Causal language model, it predicts token based on previous sequence of tokens.

Good use cases: Text generation

Exp: -GPT . BLOOM

③ Sequence-to-sequence model: Encoder-Decoder. Masks random sequence of input tokens, then the tokens are replaced with a unique sentinel token. The Decoder reconstructs the mask token sequences.

Good use cases: Translation, text summarization, Q/A Answering

Example: T5, BART

Week #2 Fine-tuning LLM's with Instruction ③

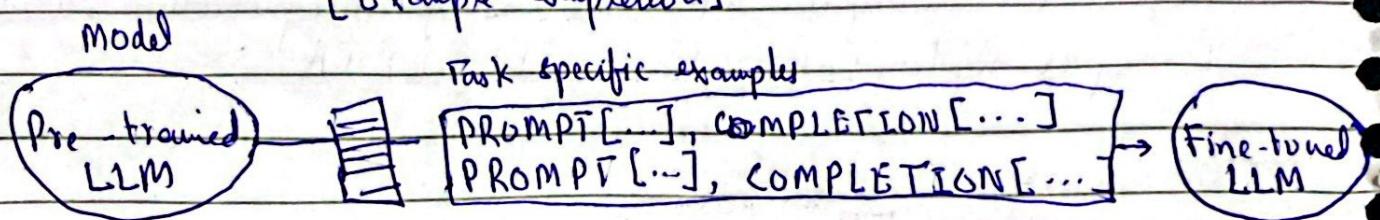
Fine-Tuning

- Using prompts to Fine-tune LLM's with instruction
 - if you want to fine tune your model to improve its summarization ability you'd build up a dataset of example within which the instructions be like:

summarize the following text:

[Example text]

[Example completion]



Catastrophic forgetting:- when a model is fully fine-tuned to a problem, then it forget its previous trained memory. Then it does not work on other problems.

How to avoid:

- Fine tune multiple tasks at the same time.
- Consider Parameter Efficient Fine-tuning (PEFT)

(9)

LLM Evaluation - Metrics:-

Rouge

- Used for text summarization
- Compares a summary to one or more reference summaries

BLEU Score

- Used for text translation
- Compares to human-generated translations.

Rouge - 1

• Referring (human)

{ It is cold outside

$$R-1 \cdot \text{Recall} = \frac{\text{unigram matches}}{\text{unigrams in reference}} \rightarrow \frac{4}{4} = 1.0$$

Generated output

It is very cold outside

$$R-1 \cdot \text{Precision} = \frac{\text{unigram matches}}{\text{unigrams in output}} \rightarrow \frac{4}{5} = 0.8$$

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

R-2: uses bigram matches instead of individual words

$$R-2-R = \frac{2}{3} = 0.67$$

$$R-2-P = \frac{2}{4} = 0.5$$

$$R-2 \cdot F1 = 2 \frac{0.335}{1.17} = 0.57$$

(2) Reference (Human)

10

ROUGE-L: uses LCS (Longest common subsequence).

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

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

$$\text{R-L F1} = \frac{2 \cdot \text{Precision} \times \text{Recall}}{\text{Precision} + \text{Recall}} = 2 \cdot \frac{0.2}{0.9} = 0.44$$

BLEU:

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

~~For fine tuning models used~~

Parameter efficient fine tuning (PEFT):

- Only updates small subset of parameters
- Most of the model weights are freezed in focus
- On fine tuning a subset of existing model parameters e.g. particular layers or components.

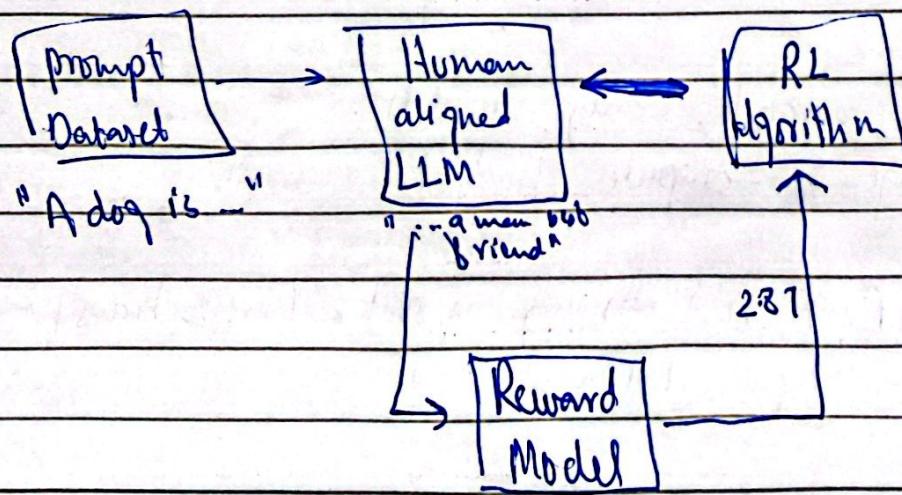
QLoRA = LoRA + quantization

(ii)

PEFT technique 2: Soft prompts

- You add additional trainable tokens to your prompt & unlike hard prompt which are direct & specific, soft prompt are more subtle & indirect.

Fine-tuning with reinforcement learning (RLHF)



- # So how fast do you need your model to generate completion?
- # What Compute budget you have available?
- # Are you willing to trade off model performance for improved inference speed or lower storage?
- # Do you intend for your model to interact with external data or other applications?
if so how will you connect to those resources?
- # How your model will be consumed

LLM optimization techniques:-

- ① Distillation: Larger teacher model train smaller student model. The student model statistically mimic the behaviour of the teacher model. either in final prediction layer or in hidden layers as well.
 - Effective for encoder only models.
 - To lower storage & compute budget.
- ② Quantization (PTQ): Post training quantization: transform weights to lower precision representation
- ③ Pruning: Remove model weights with values close or equal to zero
In theory reduces model size & improves performance. In practice only small % in LLMs are zero weights