

LLMs include BERT, GPT, LLaMA, Bloom, FLAN-T5, PaLM

previous generation can used RNN

Self-attention the ability to learn a tension in this way across the whole input significantly.

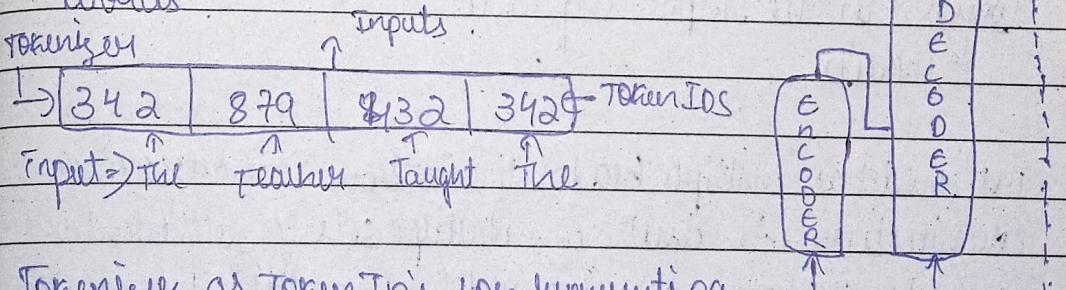
Attention is all you need by google and the university of Toronto. → gave birth to the Transformer architecture.

• Scale efficiently, parallel process • Attention to input

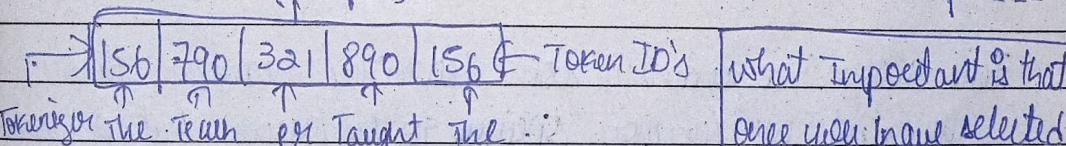
"Machine learning models are just big statistical calculators and they work with numbers not with words"

To work with words before passing text into the model to process, you must first tokenize the words. Simply put the words into numbers, with each number representing a position in a dictionary of all possible words that the model to process you must first tokenize the words.

Tokenize as Token ID's for complete words



Toneriges as Tonkin In's for supurouting parts of woods.



The model you must use the `Same Tokenizer` when you generate text.

Now the words or text converted to text, you can pass it to the embedding layer. This layer is a trainable vector embedding space, a high dimensional space where each token is represented as a vector and occupies a unique location within that space.

the words fed to the embedding layer can be plot into 3-d space and see the relationship b/w those words we can calculate the distance b/w them using angle b/w them. This gives the ability to mathematically understand language!

adding token vectors into base of encoder & decoder will add the positional encoding of the words.

→ positional encoding - "you preserve the information about the word order and don't lose the relevance of the position of the word in the sentence." → The model processes the input tokens in parallel.

summed the input tokens and positional encoding, and passing resulting vectors to the self attention layer.

this self attention layer will analyze the relationship b/w words in the input sequence.

↳ TOKENS

The model uses multiple head self attention in transformer architecture → results in multiple sets of self attention weights as heads are learned in parallel independently of each other.

This learns independently some level the meaning, others assuming we can not specify to the layer what to learn they are in general 12-100 layers in transformer.

the results of weights of self attention is passed to a fully connected feed forward network → The output of this layer is a vector of logits proportional to the probability score for each and every token in the tokenizer dictionary.

we can pass these logits to a final softmax layer where they are normalized to a probability score for each word. This output includes a probability for every single word in the vocabulary, so there's likely to be thousands of scores here one single token will have a score higher than other that token is most likely to be predicted.

There are no of ways to select from this vector probabilities

the data that leaves from the encoder to decoder is a deep representation of the structure and meaning of the input sequence this is injected at the middle of the decoder to influence the decoder's self attention mechanisms.

a start of sequence token is added to the input of the decoder. this triggers the decoder to predict the next token. it does based on the contextual understanding that is provided from the encoder.

the output goes to feed forward NN and through softmax layer and a ~~next token~~ token is generated. will continue this loop until all words. Then we will detokenize the words.

Transformer models

encoder only models
sequence to sequence (Seq2Seq)

Classification tasks.
sentiment analysis

BERT

encoder-decoder
Translation
different length
GPT
chatbots
LLM

only
decoder
GPT
chatbots
LLM

The text we give as input is known as prompt.
The act of generating text is known as inference.
The output text is known as the completion.

Context window = the limit of words that can be feed into the model.

Provide examples inside the context window is a way of asking
models to generate in a certain way. When it is called
In context learning.

Zero shot inference One shot inference few shot inference

Configuration parameters.

Generative configuration \rightarrow greedy \Rightarrow The word/token with
the highest probability is selected.

Random sampling \rightarrow Choose a random word/token randomly
by we reduce the likelihood that words will be repeated.

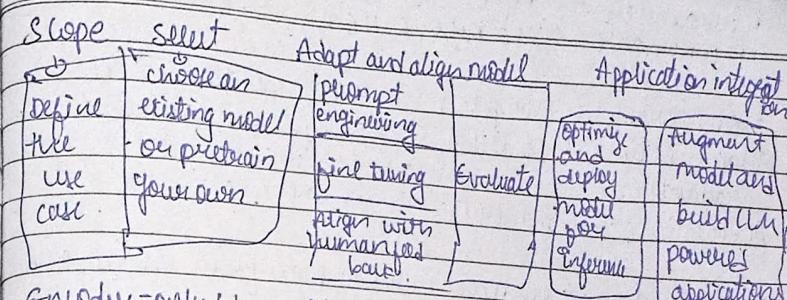
Top k - Sampling \rightarrow select an output from the top k results
applying random weighted strategy using the probabilities.

Suppose $k=10$ choose only from those top 10 tokens

Top p - Sampling = select an output using the random
weighted strategy with the top-ranked cumulative results by
probability and with a cumulative probability $c=p$.

Temperature \rightarrow ↑ randomness ↓ Temperature \rightarrow ↓ randomness
high \rightarrow random noise low \rightarrow deterministic

Generative AI project lifecycle



Encoder-only LM aka Autoencoding models
masked language modeling (MLM)

Objective: Reconstruct text ("decoding")
Good use cases: \rightarrow bidirectional
Sentiment analysis, named entity recognition,
word classification

models: BERT, ROBERTA

Decoder-only AKA Autoregressive models
causal language modeling (CLM)

Objective: predict next token
Good use cases: \rightarrow unidirectional

Text generation * other emergent behavior

models: GPT, BLOOM

Sequence to sequence models. Aka Encoder-decoder LM
span corruption (sentinel token)

Objective: reconstruct span

Good use cases:

* Translation, Text summarization, question answering
example models: T5, BART

"CUDA \rightarrow computer unified device architecture"

computational challenges of training LMs

1 parameter = 11 bytes (32-bit float)
16 parameters = 16×10^9 bytes (4 GB) \rightarrow 6 times

Model parameters (weights) \rightarrow 4 bytes per parameter.
Adam optimizer (2 states) \rightarrow +8 bytes per parameter.
Gradients \rightarrow +4 " " " "
Activations & temporary \rightarrow +8 " " " "

Quantization: reduce memory bytes from 32-bit to 16-bit or 8-bit integer

FP32 \rightarrow FP16 / BFloat16 / INT8

FP32 0 100 00000
sign 1bit
8 bits exponent
23 bits fraction.

	bits	exponent	fraction	memory needed for total one value
FP32	32	8	23	4 bytes
FP16	16	5	10	2 bytes
BFLOAT16	16	8	7	2 bytes
INT	8	-/-	2	1 byte

To improve performance we can increase the size of the dataset or by increasing the no. of parameters in the model.

constraint = compute budget (tokens, training time, cost)

1 petaflop/day = 1,000,000,000,000,000 (One quadrillion)
floating point operations per second.

Chinchilla paper \rightarrow

4.11.23

week - 2
Limitations of in context learning
In-context learning may not work for smaller models
Examples take up space in the context window.

pre-training where you train the LM using vast amounts of unstructured textual data via self-supervised learning

Fine tuning is a supervised learning process where you use a dataset of labeled examples to update the weights of the LM

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

Fine tuning is meant as Instruction fine tuning.

Sample prompt instruction template: The result contains a classification / sentiment analysis, now contains both an * Text generation
* Text summarization
instruction and the example from the dataset

you can pass the prompt from this above training dataset and this generate a completion now we can compare it with the label (training label) so the compared result distribution of the completion & label will give the cross entropy loss.

by this loss we can update our model weights in a standard backpropagation. we will do this for many batches of prompt completion pairs & several epochs, update the weights so the model's performance on the task improves.

The output of an LM is a probability distribution over tokens.

Catastrophic forgetting happens because the full fine tuning process modifies the weights of the original LM, which this leads to great performance on the single fine tuning task. it can degrade performance on other tasks.

how to avoid catastrophic forgetting.

Multitask fine-tuning is a extension of single task fine-tuning. training dataset is composed of examples inputs & outputs for multiple tasks.

Multitask instruction fine-tuning \rightarrow FLAN family
FLAN \rightarrow "fine tuned language net".

Model evaluation metrics.

$$\begin{aligned} \text{Accuracy} &= \frac{\text{Correct predictions}}{\text{Total predictions}} \\ &= \frac{\text{True Positives} + \text{True Negatives}}{\text{True Positives} + \text{False Positives} + \text{True Negatives} + \text{False Negatives}} \end{aligned}$$

Rouge and BLEU, two widely used evaluation metrics for different tasks.

Rouge for recall \rightarrow compares NLP summaries with human summaries.

BLEU \rightarrow Bilingual evaluation understudy is an algorithm designed to evaluate the quality of machine-translated text against human-generated translations.

In the anatomy of language, a unigram is a single word. A "bigram" is two words and an "n-gram" is a group of n-words.

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

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

$$\text{Rouge-1} = \frac{\text{precision} * \text{recall}}{\text{precision} + \text{recall}} = \frac{2.08}{1.8} = 0.89$$

F1 score is the harmonic mean of Recall & precision.
Rouge-2 \rightarrow bigram and soon for better we use the LCS: longest common subsequence.

BLEU score quantifies the quality of translation by checking how many n-grams in the machine generated translation match those in the reference translation.

BLEU & ROUGE are both cost effective to calculate we should not use them alone to report the final evaluation of a large language model.

Benchmarks = GLUE, SuperGLUE, Hellm.

GLUE = General language understanding evaluation
2018 - 2019 \rightarrow SuperGLUE

Massive multitask language understanding or MMU

MMU \rightarrow Holistic Evaluation of language models. multimetric approach measuring 7 metrics across 16 core scenarios.

- 1) Accuracy 2) Calibration 3) Robustness 4) Fairness 5) Bias
- 6) Toxicity 7) Efficiency.

Full fine tuning of large LMs \rightarrow is challenge so we have parameter efficient fine tuning (PEFT). only update a small subset of parameters. freeze most of the model weights & focus on fine tuning a subset of existing model parameters.

We can add new parameters or layers to the existing LMs. and finetune the new components!

full fine-tuning creates full copy of original LLM parameters
each of these is same size as the original model

PEFT methods

Sensitive-select subset of initial LLM parameters to fine-tune

reparameterization model weights using a low rank representation (LORA)

Additive: Add trainable layers or parameters to model (Adapters) (soft prompts prompt tuning)

CORA: low rank Adaption of LLMs

freeze most of the original LLM weights

inertia rank decomposition matrices

Decompose the weights of the smaller matrices

$$I \otimes A = B \times A$$

2 Add to original weights $\rightarrow I + B \times A$

In principle, you can use LORA to train for many steps, switch out the weights when you need to use them. And avoid having to store multiple full-size versions of the LLM.

Soft prompt just like prompt engineering.

The set of tunable tokens is called a soft prompt, they don't have any discrete fixed words.
we can think of them as virtual tokens that can take on any value within the continuous multidimensional embedding space.

we can train different set of soft prompts for each task and then easily sweep them out at inference time, just like LORA.

RLHF

→ Reinforcement learning from human feedback

model behaving badly

* Topic language * Aggressive responses, * providing dangerous information.

Helpful? Honest? Harmless?

Additional fine tuning with human feedback helps to better align models with human preferences and to increase this above 3 things

Instruct
fine-tuned
LLM

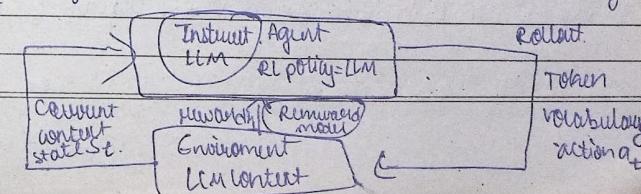
RLHF

Human aligned
LLM

RL - learning process is iterative & involves trial & error, the reward for the LLM are given assigned based on how closely the completions aligns with human preferences.

The feedback is represented in a scalar value zero or one, LLM weights are then updated iteratively to maximize the reward obtained from the human classifier.

The sequence of actions and states is called a rollout, instead of playout



collect human feedback.

- Define your model alignment criterion

→ for the prompt-response sets that you just generated, obtain human feedback through labeling workforce

The LLM will generate different completions based on the order of helpfulness from the most to least grant them using labelers.

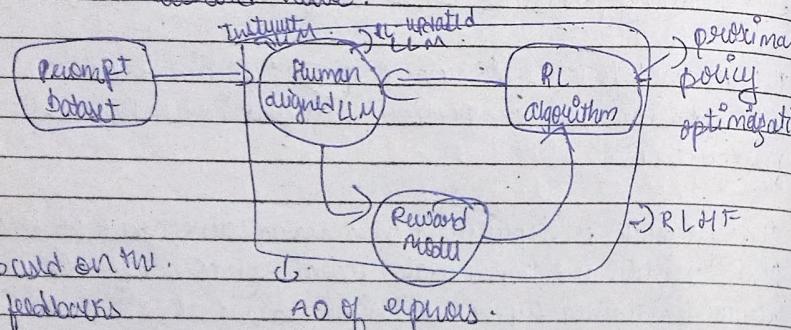
train the model with all the data that is labeled by human labelers, this data then can be classified without humans.

The feed back thumb up & thumbsdown is effective way to gather ranking.

Ranked feedback gives you more from completion data to train your reward model.

Reward model is a binary classifier to provide the positive & negative classes.

? Reward value.



build on the

feedbacks
reward

model gives reward values based on if RL algo will update or finetune weights and repeat till no longer

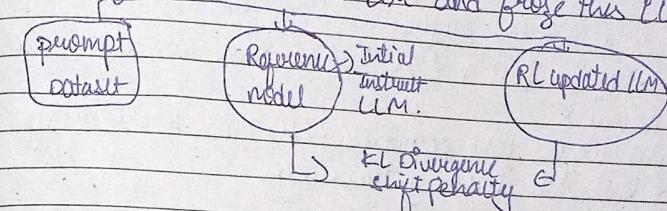
AO of epochs.

ppo → full model policy optimizations update the weights of the model (RL) which is very near to the previous version of LLM.

"keeping the changes within this small region results in a more stable learning".

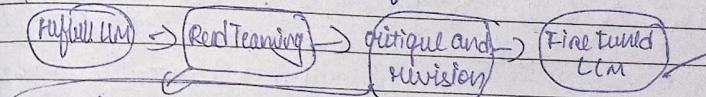
Rewarded hacking where agent learns to cheat the system by taking actions that maximize the reward received even if those actions don't align well with the objective.

To avoid the rewarded hacking we can use the reference model of a initial internet LLM and froze this LLM.



by comparing the completions of this two models we get the difference value that is called "fullback - Leibniz divergence" KL divergence

"constitutional AI is one approach of scale supervision!" (self supervision)



from the LLM to generate harmful content/completion
↳ Review or act as a critique to its own response.
by considering the constitutional principles.

Model optimizations for deployment

Distillation. → encoder only model

Train a smaller student model from a large teacher model

post training quantization (PTQ).

FP32 → FP16 / INT8, BF16 / 16

Pruning

removing model weights with values close or equal to zero.

→ full model fine-tuning; PEFT / LORA, Post training.

Quantization, Distillation and pruning all aim to reduce model size to improve model performance during inference without impacting accuracy.

Retrieval Augmented Generation.

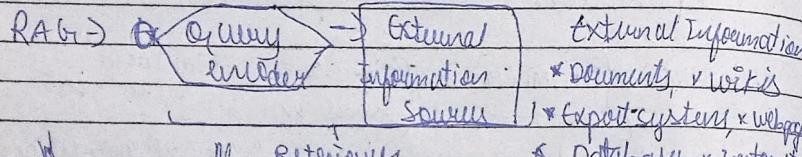
RAG → helps build LLM that can make use of external data sources.

↳ framework for providing LLMs access to data they did not see during training.

Chain of thoughts → human way of thinking.

PAL → program aided language models.

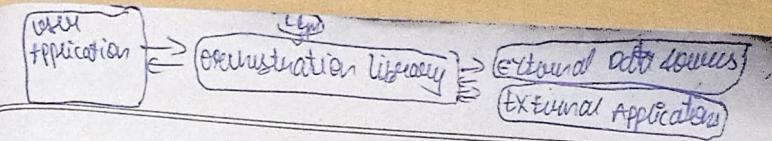
↳ framework that helps LLM to interact with other languages (Python interpreter)



↳ Lang chain → map the split of small chunks of external data to handle token limit.

* Embedding vectors.

Chain of thoughts humans take a step by step approach to solving complex problem.



ReAct: is a prompting strategy that combines chain of thought reasoning with action planning

* first, the task is defined telling the model to answer a question using the prompt structure you just explored in detail

* next, the instructions give more detail about what is meant by thought and the specifies that the action step can only be one of three types (for this example)

AWS SageMaker Jumpstart is a model hub it allows you to quickly deploy foundation models that are available within the service & integrate them into your own applications.