

Large Language Model

What is Large Language Model

→ A Language Model estimates the probability of a token sequence occurring ^{on} within a longer sequence of tokens

In language Model, the atomic unit that the model is training on. Here a token is
→ a word eg - "doge like cat"
These are the tokens
doge, like, cat
→ a character for example
Eg: Bike fish → consists of nine tokens
Space is also consists of one token

→ Subwords

→ In which single word can be single token or multiple tokens

eg → unwatched → The 3 subwords

un + watch + ed
Prefix + Root + Suffix

eg → Cat → Cat + s
Root + Suffix

Time: Token is atomic unit or smallest unit of language Model.

Here tokens are also been applied to computer vision & audio generation

Eg:-

When I hear rain on my roof, I in my kitchen
Probability

9.4%

Root soup

5.2%

warm up a kettle

3.6%

comes

2.5%

nap

2.2%

relax

In some situations it can do whole essay paragraph and sentence or can be called as sequence of tokens to complete the blank as shown in example model will pick up the token having greater percentage of probability or some threshold. Some filling the blanks can go to more higher complexity of task.

like ① Text generation

② Translating text from one language to another language

③ Summarizing documents.

N-gram Language Models

↳ N-grams are ordered sequence used to build the language models. Here N-gram refers to No of words in the sequence.

When N is 2 then model is called Bigram

If N is 5 then it called as 5-gram

Below are given in taking documents.

$4 \xrightarrow{2} b.$

\Rightarrow you are very nice.

$4 \xrightarrow{3} 2$

2-gram ① you are

② are very

③ very nice.

\Rightarrow 3-gram ① you are very

② are very nice.

Content

↳ Human can retain very long content, while machine Act 3 play you can collect or retain knowledge of act 1 etc the punch line of long conversation.

In language model Content is helpful information in front of token or after the target tokens

Eg: Content can help model to understand if orange is a fruit or color

Content can help language model to make better predictions so does 3-gram does that

so e.g. "orange is one two words" the 3-gram can get several many probabilities as annual.
→ even 3-gram model predicts lots of wrong answers. or even lot mistakes.

Longer N-grams can provide better content than short one. Different N-grams gives their relative outcome of each instance occurs. So when N becomes very large, then model has only a single instance of each outcome of N-tokens

Recurrent Neural Networks

↳ Recurrent neural networks provides more content than N-grams
So - Recurrent Neural Network is a neural network trained on sequence of tokens.

E.g. recurrent neural network can gradually learn selected content from each word in a sentence, kind of you are listening to someone speak. Large recurrent neural network can gain content from passage of several sentences.

* Though recurrent large neural network can learn more content than N-grams, but amount of useful content can intuit is still relatively limited

Recurrent neural network evaluate information token by token

⇒ Training recurrent neural network for long contents is constrained by the vanishing gradient problem.

LLMs → Large Language Models

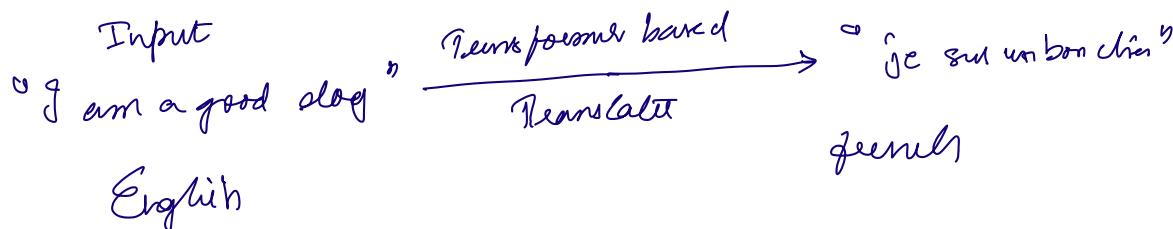
↳ It predicts a token or sequence of tokens, sometimes many paragraphs worth of predicted tokens.
LLMs make much better predictions than N-grams or

lement neural networks better. Because

- ① LMs contains far more parameters than sequence models
- ② LM's gather far more context

Transformers

This is state of art architecture for a wide variety of Language Model applications, such as translations

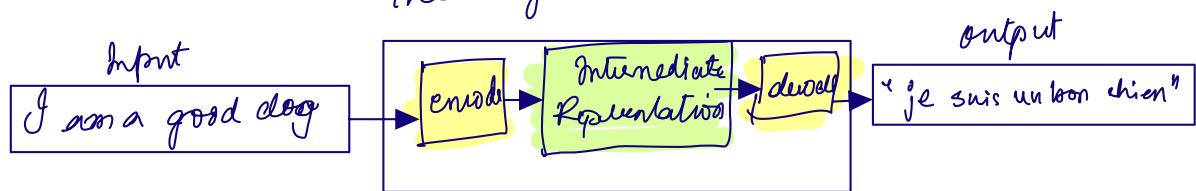


Transformers

- ① encoder → converts the text into an intermediate representation. It is enormous neural net
- ② decoder → converts that intermediate representation into useful text. A decoder is also enormous neural net.

Eg:

encoder converts the input to intermediate representation
the English



Partial Transformers

- ① ↳ Encoder only architectures
- ② ↳ Decoder only architectures

? This also exists

- ① Encoder Only architecture maps input text into an intermediate representation often an embedding layer

Voicebox

⇒ predicting ^{any} token in the input sequence (which is conventional role of language mod)

⇒ Creating a sophisticated embedding which could serve as input to the another system, such as classifier

* Decoder only architecture generates new tokens from the text already generated. This model typically generates sequences. Decoder only model can use their generation power to create continuous dialogue history and other prompts.

Self Attention

To enhance context

↳ Transfomers rely on Self attention. Effectively on behalf of each token of input, self attention asks the following questions

"How much does each other token of input affect the interpretation of this token?"

⇒ The self in self-attention refers to the input sequence. Self attention weights only the ^{importance} relationship between the sequence tokens in the input sequence.

Eg:

"The animal didn't cross the road because it was too tired."

→ Eleven words.

→ Each word is paying attention to other 10 words; wondering how much each those ten words matters to itself.

⇒ Eg. the pronoun "it" here pronouns always

refer to recent noun but in the example.

It refers 

* Self attention mechanisms determine the relevance of each nearby word to the pronoun. It.

None animal is more important than street.

* The animal didn't cross the road it was too tired

* glue self attentions would give more attention to street.

- * Some self attention mechanism are bidirectional.
i.e. they collect the content from the both sides and to generate representations of whole sequence.
- enables us bidirectional self attention
- denotes us unidirectional self attention

Multihed self-attention

↳ Each self attention layer consist of multiple self attention head.

Output layer is a mathematical operation like weighted average or dot product.

⇒ Self attention are assigned by random values
different heads can learn different relationships

between each word attended to nearby words.

* Self attention goes to learn relevance of every word in Content. so it is $O(N^2)$

$N \rightarrow$ number of tokens in the content

⇒ There consist of multiple heads per self attention layer

one multiple self attention layer.

so $O(N^2 \cdot s \cdot d) \rightarrow$ no of heads per layer.
↳ self attention layer.

LLM's

↳ variously known as.

① Generational LLM's

② base LLM's

③ pre-trained LLM's

- Foundationals are well trained on grammar, words and idiom. It can generate helpful sentences about topics it is trained on. Foundation LLM can create tasks like writing poetry.
 - ⇒ This is not solution of common ML problems such as regression or classification.
 - ⇒ It can be used as platform.
- * Fine Tuning is process where foundation LLM's can be prepared for solutions of applications. Secondary process called distillation with lesser parameters.

Fine Tuning

additional training called fine tuning is process of training Model

to predict on specified tasks

⇒ Fine Tuning trains the model with purpose of your application

does for -

⇒ Despite it is small but requires high computational power.

Fine tuning involves below tasks:

① updating weights & bias of every parameter on each back-propagation

② Smarter process called parameter-efficient tuning
Can fine tune an LLM by adjusting only a subset of parameters on each back propagation.

③ Fine tuned models are better than foundation LLM's prediction

Distillations

↳ It creates smaller versions of LLMs, the distilled LLM generates predictions much faster than and requires fewer computational & environmental power.
 models predictions are not good as foundation LLM.
 more parameters generate better predictions

Prompt Engineering

Prompt engineering enables an LLM's end user to customize the model's output. i.e end user clarify how the LLM's should respond to their prompt.

LLOis learns from example.

Showing one example to LLM is called

↑

- * One-shot prompting
- * Sometimes you have to give multiple examples called few-shot prompting.

Offline Inference

No of parameters in a LLM could be so large that online inference is too slow to be practical, real world tasks like logistics or classification.

→ Few offline inference is preferred. It called as bulk inference or static inference.

$$\begin{array}{r} 1 \ 2 \ 3 \\ + \ 5 \ 6 \\ \hline \end{array}$$