

01-09-25

Large Language Model (LLM)

Book: Hands-on LLM by O'REILLY

⇒ Language Model (Generalized Way): Models that recognize, summarize, translate, predict & generate text. A language model predicts the next word in a sentence by understanding the context & based on the probability of words following each other.

⇒ Why language Model is called Large Language Model: Because they are trained on large amount of data & have many parameters.

⇒ Large Language Model:

- It is a type of AI model based on deep learning architecture, typically the transformer architecture that is trained on massive corpora/corpus/database/dataset of natural language text.
- Its primary objective is to learn the probability distribution over sequence of words or tokens, enabling it to generate coherent & contextually relevant text.

⇒ Type of Language Model: There are two types:

- i) Representation Model: These are LLMs that don't generate text but are commonly used for task specific use cases like Classification, Clustering, Regression, etc (understands the data).
- ii) Generative Model: These are LLMs that generates text like GPT models. (generates the text).

⇒ Difference between Search Engine & ChatBot:

Search Engine

ChatBot

Purpose: It helps users to find info on the vast expanse of the internet.

Output: It returns list of relevant web-pages

• or links allowing users to browse & explore

Interaction: These provide a wide range of info.

Context: These have short-term memory.

Method: It uses algo. to crawl, index & rank web pages based on keywords.

Eg: Google, Yahoo

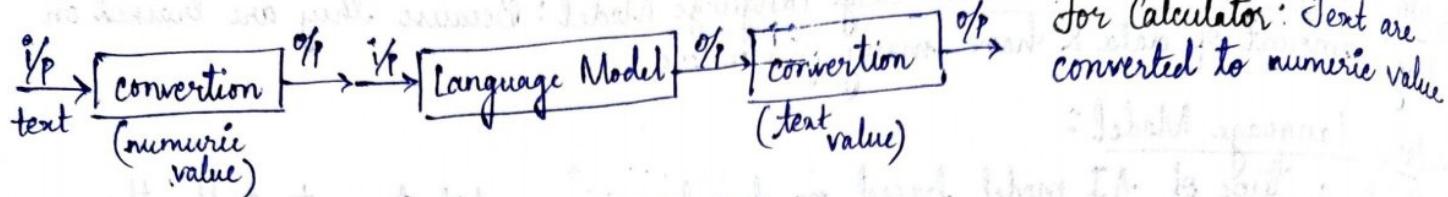
- It engages in conversation with user, answer, questions & often complete tasks.
- It provides a single, summarized answer & sometimes offering flow of questions.
- These engage in a conversation.
- These can maintain context across multiple interaction.
- It employs natural language processing & ML to understand user intent & provide responses.
- ChatGPT, Gemini, Copilot, DeepSeek

Chapter - 01

→ Language AI: It refers to sub-field of AI that focuses on developing technologies capable of understanding, processing & generating human languages.

05-09-25

→ Bag of Words (BoW) Representation: [Vocabulary of Words / Term Vocabulary / Dictionary]



S_1 : That is a cute dog. $\xrightarrow{①}$ Tokenization: S_1 : 'That', 'is', 'a', 'cute', 'dog'

S_2 : My cat is cute. $\xrightarrow{②}$ S_2 : 'My', 'cat', 'is', 'cute'.

Construct a Vocabulary

That is a cute dog My cat

S_1 : [1 1 1 1 1 0 0] → Vectors of Sentence 1

S_2 : [0 1 0 1 0 1 1] → Vectors of Sentence 2

- BoW model is a simple & widely used technique for text representation in Natural Language Processing (NLP).

- It transforms a sentence or document into a vector of word counts, ignoring grammar, word order and context. Each unique word in the corpus becomes a feature & resulting vector reflects the frequency of those words in the given text.

Limitations:

- BoW although an elegant approach but has flaws. It considers languages to be nothing more than an almost literal BoW & ignores the semantic nature / meaning of texts & also it is represented as Sparse Matrix/Vectors.

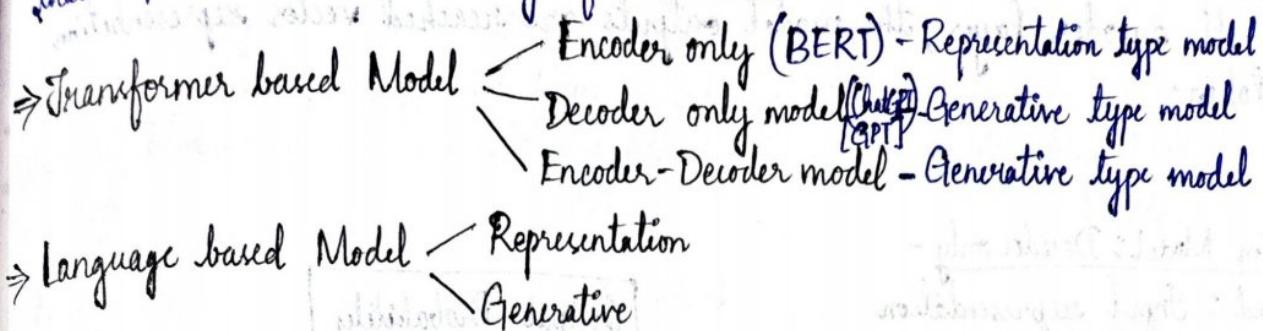
06-09-25

→ Representation of sentence as Dense Vector Embedding:

- In Natural Language Processing, a Dense vector representation of a sentence is of fixed size numerical representation that captures the semantic meaning of the entire sentence.
- Unlike, simple models like BoW, which only count word occurrences where Dense vectors encode relationship between words & their context where each word is converted into a vector using word embedding model. Word2Vec

→ Word Embedding: It represents words as numerical vector that captures the meaning & relationship between the words.

- Dense representation are foundational to have modern language models can understand & generate human language.
- A Dense Vector representation in a sentence is a compact, information rich numerical encoding that capture semantic meaning of a sentence.



⇒ Language based Model

Representation

Generative

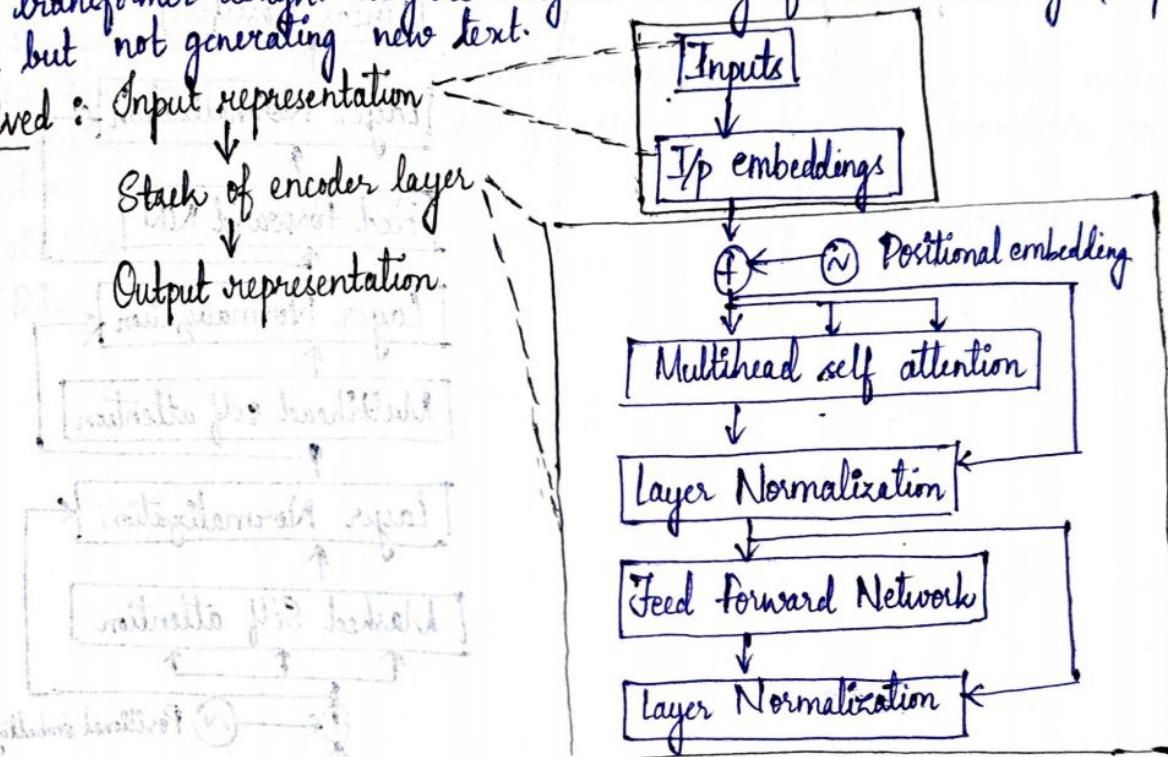
⇒ Representation Model : Encoder only -

Encoder only model are type of Transformer architecture that focus on the encoder part of the original transformer design. They are designed mainly for understanding & representing input text but not generating new text.

Steps involved : Input representation

↓
Stack of encoder layer

↓
Output representation



08-09-25

This above diagram is based on architecture of ANN : Input layer, hidden layer & output layer.

⇒ Input Representation :

- Here the input sentence is split into tokens.
- Each token is converted into a vector by combining token embedding or word embedding or positional embedding.

⇒ Stack of encoder layer :

The vectors pass through multiple identical encoder layers & each layer has -

- Multihead self attention: Here, each token attends to all tokens in the input to capture the relationship.
- Feed Forward Network: It is used to learn complex patterns & relationship in the input data.

iii) Layer Normalization: layer normalization & residual connection helps stabilizing the training & improve gradient flow.

→ Output Representation:

After passing all encoder layer, the model outputs are reached vector representation for each token.

11-09-25

→ Representation Model: Decoder only -

Steps involved: Input representation

↓
Stack of decoder layers

↓
Output representation

Output Probability

Softmax

Linear classifier

Layer Normalization

Feed forward NN

Layer Normalization

Multihead self attention

Layer Normalization

Masked Self attention

Positional embedding

Output Embedding

Output

Decoder only models are transformer architecture that use only decoder part of the original transformer. They are designed mainly for generating text (1 token at a time).

13-09-25

- Masked Self attention: It is used to ensure that the model can't attend the future tokens during training.
- Feed forward NN: It is used to transform the encoded representation back into the original input space (or) generate new data.

Linear classifier: It computes raw scores for each class.
Softmax: It converts the previous output into a probability distribution when a highest probability is chosen as predicted class.

Encoder only

- Models focus on input representation
- Working Model - I/p
- Encoder stack
- ↓
- o/p representation

Main purpose is representation learning, i.e., understanding I/p but not generating text.

Kind of representation model. Used for classification & feature extraction.

Eg: BERT, ROBERTa, DistilBERT

Decoder only

- It focus on output generation
- Working Model - I/p
- ↓
- Decoder stack
- ↓

Predict next step token

- Auto regressive text generation is their main purpose
- Kind of generative model
- Used for text generation

Eg: GPT, Llama

Encoder-Decoder

- Focus on both I/p representation & O/p generation
- Working Model - Encoder
- ↓
- Decoder

- Conditional generation (O/p text depends on I/p text) is their main purpose
- Kind of generative model
- Used for translation, summarization, Q/Aing, etc.
- Eg: T5, FLAN T5
- (Fine tuned Language Net T5)

Chapter-02

Tokens and Embeddings

⇒ Tokenization :

Tokenization is a fundamental step in natural language processing & language modelling. It involved breaking down texts into smaller units called tokens (words, sub-words, characters).

⇒ Why Tokenization is important :

- i) Raw text is not machine readable.
- ii) Bridge to embedding.
- iii) Vocabulary construction.
- iv) Handles linguistic variability.
- v) Supports positional encoding.
- vi) Context length.

15-09-25

⇒ Categ. of Tokenization (used in LMs) :

- 1) Word level Tokenization → 'large' 'language' 'Model'
- 2) Sub-word level Tokenization → 'la' 'nge' 'lang' 'age' 'Mo' 'del'
- 3) Characters level Tokenization → 'l' 'a' 'r' ... 'd' 'e' 'l'
- 4) Sentence level Tokenization → 'large language Model'

⇒ Types of Tokenization

1) Byte Pair Encoder (BPE) :

It splits words into common sub-words units, basically used in GPT models

2) Word piece Tokenization :

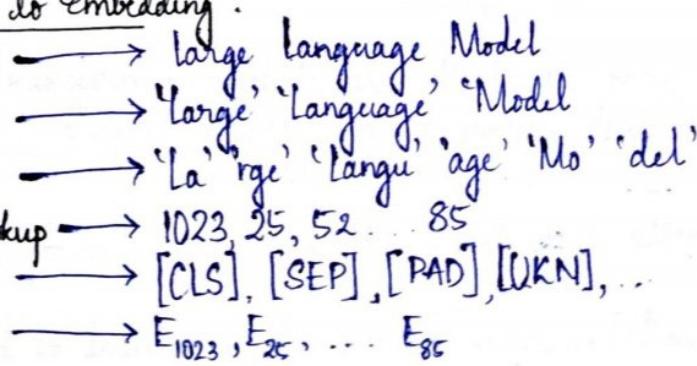
It is used by BERT, which is similar to BPE but with different rules for splitting.

3) Sentence piece Tokenization :

It is used by model like T5, Google, used by also Unigram Language Model.

⇒ Basic Pipeline of Tokenization to Embedding :

1. Text Input
2. Tokenization
3. Subword generation
4. Token mapping/vocabulary lookup
5. Special Token
6. Embedding layer



Subword Tokenization:

It is a technique that breaks words into smaller, frequent units called Subwords. This approach helps to handle out of vocabulary words & reduce vocabulary size while preserving semantic structure.

BPE Algorithm:

BPE is a data compression inspired algorithm used for Subword Tokenization. It begins with character level tokens & iteratively merges the most frequent adjacent pairs to form subwords.

(1) Pretokenization 1.1 Algorithm starts with a corpus of text data (Raw IP)

1.2 Corpus is pretokenized to create a list of all unique words. (List of unique word)

Basic Vocabulary

2.1: It is initialized with all unique characters found in list of unique words (Cont. initial var)

2.2: Each word in corpus is then represented as a sequence of symbols from this vocabulary. (Character token)

Pair Merging

3.1: Count the frequency of adjacent symbol pairs (Bigram counts) (Bigm)

3.2: Most frequent Bigram is then merged into a symbol & words in the corpus are updated to reflect this merge. (Merging)

This process continues iteratively until the desired vocabulary is reached.

Eg: 1.1: hello
1.2: 'hello'

2.1: 'h' 'e' 'l' 'l' 'o'

2.2: 'h' 'e' 'l' 'l' 'o'

3.1: ('h', 'e') ('e', 'l') ('l', 'l') ('l', 'o')

3.2: ('l') → 'l'

→ 'h' e l o → 'h' e 'l' o

Bigm	Freq
h	1
e	1
l	2
o	1

Eg: aa abc abc

Iteration 1 :-
1.1: aa abc abc
1.2: 'aa' 'bc' 'bc'

2.1: 'a' 'a' 'b' 'c'

2.2: 'a' 'a' 'a' 'b' 'c' 'a' 'b' 'c'

3.1: ('a', 'a') ('a', 'b') ('b', 'c') ('a', 'b') ('b', 'c')

3.2: ('a', 'b') → 'ab'

Bigm	Freq
aa	1
ab	2
bc	2

2.1: 'a' 'b' 'c' 'ab'

2.2: 'a' 'a' 'ab' 'c' 'ab' 'c'

3.1: ('a', 'a') ('ab', 'c') ('ab', 'c')

3.2: ('ab', 'c') → 'abc'

Bigm	Freq
aa	1
abc	2

2.1: 'a' 'b' 'c' 'ab' 'abc'

2.2: 'a' 'a' 'abc' 'abc'

3.1: ('a', 'a')

3.2: ('a', 'a') → 'aa'

Bigm	Freq
aa	1

Iteration 3 :-

18-09-25

Q) Consider the following corpus with pre-tokenized text with its freq. Using BPE algo. generate the sub-words & display final vocab

Ans: For eg: cat bat bag tag cats

1.1: cat bat bag tag cats

1.2: 'cat' 'bat' 'bag' 'tag' 'cats'

3.1: ('c', 'a') ('a', 't') ('a', 'b') ('b', 'a') ('a', 't') ('t', 'a') ('t', 'b') ('b', 'a') ('a', 'g') ('g', 't') ('t', 'a') ('a', 'g') ('g', 'c')

('c', 'a') ('a', 't') ('t', 's')

Bigm	Freq
ca	2
at	3
tb	2
ba	2
ag	2
gt	1
ta	1

3.2: ('a', 't') → 'at'

- It-2:
- 2.1: 'a' 'b' 'c' 'd' 'g' 's' 't' 'at'
 - 2.2: 'c' 'at' 'b' 'at' 'b' 'a' 'g' 't' 'a' 'g' 'c' 'at' 's'
 - 3.1: ('e' 'at') ~~f~~ 'at' ('b' 'at') ('b' 'a') ('a' 'g') ('t' 'a') ('a' 'g') ('a' 'at') ('at' 's')

BiGm	Freq
c,at	2
b,at	1
b,a	1
a,g	2
t,a	1
ats	1

3.2: 'c',at → 'cat'

- It-3:
- 2.1: 'a' 'b' 'c' 'd' 'g' 's' 'f' 'at' 'cat'
 - 2.2: 'cat' 'b' 'at' 'b' 'a' 'g' 't' 'a' 'g' 'cat' 's'
 - 3.1: ~~cat~~, ('b' 'at') ('b' 'a') ('a' 'g') ('t' 'a') ('a' 'g') ('cat' 's')

BiGm	Freq
b,at	1
b,a	1
a,g	2
t,a	1
eats	1

3.2: 'a',g → ag

BiGm	Freq
ea	15
at	20
#ba	17
ag	16
ts	1

It-2

BiGm	Freq
cat	15
b,at	5
ba	12
ag	16
ta	4
ats	5

3.2: 'b','g' → 'ag'

BiGm	Freq
c,at	15
b,at	5
b,ag	12
t,ag	4
ats	5

It-4

BiGm	Freq
bat	5
b,ag	12
t,ag	4
cats	5

It-5

BiGm	Freq
t,ag	4
cats	5

BiGm	Freq
b,at	5
t,ag	4
cats	5

3.2: 'bat' → bat

BiGm	Freq
t,ag	1

3.2: 't','ag' → tag

M-09-25
Word piece alg.:
 It is a sub-word tokenization algorithm developed by Google, widely used in models like BERT. Its primary goal is to handle rare or unseen words more effectively than traditional word level tokenization. Instead of treating each word as a token, wordpiece breaks words into smaller, meaningful subword units such as prefix, suffix or roots based on statistical patterns in the training corpus.

(1) Pre tokenization: 1.1: Raw text

1.2: Generate the tokens

(2) Base vocabulary: 2.1: Construction of initial vocabulary

2.2: Generate Character tokens Eg: token → t, ##0, ##k, ##e, ##n

(3) Pair merging: 3.1: Generate Bigram

3.2: Merge
 $x, \#\#y \rightarrow xy$
 $\#\#x, \#\#y \rightarrow \#\#xy$

score = freq. of pair
 freq of 1st token × freq of 2nd token

Eg: hello

It-1 1.1: hello
 1.2: 'hello'

2.1: [h, ##e, ##l, ##0]

2.2: {h, ##e, ##l, ##l, ##0}

3.1: (h, ##e), (##e, ##l), (##l, ##l), (##l, ##0)

pair	score
h, ##e	$y_1 \times 2 = y_2 = 0.5$ h, ##e → he
##e, ##l	$y_2 \times 4 = y_8 = 0.12$
##l, ##l	$y_4 \times 4 = y_{16} = 0.06$
##l, ##0	$y_4 \times 1 = y_9 = 0.25$

It-2: 2.1: [h, ##e, ##l, ##0, he]
 2.2: [he, ##l, ##l, ##0]

pair	score
h, ##l	$y_1 \times 4 = y_4 = 0.25$ he, ##l → hel
##l, ##l	$y_4 \times 4 = y_{16} = 0.06$
##l, ##0	$y_4 \times 1 = y_9 = 0.25$

It-3: 2.1: [h, ##e, ##l, ##0, he, hel] 3.1: (hel, ##l), (##l, ##0)

2.2: [hel, ##l, ##0]

pair	score
hel, ##l	$y_1 \times 2 = y_2$ hel, ##l → hell
##l, ##0	$y_2 \times 1 = y_2$

It-4: 2.1: [h, ##e, ##l, ##0, he, hel, hell]. 3.1: (hell, ##0)

2.2: [hell, ##0]

pair	score
hell, ##0	$y_1 \times 1 = 1$ hell, ##0 → hello

Eg: aa abc abc

It-1 1.1: aa abc abc

1.2: aa abc abc

2.1: a, ##b, ##c

2.2: a, ##a, d, ##b, ##c, a, ##b, ##c

3.1: (##a), (a, ##b), (##b, ##c), (a, ##b), (##b, ##c)

3.2: pair score

a, ##a	$y_3 \times 1 = y_3 = 0.35$
a, ##b	$y_3 \times 4 = 2/12 = y_6 = 0.16$
##b, ##c	$2/4 \times 2 = y_4 = 0.25$

<u>It-2:</u>	2-1: [a, #b, #c, aa]	3.1: (a, #b), (#b, #c), (a, #b), (#b, #c)						
	2-2: {aa, a, #b, #c, a, #b, #c}	3.2: pair score <table border="1"> <tr> <td>a, #b</td> <td>$\frac{1}{2} \times \frac{1}{4} = \frac{1}{4} = 0.25$</td> <td>a, #b \rightarrow ab</td> </tr> <tr> <td>#b, #c</td> <td>$\frac{1}{4} \times \frac{1}{2} = \frac{1}{4} = 0.25$</td> <td></td> </tr> </table>	a, #b	$\frac{1}{2} \times \frac{1}{4} = \frac{1}{4} = 0.25$	a, #b \rightarrow ab	#b, #c	$\frac{1}{4} \times \frac{1}{2} = \frac{1}{4} = 0.25$	
a, #b	$\frac{1}{2} \times \frac{1}{4} = \frac{1}{4} = 0.25$	a, #b \rightarrow ab						
#b, #c	$\frac{1}{4} \times \frac{1}{2} = \frac{1}{4} = 0.25$							
<u>It-3:</u>	2-1: [a, #b, #c, aa, ab].	3.1: (ab, #c), (ab, #c)						
	2-2: {aa, ab, #b, #c, ab, #b, #c}	3.2: pair score <table border="1"> <tr> <td>ab, #c</td> <td>$\frac{1}{2} \times \frac{1}{2} = \frac{1}{2} = 0.5$</td> <td>ab, #c \rightarrow abc</td> </tr> <tr> <td>ab, #c</td> <td>$\frac{1}{2} \times \frac{1}{2} = \frac{1}{2} = 0.5$</td> <td></td> </tr> </table>	ab, #c	$\frac{1}{2} \times \frac{1}{2} = \frac{1}{2} = 0.5$	ab, #c \rightarrow abc	ab, #c	$\frac{1}{2} \times \frac{1}{2} = \frac{1}{2} = 0.5$	
ab, #c	$\frac{1}{2} \times \frac{1}{2} = \frac{1}{2} = 0.5$	ab, #c \rightarrow abc						
ab, #c	$\frac{1}{2} \times \frac{1}{2} = \frac{1}{2} = 0.5$							

20-09-25

3) Sentence Piece Algorithm:

Note: In Word Piece it isn't possible to find a sub-word in the vocabulary, the whole sub-word is tokenized as [UNK] (Unknown Token). Eg: Funny \rightarrow F, #u, #n, #n, #y
Funny \leftarrow [UNK]

whereas in BPE algo, only missing token is replaced by UNK.

3) Sentence Piece Algorithm:

It is a language independent sub-word tokenization algorithm developed by Google, designed to train directly on raw text without requiring pre-tokenization.

Unlike Word Piece or BPE, which assume space separated words as input, Sentence Piece treats the input as a raw stream of Unicode characters & learns subword units using either BPE or Unigram Language Model. This approach enables consistent handling of white-space punctuations & multi-lingual corpora (appr which works on multiple languages).

(1) Initializing the Base vocabulary: It includes all the characters & the most frequent substrings found in the corpus.

(2) Log likelihood loss computation: The algo. calculates log likelihood loss over the entire training data based on the current vocabulary. It aims to find the smallest & the most effective vocabulary that still adequately represents the training data.

$$\text{Loss} = \sum \text{Freq} * (-\log(P(\text{word})))$$

(3) Finding candidates for removal: Symbols that contribute less to the overall rep of the data are marked as candidates for removal.

(4) Progressive Vocabulary Pruning: The algo. removes a small percentage of symbols that have the least impact on the log likelihood loss. These are the symbols for which the increase in training loss is the lowest if they are removed. This process is repeated iteratively. After each round of pruning, the log likelihood loss is recalculated with the updated vocabulary & process continues until the vocabulary reaches the desired size.

Q) Consider the following corpus

word	Freq	Using Sentence Piece algo, construct the vocabulary & display it
run	3	
bug	5	
fun	13	
sun	10	

Step 1: Possible word segments - for each word are run → r, u, n, ru, un

bug → b, ug, bu, ug

fun → f, un, fu, un

sun → s, un, su, un

Initial vocabulary: [r, u, n, ru, un, b, g, bu, ug, f, fu, s, su]

Compute Freq. of each tokens:

Token	r	u	n	ru	un	b	g	bu	ug	f	fu	s	su
Freq.	3	31	26	3	26	5	5	5	5	13	13	10	10

22-09-25

Q) Let us consider the following corpus

Word	Freq	[continuation of prev qns]
run	3	
bug	5	
fun	13	
sun	10	

Step 1: Initialize the Bas Vocabulary

1.1: Generate the tokens/substrings

run → r, u, n, ru, un

bug → b, ug, bu, ug

fun → f, un, fu, un

sun → s, un, su, un

1.2: Initialize the vocabulary, InVar → {r, u, n, b, g, ru, un, bu, ug, f, fu, s, su}

1.3: Token frequency

Token	r	u	n	ru	un	b	g	bu	ug	f	fu	s	su
Freq.	3	31	26	3	26	5	5	5	5	13	13	10	10

1.4: Compute the probability of each token

$$P(t) = \frac{\text{freq}(t)}{\text{Total freq.}}$$

Token	r	u	n	ru	un	b	g	bu	ug	f	fu	s	su
Prob	0.0193	0.2	0.1677	0.0193	0.1677	0.0323	0.0322	0.0322	0.0322	0.0338	0.0338	0.013	0.013

1.5: Find most probable segmentation.

Suppose, run, possible segmentation are (r, u, n) (ru, n) (r, un)

$$p(r, u, n) = p(r) \times p(u) \times p(n) = 0.0193 \times 0.2 \times 0.1677 = 0.00064$$

$$p(ru, n) = p(ru) \times p(n) = 0.0193 \times 0.1677 = 0.00323 \quad [\text{max}]$$

$$p(r, un) = p(r) \times p(un) = 0.0193 \times 0.1677 = 0.00323$$

$$p(b, u, g) = p(b) \times p(u) \times p(g) = 0.0322 \times 0.2 \times 0.0322 = 0.00020$$

$$p(bu, g) = p(bu) \times p(g) = 0.0322 \times 0.0322 = 0.00103 \quad [\text{max}]$$

$$p(b, ug) = p(b) \times p(ug) = 0.0322 \times 0.0322 = 0.0013$$

25/09/25

$$P(F, u, n) = P(F) \times P(u) \times P(n) = 0.0838 \times 0.2 \times 0.1677 = 0.00281$$

$$P(fu, n) = P(fu) \times P(n) = 0.0838 \times 0.1677 = 0.01405 \text{ [max]}$$

$$P(F, un) = P(F) \times P(un) = 0.0838 \times 0.1677 = 0.01405$$

$$P(s, u, n) = P(s) \times P(u) \times P(n) = 0.0645 \times 0.2 \times 0.1677 = 0.00216$$

$$P(su, n) = P(su) \times P(n) = 0.0645 \times 0.1677 = 0.01081 \text{ [max]}$$

$$P(s, un) = P(s) \times P(un) = 0.0645 \times 0.1677 = 0.01081$$

word	Freq.	Split/most probable	Score
sun	3	(u, n)	0.00325388
bug	5	(bu, g)	0.00103
fun	13	(fu, n)	0.01405
sun	10	(su, n)	0.01081

Step 2: Calculate Log Likelihood loss.

$$\text{loss} = \sum \text{Freq.} * (-\log(p(\text{csplit})))$$

$$= 3 * (-\log(0.00325388)) + 5 * (-\log(0.00103)) + 13 * (-\log(0.01405)) + 10 * (-\log(0.01081))$$

$$= 66.102$$

Step 3: Removal of candidate from vocabulary.

To determine which token to remove, we need to calculate the associated loss for each token in the vocabulary. However, that is not an elementary token, then compare the loss.

Suppose, you want to eliminate 'un' -

St-2: $\text{Var} = \{r, u, n, ru, b, bu, ug, f, fu, s, su\}$

After eliminating 'un', again loss value will be 66.102.

27/09/25 Suppose, you want to eliminate 'ru' -

$\text{Var} = \{r, u, n, un, b, bu, ug, f, fu, s, su\}$

After eliminating 'ru', again loss value will be 66.102.

Substring | Ass. loss As, all associated loss are same, pick anyone for removal.

Substring	Ass. loss
ru	66.102
un	66.102
bu	66.102
ug	66.102
fu	66.102
su	66.102

Let us assume, we removed 'un'.

St-2: $\text{Var} = \{r, u, n, ru, b, bu, ug, f, fu, s, su\}$

1.3:

Token	-
Freq	-

 Everything except 'un'

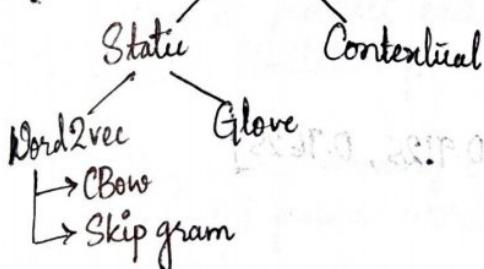
Total = 29

1.4: Similarly without 'un'

Differentiate between BPE, WP & SP:

	<u>BPE</u>	<u>WP</u>	<u>SP</u>
Features			
Core idea/concept	- It uses greedy merging of frequent char. pairs.	- It uses likelihood based sub-word merge.	- It uses unigram lang. model or BPE over raw text.
Training I/P	- It uses pre-tokenized text.	- It uses pre-tokenized text.	- It uses raw text (no preprocessing needed)
Sampling capacity	- It is deterministic.	- It is deterministic.	- It can sample multiple tokeniz.
Multi-lingual support	- Limited.	- Limited.	- Strongly supports.
Used in LLM (Eg)-	GPT with Falcon	BERT, DistilBERT, ALBERT	T5, XLM, mT5, PaLM. (developed by Google)

Word/Token embedding:



- Cbow (Continuous Bag of Word) : It predicts the current word based on surrounding context words.

- Skip gram : It predicts surrounding context word by given current word.

Eg: Generate a dense vector representation for following sentence using Cbow model.

Sentence : The cat sits on the mat.

Target word : sits

Context window : 2

Step 1: Construct a vocabulary : {the, cat, sits, on, mat}

Step 2: One-hot vector representation - the : [1, 0, 0, 0, 0]

cat : [0, 1, 0, 0, 0]

sits : [0, 0, 1, 0, 0]

on : [0, 0, 0, 1, 0]

mat : [0, 0, 0, 0, 1]

Step 3: Consider content words (the, cat, on, th) to get a combined I/P of Cbow that averages one-hot vectors

$$\text{Avg-vector} = \frac{1}{4} ([\underbrace{1, 0, 0, 0, 0}_{\text{the}}] + [\underbrace{0, 1, 0, 0, 0}_{\text{cat}}] + [\underbrace{0, 0, 1, 0, 0}_{\text{on}}] + [\underbrace{1, 0, 0, 0, 0}_{\text{th}}]) = \\ = \frac{1+0+0+1}{4}, \frac{0+1+0+0}{4}, 0, 0.25, 0 = [0.5, 0.25, 0, 0.25, 0]$$

Step 4: Initialize the weight matrix

	dim1	dim2	dim3	dim4	dim5
the	0.1	0.2	0.3	0.4	0.5
cat	0.6	0.7	0.8	0.9	0.1
sit	0.2	0.3	0.4	0.5	0.6
on	0.1	0.8	0.9	0.1	0.2
mat	0.3	0.4	0.5	0.6	0.7

Hidden layer

Step 5: Hidden layer o/p : By multiplying Avg. vector with weight matrix

$$\begin{aligned} & \rightarrow [0.6, 0.25, 0, 0.25, 0] \times [0.1, 0.2, 0.3, 0.4, 0.5] \\ & \rightarrow [0.375, 0.525, 0.725, 0.525, 0.3] \end{aligned}$$

Step 6: Initialize weight matrix for o/p layers.

let w' =	0.8	0.4	0.2	0.1	0.5
	0.3	0.6	0.9	0.8	0.2
	0.9	0.5	0.7	0.3	0.1
	0.9	0.1	0.2	0.6	0.8
	0.5	0.7	0.3	0.4	0.6

Step 7: O/p layer calculation : Hidden layer o/p $\times w'$.

$$\Rightarrow \text{Output(O/p layer)} = [1.4325, 0.9525, 1.4175, 0.9125, 0.7625]$$

Step 8: Softmax Activation

To apply softmax, need to exponentiate

$$\begin{aligned} \text{exp. scores} &= [\exp(1.4325), \exp(0.9525), \dots] \\ &= [4.1903, 2.5926, 4.1239, 2.4906, 2.1435] \end{aligned}$$

$$\begin{aligned} \text{sum. exp. scores} &= 4.1903 + 2.5926 + 4.1239 + 2.4906 + 2.1435 \\ &= 15.5409 \end{aligned}$$

Step 9: Predict the probability,

$$P = \left[\frac{4.1903}{15.5409}, \frac{2.5926}{15.5409}, \frac{4.1239}{15.5409}, \frac{2.4906}{15.5409}, \frac{2.1435}{15.5409} \right]$$

$$\downarrow \text{the} = 0.2697$$

$$\text{cat} = 0.1668$$

$$\text{sit} = 0.2653$$

$$\text{on} = 0.1603$$

$$\text{mat} = 0.1379$$

Looking inside LLMs

07-10-25

Transformer:

- A transformer is a type of AI model that learns to understand & generate human-like text by analyzing patterns in large amounts of text data.
- It processes whole sentence in parallel.

- Transformer architecture uses self attention to transform one whole sentence into a single sentence.

This is useful where older models work step by step & it helps to overcome the challenges seen in the models like RNN, LSTM (long short term memory).

Traditional models like RNN suffers from the vanishing gradient problem which leads to long term memory loss & RNN process takes one word at a time sequentially.

Eg: In the sentence "XYZ went to France in 2019 when there were no cases of covid & there he met the President of that country."

Here that country refers to France. However RNN struggle to link that country to France since it processes each word in sequence leading to loosing context over long sentences. This limitation prevents RNN from understanding the full meaning of that sentence.

While adding more memory cells in LSTM that help to address the vanishing gradient issues but still process words one by one.

This sequential process means LSTM cannot understand or analyze an entire sentence at once.

Transformers were inspired by Encoder-Decoder architecture found in RNN. However, instead of recurrence transformer model is completely based on attention mechanism.

The model doesn't generate the text all in one operation rather it generates one token at a time. After each token generation, the input prompt for the next generation step by appending the output token to the end of input prompt that means the model consumes the earlier predictions to make later prediction, are called Auto regressive models.

10-10-25

- The method of choosing a single token from the probability is called the decoding strategy.
- The easiest decoding strategy to pick the token with the highest probability score. And choosing the highest scoring token everytime is called greedy decoding.

Transformer is made up of two successive components -

- Attention layer i.e., It is mainly concerned with incorporating with relevant information from other input tokens & positions.
- The feed forward neural - token representation individually, helping the model to capture complex patterns. After attention has modelled relationship between tokens.

Attention Mechanism is a technique used in machine learning & AI to improve the performance of the model.

There are two steps involved in attention mechanisms -

(i) A way to score how relevant input of the previous token into the current token being processed.

(ii) Using those scores combine the information from various position

→ Advantages : (Attention Mechanism)

- It enhances the performance of neural machine translate system.
- Ability to address the bottle neck problem where all the information from a long input sequence needs to be compressed into a single fixed length vector.
- With attention the decoder can analyze all the encoder's hidden state with attention & dynamically weight their importance.
- It also helps to reduce the managing gradient problem by connecting all encoders to decoders.

11-10-25

Q) Consider translating the English phrase "I love bananas" into Spanish as "amo las bananas". The hidden states of the encoders are $e_1 = [0.1, 0.2]$, $e_2 = [0.3, 0.8]$, $e_3 = [0.5, 0.6]$. Compute the attention scores for the decoder, which outputs the Spanish word 'las' having a hidden state of $[0.4, 0.6]$.

The dot product (Attention score) i.e. similarity between Encoder's hidden state & Decoder's hidden state are computed as:

$$I = 0.1 \times 0.4 + 0.2 \times 0.6 = 0.04 + 0.12 = 0.16.$$

$$\text{love} = 0.3 \times 0.4 + 0.8 \times 0.6 = 0.12 + 0.48 = 0.6.$$

$$\text{bananas} = 0.5 \times 0.4 + 0.6 \times 0.6 = 0.2 + 0.36 = 0.56$$

As similarity between different hidden states are not directly comparable, therefore, normalise these score to convert the similarity score to probabilities score by using softmax function.

$$\text{Softmax}(z_i) = \frac{e^{z_i}}{\sum_{j=1}^t e^{z_j}}$$

$$\begin{aligned} \text{Softmax}(I) &= \frac{e^{0.16}}{e^{0.16} + e^{0.6} + e^{0.56}} \\ &= \frac{e^{0.16}}{1.17} \\ &= 0.25. \end{aligned}$$

$$\begin{aligned} \text{Softmax}(\text{love}) &= \frac{e^{0.6}}{e^{0.16} + e^{0.6} + e^{0.56}} \\ &= \frac{e^{0.6}}{1.82} \\ &= 0.38. \end{aligned}$$

$$\begin{aligned} \text{Softmax}(\text{bananas}) &= \frac{e^{0.56}}{e^{0.16} + e^{0.6} + e^{0.56}} \\ &= \frac{e^{0.56}}{1.75} \\ &= 0.37. \end{aligned}$$

The Attention dist/output i.e., $\alpha(t)$ gives us the fixed length input context vector for the decoder state o_t are as follows:

$$\alpha(t) = \sum_{i=1}^t \text{softmax}(i) \cdot e^{(i)} = 0.25 \times [0.1, 0.2] + 0.38 \times [0.3, 0.8] + 0.37 \times [0.5, 0.6] = [0.25 \times 0.1, 0.25 \times 0.2] + ...$$

Final decoder hidden state for 'las', after concatenation with initial decoder hidden state, o_{t-1} as $[0.32, 0.51, 0.4, 0.6] = [0.025, 0.05] + [0.114, 0.304], [0.185, 0.222], [0.324, 0.516]$

Self Attention:
It is a mechanism to capture dependencies & relationships within input sequences. It allows the model to identify & weight the importance of different parts of the i/p sequence by attending itself.

Why Self Attention is required/important:

i) Long range dependencies: It allows the model to capture relationship between distant element in a sequence.

ii) Contextual understanding: It helps the model to understand the context & assign appropriate weights to each element based on its relevance.

iii) Parallel computation: It can be computed in parallel for each element, parallel in sequence.

B-10-25

How Self attention works?

S-1: The 1st step is to generate the query (Q), key (K) & Value (V).

$$\begin{aligned} Q &= X \cdot W_Q \\ K &= X \cdot W_K \\ V &= X \cdot W_V \end{aligned}$$

S-2: Calculate attention score = QK^T

S-3: Scale the scores = $\frac{QK^T}{\sqrt{d_K}}$

S-4: Apply the softmax function to calculate attention weight = $\text{softmax}\left(\frac{QK^T}{\sqrt{d_K}}\right)$

S-5: Calculate final output (O) = Attention weight $\times V$

Eg (Single head without mask): "the cat sat". Consider the sentence, compute the self attention representation for the word 'cat'.

S-1: Let embedding size, $d_{\text{model}} = 3$

head size, $d_h = d_v = 2$

Consider, $X = \begin{bmatrix} \text{dim}_1 & \text{dim}_2 & \text{dim}_3 \\ \text{the} & \text{cat} & \text{sat} \\ 0.2 & 0.5 & 0.1 \\ \text{at} & 0.9 & 0.7 \\ 0.1 & 0.3 & 0.8 \\ \text{sat} & 0.4 & 0.7 \\ 0.8 & 0.9 \end{bmatrix}$

$$W_Q = \begin{bmatrix} 0.1 & 0.3 \\ 0.2 & 0.7 \\ 0.5 & 0.6 \end{bmatrix}, W_K = \begin{bmatrix} 0.4 & 0.8 \\ 0.9 & 0.1 \\ 0.3 & 0.5 \end{bmatrix}, W_V = \begin{bmatrix} 0.6 & 0.2 \\ 0.1 & 0.9 \\ 0.7 & 0.3 \end{bmatrix}$$

$$\begin{aligned} Q &= X \cdot W_Q \\ &= \begin{bmatrix} 0.2 & 0.5 & 0.1 \\ 0.9 & 0.1 & 0.3 \\ 0.9 & 0.7 & 0.8 \\ 0.17 & 0.47 \\ 0.26 & 0.52 \\ 0.58 & 1.09 \end{bmatrix} \begin{bmatrix} 0.1 & 0.3 \\ 0.2 & 0.7 \\ 0.5 & 0.6 \end{bmatrix} \\ &= \begin{bmatrix} 0.24 & 0.52 \\ 0.76 & 0.36 \\ 0.81 & 0.95 \end{bmatrix} \end{aligned}$$

$$\begin{aligned} K &= X \cdot W_K \\ &= \begin{bmatrix} 0.2 & 0.5 & 0.1 \\ 0.9 & 0.1 & 0.3 \\ 0.4 & 0.7 & 0.8 \\ 0.56 & 0.26 \\ 0.54 & 0.88 \\ 1.03 & 0.79 \end{bmatrix} \begin{bmatrix} 0.4 & 0.8 \\ 0.9 & 0.1 \\ 0.3 & 0.5 \end{bmatrix} \\ &= \begin{bmatrix} 0.24 & 0.52 \\ 0.76 & 0.36 \\ 0.81 & 0.95 \end{bmatrix} \end{aligned}$$

$$\begin{aligned} V &= X \cdot W_V \\ &= \begin{bmatrix} 0.2 & 0.5 & 0.1 \\ 0.9 & 0.1 & 0.3 \\ 0.4 & 0.7 & 0.8 \\ 0.24 & 0.52 \\ 0.76 & 0.36 \\ 0.81 & 0.95 \end{bmatrix} \begin{bmatrix} 0.6 & 0.2 \\ 0.1 & 0.9 \\ 0.7 & 0.3 \end{bmatrix} \\ &= \begin{bmatrix} 0.24 & 0.52 \\ 0.76 & 0.36 \\ 0.81 & 0.95 \end{bmatrix} \end{aligned}$$

S-2: Attention Score $\Rightarrow QK^T$

$$\begin{bmatrix} 0.17 & 0.47 \\ 0.26 & 0.52 \\ 0.58 & 1.09 \end{bmatrix} \cdot \begin{bmatrix} 0.56 & 0.54 & 1.03 \\ 0.26 & 0.88 & 0.79 \end{bmatrix} = \begin{bmatrix} 0.2174 & 0.5059 & 0.5469 \\ 0.2808 & 0.598 & 0.6786 \\ 0.6082 & 1.2724 & 1.4585 \end{bmatrix}$$

S-3: $\text{score} = \frac{QK^T}{\sqrt{d_k}} = \begin{bmatrix} 0.1537 & 0.3573 & 0.3863 \\ 0.1985 & 0.4228 & 0.4798 \\ 0.43 & 0.8997 & 1.0313 \end{bmatrix}$

S-4: Attention Weight $= \text{softmax}\left(\frac{QK^T}{\sqrt{d_k}}\right) = \text{the} \begin{bmatrix} x & y & z \\ \text{cat} & \text{cat} & \text{cat} \end{bmatrix} \quad x = \frac{e^x}{e^x + e^y + e^z} \Rightarrow y = \frac{e^y}{e^x + e^y + e^z}$

$$= \begin{bmatrix} 0.2867 & 0.3514 & 0.3618 \\ 0.2796 & 0.3499 & 0.3704 \\ 0.2260 & 0.3615 & 0.4124 \end{bmatrix} \begin{array}{l} \text{softmax(the)} \\ \text{softmax(cat)} \\ \text{softmax(cat)} \end{array}$$

S-5: Output (O) = Attention weight $\times V$

$$= \begin{bmatrix} 0.2867 & 0.3514 & 0.3618 \\ 0.2796 & 0.3499 & 0.3704 \\ 0.2260 & 0.3615 & 0.4124 \end{bmatrix} \times \begin{bmatrix} 0.24 & 0.52 \\ 0.76 & 0.36 \\ 0.87 & 0.95 \end{bmatrix} = \begin{bmatrix} 0.6506 & 0.6192 \\ 0.6552 & 0.6232 \\ 0.6877 & 0.6394 \end{bmatrix}$$

Self Attention representation for word cat is $[0.6552 \ 0.6232]$.

16-10-25

→ Multi-head Attention:

- It involves repeating the self-attention process multiple times, with each repetition using different linear projections of the input data.
- Using multiple attention heads allow the model to capture a richer set of dependencies in the input sequence.

Eg: One head might focus on the overall sentence structure, while another zooms in on specific details, while another zooms in on specific details.

- By combining diverse perspectives, multi-head attention provides a more comprehensive understanding of the input, much like how humans consider multiple aspects of information simultaneously.

→ Algorithmic Step(s) Workflow of Multi-head Attention :

S-1: Linear Projection: For each attention head the input sequence is linearly projected into separate queries (Q), Keys (K), Values (V) using distinct learn weight matrices.

$$[Q = XW_Q, \ K = XW_K, \ V = XW_V]$$

S-2: Attention Calculation per head (or) Scaled Dot Product Attention.

$$\text{Attention}_i(Q_i, K_i, V_i) = \text{softmax}\left(\frac{Q_i K_i^T}{\sqrt{d_k}}\right) V_i$$

- Attention score $= Q_i K_i^T$
- Scale the score $= \frac{Q_i K_i^T}{\sqrt{d_k}}$
- Apply softmax $= \text{softmax}\left(\frac{Q_i K_i^T}{\sqrt{d_k}}\right)$

$$\text{Att. wt} = \text{softmax}\left(\frac{Q_i K_i^T}{\sqrt{d_k}}\right)$$

S.3: After processing through the individual attention head, the outputs of all heads are concatenated together i.e.,

$$\text{Concatenated } \mathbf{O}_p = \text{Concat}(\text{attention}_1, \text{attention}_2, \dots, \text{attention}_n)$$

S.4: Final Linear Projection: The concatenated output is passed through another linear projection using a weight matrix to produce the final output (\mathbf{FO}).

$$\mathbf{FO} = \text{Concatenated } \mathbf{O}_p \times \mathbf{W}_o$$

Eg (Multi-head Attention): "the cat sat". Consider the sentence & compute the attention score

S.1: Consider the model with $d_{\text{model}} = 4$, no. of heads = 2, $d_k = d_v = \frac{d_{\text{model}}}{\text{heads}} = \frac{4}{2} = 2$.

$$X = \begin{bmatrix} \text{the} & 1 & 0 & 1 & 0 \\ \text{cat} & 0 & 1 & 1 & 0 \\ \text{sat} & 0 & 0 & 1 & 1 \end{bmatrix}$$

$$\mathbf{W}_Q = \mathbf{W}_K = \mathbf{W}_V = \begin{bmatrix} 1 & 0 \\ 0 & 1 \\ 0 & 0 \\ 0 & 0 \end{bmatrix}$$

$$\mathbf{Q} = X \mathbf{W}_Q$$

$$= \begin{bmatrix} 1 & 0 & 1 & 0 \\ 0 & 1 & 1 & 0 \\ 0 & 0 & 1 & 1 \end{bmatrix} \begin{bmatrix} 1 & 0 \\ 0 & 1 \\ 0 & 0 \\ 0 & 0 \end{bmatrix} = \begin{bmatrix} 1 & 0 \\ 0 & 1 \\ 0 & 0 \end{bmatrix}$$

$$\mathbf{K} = X \mathbf{W}_K$$

$$= \begin{bmatrix} 1 & 0 & 1 & 0 \\ 0 & 1 & 1 & 0 \\ 0 & 0 & 1 & 1 \end{bmatrix} \begin{bmatrix} 1 & 0 \\ 0 & 1 \\ 0 & 0 \\ 0 & 0 \end{bmatrix} = \begin{bmatrix} 1 & 0 \\ 0 & 1 \\ 0 & 0 \end{bmatrix}$$

$$\mathbf{V} = X \mathbf{W}_V$$

$$= \begin{bmatrix} 1 & 0 & 1 & 0 \\ 0 & 1 & 1 & 0 \\ 0 & 0 & 1 & 1 \end{bmatrix} \begin{bmatrix} 1 & 0 \\ 0 & 1 \\ 0 & 0 \\ 0 & 0 \end{bmatrix} = \begin{bmatrix} 1 & 0 \\ 0 & 1 \\ 0 & 0 \end{bmatrix}$$

S.2: Attention score = $\mathbf{Q}\mathbf{K}^T$ (for heads)

$$= \begin{bmatrix} 1 & 0 \\ 0 & 1 \\ 0 & 0 \end{bmatrix} \begin{bmatrix} 1 & 0 & 0 \\ 0 & 1 & 0 \\ 0 & 0 & 0 \end{bmatrix}^T = \begin{bmatrix} 1 & 0 & 0 \\ 0 & 1 & 0 \\ 0 & 0 & 0 \end{bmatrix}$$

$$\text{Scale the score} = \frac{\mathbf{Q}\mathbf{K}^T}{\sqrt{d_k}} = \begin{bmatrix} 0.7071 & 0 & 0 \\ 0 & 0.7071 & 0 \\ 0 & 0 & 0 \end{bmatrix}$$

$$\text{Apply softmax} = \begin{bmatrix} 0.5039 & 0.2482 & 0.2482 \\ 0.2482 & 0.5034 & 0.2482 \\ 0.3333 & 0.3333 & 0.3333 \end{bmatrix}$$

$$e^{0.7071} = 2.028$$

$$\text{Attention weight} = \begin{bmatrix} 0.5039 & 0.2482 & 0.2482 \\ 0.2482 & 0.5034 & 0.2482 \\ 0.3333 & 0.3333 & 0.3333 \end{bmatrix} \begin{bmatrix} 1 & 0 \\ 0 & 1 \\ 0 & 0 \end{bmatrix} = \begin{bmatrix} 0.5034 & 0.2482 \\ 0.2482 & 0.5034 \\ 0.3333 & 0.3333 \end{bmatrix}$$

7-10-25

$$\text{For head 2, } \mathbf{X} = \begin{bmatrix} 1 & 0 & 1 & 0 \\ 0 & 1 & 1 & 0 \\ 0 & 0 & 1 & 1 \end{bmatrix} \quad \mathbf{W}_Q = \mathbf{W}_K = \mathbf{W}_V = \begin{bmatrix} 0 & 0 \\ 0 & 0 \\ 0 & 1 \\ 1 & 0 \end{bmatrix}$$

$$Q = XW_Q \quad [Q = K = V = XW_K = XW_V]$$

$$= \begin{bmatrix} 1 & 0 & 1 & 0 \\ 0 & 1 & 1 & 0 \\ 0 & 0 & 1 & 1 \end{bmatrix} \begin{bmatrix} 0 & 0 \\ 0 & 0 \\ 0 & 1 \\ 1 & 0 \end{bmatrix} \quad QK^T = \begin{bmatrix} 0 & 1 \\ 0 & 1 \\ 1 & 1 \end{bmatrix} \begin{bmatrix} 0 & 0 & 1 \\ 1 & 1 & 1 \end{bmatrix} = \begin{bmatrix} 1 & 1 & 1 \\ 1 & 1 & 2 \end{bmatrix}$$

$$= \begin{bmatrix} 0 & 1 \\ 0 & 1 \\ 1 & 1 \end{bmatrix} \quad \frac{QK^T}{\sqrt{d_K}} = \begin{bmatrix} 0.7071 & 0.7071 & 0.7071 \\ 0.7071 & 0.7071 & 0.7071 \\ 0.7071 & 0.7071 & 0.4142 \end{bmatrix}$$

Apply softmax -

$$\begin{bmatrix} 0.3333 & 0.3333 & 0.3333 \\ 0.3333 & 0.3333 & 0.3333 \\ 0.2482 & 0.2482 & 0.5034 \end{bmatrix}$$

Attention weight =

$$\begin{bmatrix} 0.3333 & 0.3333 & 0.3333 \\ 0.3333 & 0.3333 & 0.3333 \\ 0.2482 & 0.2482 & 0.5034 \end{bmatrix} \begin{bmatrix} 0 & 1 \\ 0 & 1 \\ 1 & 1 \end{bmatrix} = \begin{bmatrix} 0.3333 & 0.9999 \\ 0.3333 & 0.9999 \\ 0.5034 & 0.9998 \end{bmatrix}$$

S-3: Head Concatenation

Concatenated o/p =

$$\begin{bmatrix} 0.5034 & 0.2482 & 0.3333 & 0.9999 \\ 0.2482 & 0.5034 & 0.3333 & 0.9999 \\ 0.3333 & 0.3333 & 0.5034 & 0.9998 \end{bmatrix}$$

S-4: Final Output = Concatenated o/p x W_o (Generate FO)

$$\begin{bmatrix} 0.5034 & 0.2482 & 0.3333 & 0.9999 \\ 0.2482 & 0.5034 & 0.3333 & 0.9999 \\ 0.3333 & 0.3333 & 0.5034 & 0.9998 \end{bmatrix} \begin{bmatrix} 1 & 0 & 0 \\ 0 & 1 & 0 \\ 0 & 0 & 1 \\ 0 & 0 & 0 \end{bmatrix}$$

$$= \begin{bmatrix} 0.5034 & 0.2482 & 0.3333 \\ 0.2482 & 0.5034 & 0.3333 \\ 0.3333 & 0.3333 & 0.5034 \end{bmatrix}$$

23-10-25

Difference between Single Head, Multi Head & Masked Attention:

Single-Head

- It uses single-attention mechanism to compute attention weight & generate a context vector.
- It may struggle to capture complex relationships between input elements.
- Only one set of query, key & value weight matrices is used.
- Attention $(Q, K, V) = \text{Softmax}\left(\frac{QK^T}{\sqrt{d_k}}\right) \times V$

v) Eg: In the sentence, 'The cat sat on the mat', Single-head attention might focus mainly on Subject + verb i.e., cat-cat but miss out 'on-mat' relationship.

Multi-Head

- It uses multiple attention mechanism in parallel, each computing attention weight & generating a context vector.
- It captures more complex relationship between input element by combining the outputs of multiple attention heads.
- Multiple attention heads use separate set of query, key & value weight matrices & they operate in parallel & their ops. are concatenated at the end.
- Multihead $(Q, K, V) = \text{Concat}(h_1, h_2, \dots, h_n) \cdot W_o$

✓) Eg: In this case, One head might focus on Subject + verb i.e., cat-sat relationship, another head focused on prepositions like 'on-mat', another might learn co-references like 'The-cat'.

Masked Attention

- It is used in auto-regressive model where future are masked to prevent the model from attending to them.
- It is used to prevent attending to future tokens while single & multi-head attention do not use masking by-default.
- It is same as multi-head attention last before applying Softmax, a mask is added that sets scores of future position to -∞.
- Masked attention $(Q, K, V) = \text{Softmax}\left(\frac{QK^T}{\sqrt{d_k}} + \text{mask}\right) \times V$

Q) Why the Mask is needed?

- In Masked Attention, the goal is to prevent a word from 'seeing' future words during training / generation.
- This is important because the model must learn to predict the next token only from the tokens it has already seen but not by looking ahead at future tokens.

Eg. "The cat sat on the mat"

Step Current Token being generated

1	the
2	cat
3	sat
4	on
5	the
6	mat

Allowed to attend

- only itself
[the, cat]
- [the, cat, sat]
- [the, cat, sat, on]
- [the, cat, sat, on, the]
- [the, cat, sat, on, the, mat]

Masked effect

- nothing to be masked yet
- can't look at 'sat', 'on', 'the', 'mat'
- can't look at 'on', 'the', 'mat'
- can't look at 'the', 'mat'
- can't look at 'mat'

Final

⇒ Mathematical Representation of Masked Attention :

the	cat	sat	on	the	mat
the	0	-∞	-∞	-∞	-∞
cat	0	0	-∞	-∞	-∞
sat	0	0	0	-∞	-∞
on	0	0	0	0	-∞
the	0	0	0	0	0
mat	0	0	0	0	0

0 → Visible

-∞ → Not Visible

24-10-25

Chapter-01 Text Classification

Part-2

→ Differentiate Text Classification with Representation type model & Generative type model:

Representation type Model

- i) It uses models primarily designed to create a numerical (meaningful) representation of input text. These representation/embedding capture the semantic meaning & content of text.
- ii) In this model, an encoder is used to generate contextualised embedding vector & then pass through a task specific layer also called classification head to map numerical vector to final class probability.
- iii) To find the optimal decision boundary in the high dimensional embedding space that cleanly separate the classes, this makes them discriminative model. (GOAL)

Eg: y_p : "Congratulations, You won a free iPhone."
 y_o : Spam

Models Used: BERT, ROBERTA

25-10-25
→ Text Classification with Representation Models:

Task Specific Model (label)

- i) It is trained to perform a particular NLP task such as Sentiment Analysis, Text classification, Summarization, etc. It is fully trained or fine-tuned on label data for that specific task.
- ii) The main purpose is to directly output the final result or decision for the specific application.
- iii) The goal is to perform a defined task.
- iv) It is optimized for one task.

→ Supervised Classification in Representation Models:

It refers to a process of training a model to map input data to meaningful representation (features/embeddings) & then assigning class labels based on those representations.

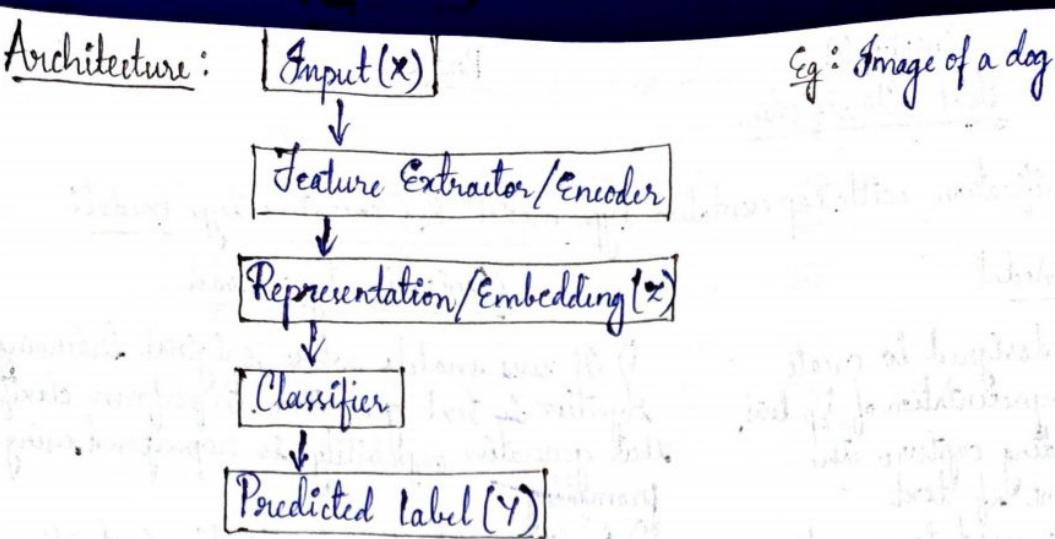
Generative type Model

- i) It uses models whose original training objective is text generation. To perform classification, their generative capability is repurposed using prompting.
 - ii) In this model, the classification instruction is embedded directly into the input text, a process called prompt engineering.
 - iii) To generate a textual response, that fulfills the instruction in the prompt. (GOAL)
- Eg: y_p : "Congratulations, You won a free iPhone"
 y_o : This msg looks like spam or fishing attempt.

Models Used: GPT, T5

Embedding Model (vector)

- i) It is trained to produce dense-vector representation for words, sentence & documents capturing semantic meanings.
- ii) The main purpose is to convert text into numerical form that can be reused for many tasks.
- iii) The goal is to represent meaning in vector form.
- iv) It is reusable for many tasks.



27-10-25

Q) What if we do not have a label data?

A: If we do not have label data we can't train a supervised classifier directly rather we can use supervised, unsupervised or transfer learning to train or reuse a representation model without labels and optionally add a small amount of labels or weak labels for final classification.

1 Solution to this problem is - use of zero-shot model.

+ one-shot model
+ few-shot model

→ How Zero-shot model works?

- Most modern zero-shot model is a model that can make prediction on new classes/tasks without having been explicitly trained on them.
- Most modern zero-shot models (eg. clip) [one-shot model] (eg. Flamingo) [few-shot model] are trained on very large dataset of images & texts together (eg., an img of a dog & the caption - the dog is running).
- During training they learn to align image features & text features in the same embedding space, this allows the model to recognise classes without seeing labeled training data.

→ Step 1: It encodes the image into an image embedding.

Step 2: It encodes the text label into text embedding.

Step 3: It computes cosine similarity between image & text embedding.

Step 4: It selects closest text as the predicted class.

→ How cosine similarity between image & text embedding works?

$$\text{cosine similarity } (A, B) = \frac{A \cdot B}{\sqrt{\sum_{i=1}^n A_i^2} \cdot \sqrt{\sum_{j=1}^m B_j^2}}$$

Eg: Consider the following images & calculate similarity among them to provide label info.

Image 1 (dog), i.e., $A = [0.2, 0.5, 0.8, 0.1]$

Image 2 (another dog), i.e., $B = [0.3, 0.1, 0.7, 0.2]$

Image 3 (cat), i.e., $C = [0.9, 0.1, 0.2, 0.3]$

$$\text{cosine sim}(A, B) = \frac{0.2 \times 0.3 + 0.5 \times 0.1 + 0.8 \times 0.7 + 0.1 \times 0.2}{\sqrt{0.2^2 + 0.5^2 + 0.8^2 + 0.1^2} \times \sqrt{0.3^2 + 0.4^2 + 0.7^2 + 0.2^2}} = \frac{0.81}{0.96 \times 0.8839} = 0.98 = 98\%$$

$$\text{cosine sim}(A, C) = \frac{0.42}{0.9695} = 0.45 = 45\% \quad | \quad \text{cosine sim}(B, C) = \frac{0.51}{0.88 \times 0.97} = 0.59 = 59\%$$

Text clustering & Topic modelling

Text clustering & Topic modelling are unsupervised analysis technique that aims to discover structure & pattern within large collection of documents while they share the same goal i.e grouping or organizing text data but they differ in their approach, representation & interpretation.

→ Difference between Text Clustering & Topic Modelling :Text Clustering

- i) It aims to group similar text based on their semantic content, meaning & relationship.
- ii) The goal is to group similar text.
- iii) It is a distance / similarity based approach that groups documents into clusters.
- iv) It is usually a hard clustering method where each document belongs to one cluster.

Topic Modelling

- i) It refers to process of representing, generating or understanding text using statistical / neural models.
- ii) The goal is to represent/generate text patterns.
- iii) It is a probabilistic approach that assumes each document is a mixture of multiple topics. Each topic is a distribution of words. It covers latent (hidden) semantic structure.
- iv) It can be viewed as a soft clustering method where document belongs to multiple topics with varying probabilities.

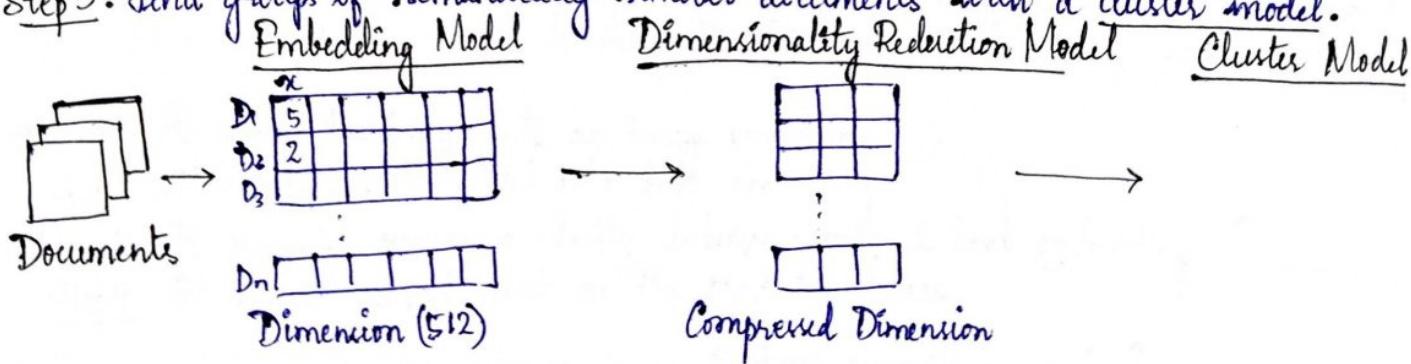
→ A common pipeline for Text Clustering :

There are 3 steps in common pipeline for Text Clustering :

Step 1: Convert the input documents to embedding with an embedding model.

Step 2: Reduce the dimensionality of embeddings with a dimensionality reduction model.

Step 3: Find groups of semantically similar documents with a cluster model.

→ Dimensionality Reduction Model

As the no. of dimensions increases, there is an exponential growth in no. of possible values within each dimension. Finding all sub-spaces within each dimension becomes increasingly complex, as a result it is difficult to identify meaningful clusters.

Uniform Manifold Approximation & Projection (UMAP)

UMAP is a non-linear dimensionality reduction technique that helps to represent high dimensional data in a lower dimensional space while preserving the structure & relationship in the data as much as possible. For mathematical transformation it uses non-linear optimisation, graph based or distance preserving method.

The main core mechanism of UMAP is to built a graph representation of data in high dimensions & then optimizes a low dimensional projection that maintains both -

- (a) Local structure - nearby points stay close.
- (b) Global structure - overall data distribution remains meaningful.

UMAP works in two phases -

- (a) Constructing a high dimensional graph.
- (b) Optimizing the lower dimensional embeddings.

Ex: Consider the following 6 points in 3D to form two clusters where each point A,B,C,D,E,F represents a data sample in HD space.

Points	x	y	z
A	0.0	0.0	0.0
B	0.2	0.1	0.0
C	0.1	-0.2	0.1
D	3.0	3.1	3.0
E	3.2	2.9	3.1
F	2.9	3.2	3.2

Step 1: Compute Euclidean distance

$$D(A,B) = \sqrt{(x_2 - x_1)^2 + (y_2 - y_1)^2 + (z_2 - z_1)^2}$$

Point1	A	B	C	D	E	F
A	0	0.2236	0.2449	5.2545	5.3160	5.3749
B	0.2236	0	0.3317	5.0833	5.1429	5.2096
C	0.2449	0.3317	0	5.2640	5.3122	5.3860
D	5.2545	5.0833	5.2640	0	0.3	0.2449
E	5.3160	5.1429	5.3122	0.3	0	0.4358
F	5.3749	5.2096	5.3860	0.2449	0.4358	0

Step 2: Find k-Nearest Neighbour where k=2.

Points	1 st Nearest	2 nd Nearest
A	B(0.2236)	C(0.2449)
B	A(0.2236)	C(0.3317)
C	A(0.2449)	B(0.3317)
D	F(0.2449)	E(0.3)
E	D(0.3)	F(0.4358)
F	D(0.2449)	E(0.4358)

03-11-25

Step 3: Compute local parameters to check local connectivity strength.

UMAP uses two parameters for each point:

- i) $f(\text{rho})$ - that measures the distance to the nearest neighbour.
- ii) σ (sigma) - scaling factor to make total local probabilities i.e. $\sigma = 0.5$ (Here).

To compute the local connectivity strength i.e. affinity from i to j is

$$a_{ij} = \begin{cases} 1 & \text{if } d_{ij} \leq f_i \\ e^{-(\frac{(d_{ij}-f_i)}{\sigma})} & \text{if } d_{ij} > f_i \end{cases}$$

For Cluster-1: For A, $f_A = 0.2236$

For A \rightarrow B, $a_{AB} = \frac{1}{e^{-(\frac{(0.2119-0.2236)}{0.5})}} = 0.2149 > f_A$ [i.e. 0.2236]

For A \rightarrow C, $a_{AC} = e^{-\frac{0.5}{0.961}} = 0.961$

For B, $f_B = 0.2236$

For B \rightarrow A, $a_{BA} = \frac{1}{e^{-(\frac{0.2119-0.2236}{0.5})}} = 0.2149$ ($d=0.22$)

For B \rightarrow C, $a_{BC} = e^{-\frac{0.5}{0.802}} = 0.802$ ($d=0.33$)

For C, $f_C = 0.2449$

For C \rightarrow A, $a_{CA} = 1$ ($d=0.24$)

For C \rightarrow B, $a_{CB} = e^{-\frac{0.5}{0.836}} = 0.836$ ($d=0.33$)

06-11-25

For Cluster-2: For D, $f_D = 0.2449$

For D \rightarrow E, $a_{DE} = 1$

For D \rightarrow F, $a_{DF} = e^{-\frac{(0.3-0.2449)}{0.5}} = 0.895$

For E, $f_E = 0.3$

For E \rightarrow D, $a_{ED} = 1$

For E \rightarrow F, $a_{EF} = e^{-\frac{(0.4358-0.3)}{0.5}} = 0.7621$

For F, $f_F = 0.2449$

For F \rightarrow D, $a_{FD} = 1$

For F \rightarrow E, $a_{FE} = e^{-\frac{(0.4358-0.2449)}{0.5}} = 0.682$

Step 4: Symmetrize the Affinities

UMAP symmetrizes directed affinities to get an undirected fuzzy graph using:

$$W_{ij} = a_{ij} + a_{ji} - a_{ij} \times a_{ji}$$

For Cluster 1: Pair

A-B

formula

$$1+1-1 \times 1 = 1.00$$

A-C

$$0.961+1-0.961 \times 1 = 1.00$$

B-C

$$0.802+0.836-0.802 \times 0.836 = 1.638-0.671 = 0.967$$

24-11-25

Chapter-06

Prompt Engineering

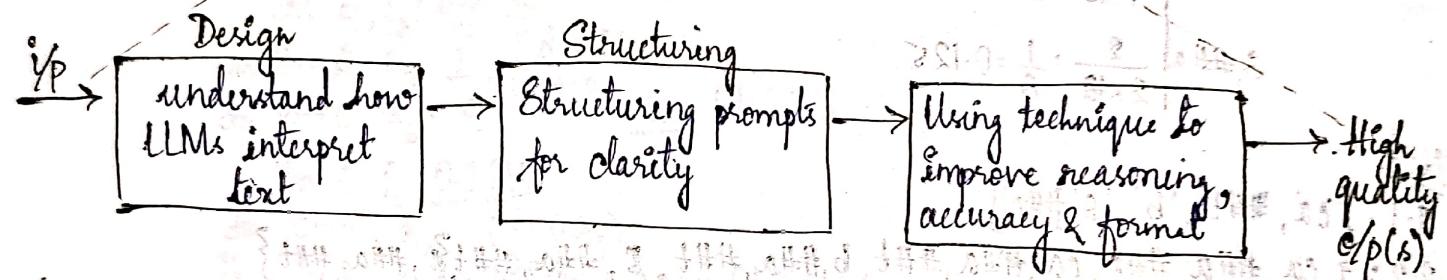
⇒ What is prompt?

A prompt is an input that user ~~to~~ give ^{to} an AI model to make it generate a response.

⇒ Prompt Engineering:

It is a process of designing, structuring & optimizing prompts to LLM so that the model produces accurate, relevant & high quality outputs.

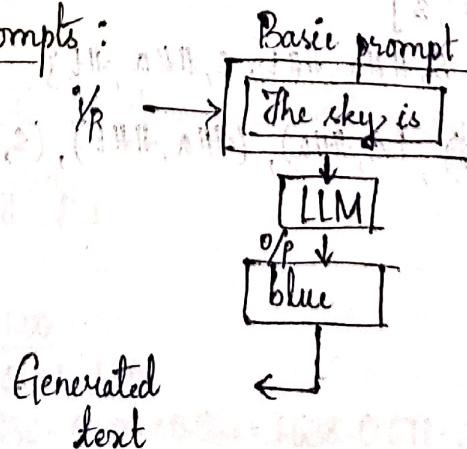
Prompt Engineering



⇒ Why Prompt Engineering is important?

- For Better output: Good prompts lead to more accurate, relevant & useful responses.
- Clearer communication: It helps the model to understand what you want, by reducing misunderstandings.
- Efficient interaction: Effective prompts save time & efforts, i.e., to get the information that you need faster.
- Unlocking capabilities: It can tap into the model's full potential, revealing hidden strengths & abilities.
- Reducing error: Thoughtful prompts minimize mistakes, biases & off-topic responses.

⇒ Components of an effective prompt:



- 1) Introduction: Tells the model exactly what you want it to do.
- 2) Context: Give any background, constraints that shape the answer.
- 3) Input data: Provide the text of data the model should work on.
- 4) Output Format: Specify how do you want the results.
- 5) Example: A quick help to the model match the style hereafter.

27-11-25

Q) How to control model's output?

A: To control the model's output, adjust the model's parameter like Temperature, Top k, Top P.

Temperature: It controls random or deterministic the model's next token.

- It controls the degree of randomness in token selection. It defines how likely it is to choose tokens that are less probable.
- Lower temperatures are good for prompts that expect a more deterministic response, while higher temperatures can lead to more diverse or unexpected results.
- Use a low temperature for factual answers, while a higher one for stories or brainstorming.

Top k: It controls the top k probable tokens before sampling.

- Top k selects the top k most likely tokens from the model's predicted distribution.
- The higher top k is the more creative & varied model's output, while the lower top k is more restrictive & factual the model's output.
- A top k of one is equivalent to greedy decoding.

Top P: It controls to choose tokens from the smallest set whose cumulative probability $\geq P$

- Top P sampling selects the top tokens whose cumulative probability doesn't exceed a certain value P .
- Instead of a fixed K , set a cumulative probability threshold P .
- The model selects a smallest set of token whose total probability $\geq P$.

) Input = "Our new smartwatch is packed with features including.."

Model's raw token probabilities (Top 6) are :

Token	Prob
Health	0.35
Fitness	0.22
Sleep	0.18
Notification	0.12
GPS	0.08
Battery	0.05

let $P = 0.5 \Rightarrow 0.57$ (Cumulative score of Health & Fitness)

let $P = 0.8 \Rightarrow 0.87$ (Cum. score of Health, Fitness, Sleep, Notification)

26-11-25 ⇒ Advanced Prompt Engineering :

- It is a process of designing smart, multi-stepped & structured prompts that help LLMs to solve complex problems with higher accuracy & better reasoning.
- Prompt Engineering becomes advanced when it goes beyond simple instructions & uses strategic, multi-step & structure techniques to perform LLM to perform complex, accurate tasks instead of giving commands, advanced prompt engineering focuses on how the model thinks & verifies the output.
- Advanced Prompt Engineering uses methods like -
 - i) Zero-shot prompting: It is the simplest type of prompt that provides a description of a task & some texts for the LLM to get started with.
 - ii) Few-shot prompting: Two or more examples are provided to demonstrate desired behavior.
 - iii) Instruction prompting: Providing a clear & explicit set of rules or commands. For e.g. Summarise the text in 3 points within 50/100 words.
 - iv) Role based prompting: It is a technique in prompt engineering that involves assigning a specific role to the AI model, this can help the model to generate more relevant & informative output.
 - v) Chain of Thought prompting (CoT): It is a technique for improving the reasoning capabilities of LLMs by generating intermediate reasoning steps, this helps the LLM to generate more accurate answers.
 - vi) Tree of Thought prompting (ToT): It encourages to explore the multiple solution paths.
 - vii) Self Consistency prompting: Using the same prompt multiple times can lead to different results. To resolve this degree of randomness & improve the performance of generative models, Self Consistency concept was introduced. This method asks the generative model the same prompt multiple times & takes the majority result as final answer.

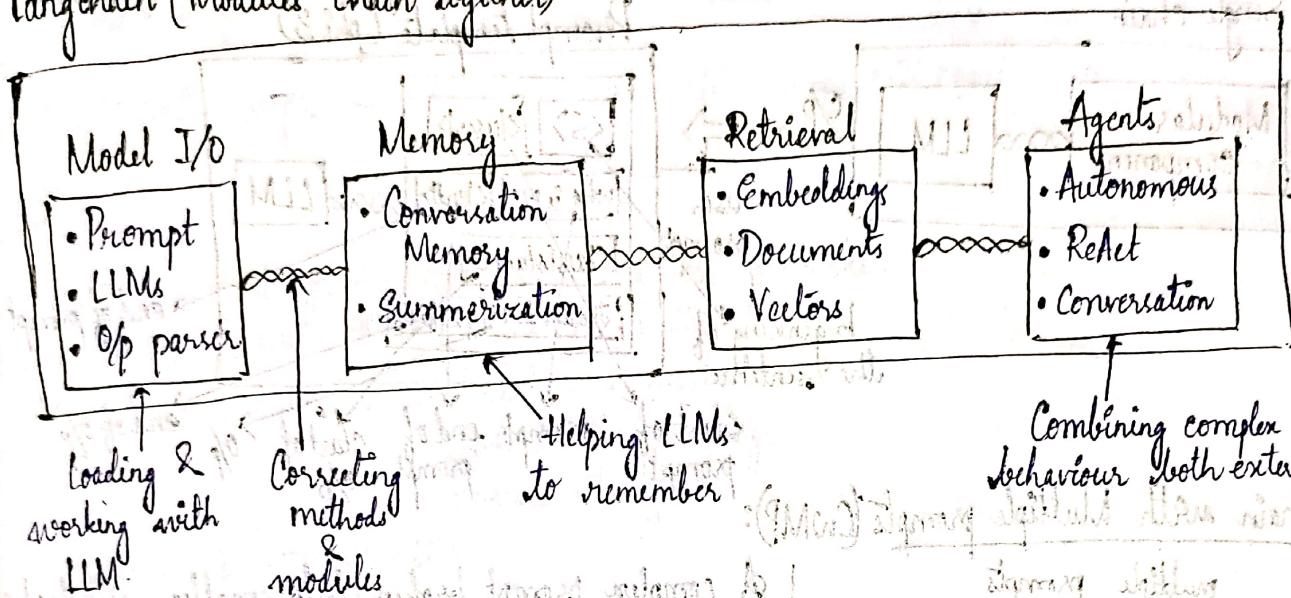
46.9	Marl
28.3	world
31.0	good?
21.0	strategic
30.0	99%
30.0	million

Advanced Text generation Technique & Tools

This chapter explores the methods to improve the quality of the generated text.

→ Framework

Langchain (Modules chain together)



⇒ Model I/O: Loading quantized models with Langchain

i) It refers to how we can load, train & manage input/output of language models.

ii) It deals with

- a) Loading models from local files or hubs.
- b) Using transformers library to handle tokenizers & model inference.
- c) Co-ordinating I/p text → tokens → Model → O/p tokens → Final text.

iii) Quantized models are machine learning or LLM models that has been compressed by reducing the precision of its numeric values (weights) to make the model smaller, faster & cheaper to run.

Eg: Original weight (FP32) = 0.092345789

$$\begin{aligned} &\text{Quantized into INT8} \\ &= \underline{\underline{0.09}} \end{aligned}$$

iv) Quantization is the process of reducing model with precision from 32-bit floats or FP32 to lower bit format i.e., INT8, INT4, INT2.

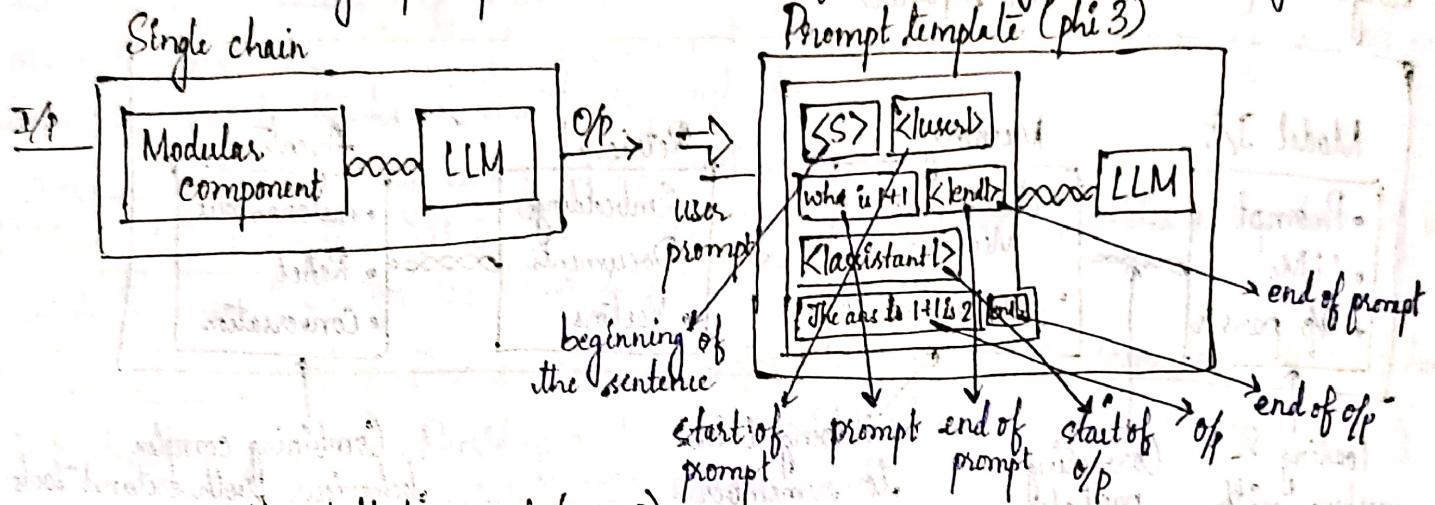
01-12-25

⇒ Chain: Extending the capabilities of LLM:

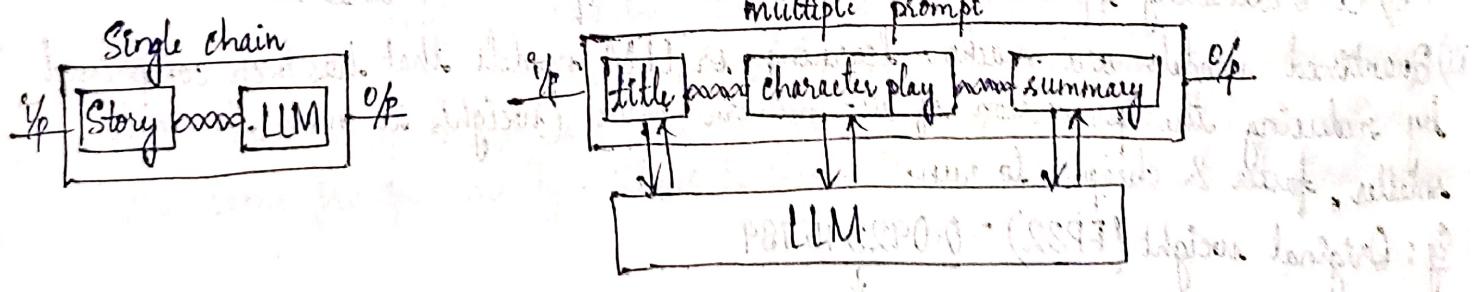
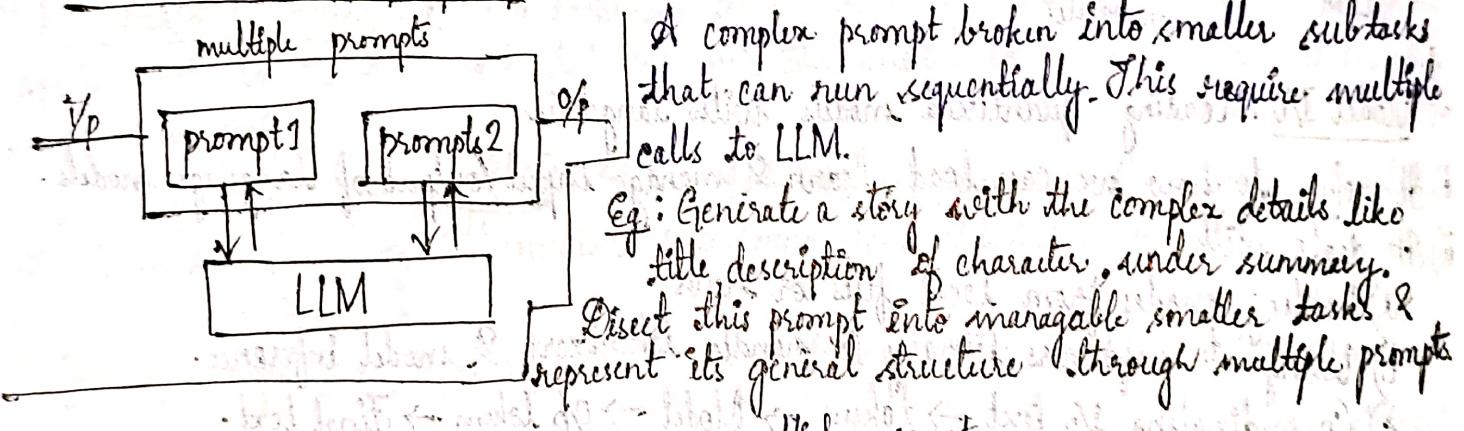
Chain → most imp/core functionality of Langchain

main method

- LangChain is named after its main method i.e., Chain.
- Chains are multi-step sequences that connect several operations together to complete a task using an LLM. Each step can be a prompt, a model call, a function call or tool, a data preprocessing step, etc. A chain allows LLM to solve complex tasks that cannot be solved in a single prompt. The most basic form of LangChain is a Single chain.



A Chain with Multiple prompts (CwMP):



04-12-25 → Differentiate between CoT, Chain with Multiple prompt, Prompt Chaining :

COT
Def: It is a method where prompt produces intermediate reasoning steps to reach final answer.

Purpose:
• Tasks model to think step by step.

No. of prompts:
• 1, One prompt
No. of model calls:
• 1, One call

<u>CwMP</u>	<u>PC</u>
A single task is completed using multiple prompts, sent in separate steps & the O/P of one model call or prompt template becomes the I/P of the next.	A full task is broken into a sequence of steps & the O/P of one model call becomes the I/P of the next.
• Multiple prompts stored in 1 template	• Split complex tasks into steps
• Multiple prompts for different steps	• Multiple prompts
• 1, one call	• 1, one call

COT

Structure:
 1. prompt (think step-by-step)
 ↓ one call
 LLM generate COT + Final o/p

CwMP

T₁: prompt 1 + prompt 2 (concatenated)
 ↓ combined call
 LLM o/p

PC

T₁: prompt 1 → LLM → o/p 1
 T₂: prompt 2 (using o/p 1) → LLM → o/p 2
 T₃: prompt 3 (using o/p 2) → LLM → Final o/p

Pros:

- It is simple & fast.

Cons:
 Harder to control steps.

- Easy to iterate & interact.

- Depends on conversation memory.

- Highly controlled & good for automation.

- It requires design effort.

05-12-25

Q) Generate a student grade report for the following marks:
 Math - 45/50, Science - 40/50, using COT, CwMP & PC.

A: Student grade report :
 1. Calculate total mark
 2. Findout the percentage
 3. Generate the grade

$$\text{Total mark} = 50 + 50 = 100$$

$$\text{Marks obtained} = 45 + 40 = 85$$

$$\text{Percentage} = \frac{85}{100} \times 100 = 85\% \rightarrow \text{B grade}$$

COT

prompt → Calculate the tot-mark
 % & grade for math - 45/50, math - 45/50, sci - 40/50

↓ think step by step & show
 the reasoning steps.

LLM o/p: Step1 - Total mark = 85

Step2 - Maximum mark = 100

Step3 - Percentage = 85%

Step4 - Grade = B

Final answer - Total = 85, Percentage = 85

Grade = B

CwMP

<task> $\langle \text{back} \rangle \langle \text{marks} \rangle$
 Generate grade report for marks

$\langle \text{instruction} \rangle$
 1. Compute Tot mark
 2. Find out the percentage
 3. Assign grade

$\langle / \text{instruction} \rangle$

LLM o/p: Total mark = 85

Percentage = 85%

Grade = B

Call 1: Extract raw data

prompt1: Extract the marks of
 math 45/50 & sci 40/50

LLM o/p 1: math - 45

sci - 40

Call 2: Calculate tot-mark & %

prompt2: math - 45

sci - 40

Calculate tot-mark & %

LLM o/p 2: Total marks = 85

Percentage = 85%

Call 3: Assigning grade

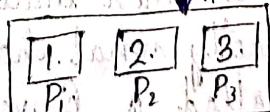
prompt3: Assigning grade for totmark

LLM o/p 3: Percentage = 85%
 Grade = B

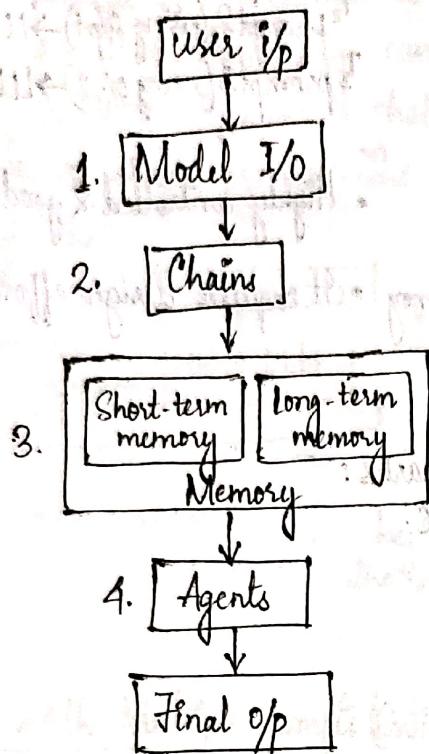
Final o/p: Total marks = 85

Percentage = 85%

Grade = B



06-12-25



3.

→ Memory: Memory stores & retrieves information, so that LLM can maintain content across steps, sessions or tasks in langchain. There are two types of memory:

- i) Short-term memory: It helps the LLM to remember past earlier discussions (past few min/hours). It stores conversation history & recent content.
- ii) Long-term memory: It helps to retrieve the facts & the knowledge. It doesn't forget the conversation after the chat ends. It uses vector database for storing of information (RAG).

06-12-25

(continuation to previous content - Long term mem)

There are 2 common methods for helping LLMs to remember conversations.

- i) Conversation-Buffer: One of the most intuitive forms of giving LLMs memory is simply reminding them exactly what has happened in the past. In Langchain this form of memory is called a Conversation-Buffer-Memory.
 - a) Windowed Conversation Buffer: It is the method used in ChatBots or LLM appⁿ to store only the recent part of conversations instead of entire chat history. It is mainly used to reduce memory usage, control context length & improve performance.
 - ii) Conversation Summary: It is a technique used to manage long chats by compressing earlier parts of the conversation into a short meaningful summary. This helps the model stay consistent without exceeding token limits.
- Why it is needed: The Conversation Summary help with:

- i) Joken efficiency: Summaries are far shorter than full history.
- ii) Maintaining coherence: The model remember user goals & previous context.
- iii) Scaling to long conversations: It enables hrs, days of conversation history.
- iv) Supporting agent based system: Agents rely heavily on memory to plan multi-stepped tasks.

4. Agents in LangChain:
→ agents: An Agent is a LLM powered decision maker that chooses what to do next based on the user's i/p & the tools available.
• An Agent in langChain is LLM + reasoning loops that lets the model choose the next action dynamically.

→ Why agents are needed?

Agents solve the problem by:

- i) Understanding the user's query
- ii) Deciding which tool to use
- iii) Calling that tool
- iv) Looking at the result
- v) Deciding the next step
- vi) Repeating until the task is completed

This makes the agent dynamic, flexible & intelligent.

11-12-25

→ Components of Agent in langChain:

- i) LLM (Brain): The language model that thinks & decides the next steps.
- ii) Tools (Actions): Calculated tool, Websearch tool, database tool, python tools.
- iii) Agent Executor: The controller that loops like think → act → observe → think.

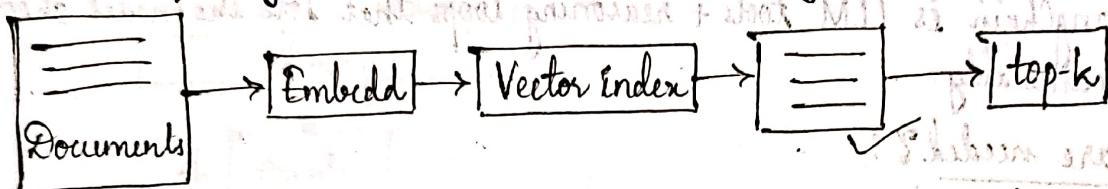
→ ReAct: It is combination of two words - Reasoning & Action.

It is a process or framework or methodology that defines how a LLM should behave when solving problems, that requires step by step reasoning, calling external tools or actions, reading results then decide the next action. Without LLM, ReAct cannot work & without ReAct, LLM cannot dynamically act with tools.

Togetherly LLM & ReAct create an autonomous problem solving system & it tells the LLM how to think & act & the LLM provides the intelligence to follow the ReAct workflow.

Semantic search & Retrieved Augmented Generation (RAG)Semantic Search:

- It is a way of retrieving information based on meaning rather than just keyword matching.
- It converts queries & documents into vector embeddings that capture conceptual similarity.
- The system then compares these vectors & returns results that are semantically related, even if they use different bodies of wordings.

Difference between Keyword Search & Semantic Search:Keyword SearchFeatures:

- Purpose • Looks for exact or stem words or phrases
matches in the index.

- Tech. used • BoW, TF-IDF, BM25, Inverted Index, Keyword lookup.

- Understand synonyms • It doesn't understand synonym of the word.

- Strength • Simple & fast to compute; transparent, & works well with the rare terms.

- Weakness • It misses related concepts; sensitive to spelling or word choice.

- Search space • Lexical or words.

- Qf quality • Depends on exact word overlap.

Semantic Search

- It converts text into dense vectors & compares the meaning using similarity score.

- Embeddings, Vector Similarity, Neural Models.

- It understands synonyms of the word.

- It captures synonym, paraphrases & context; handles fuzzing intent.

- It requires a good embedding model & more storage.

- Semantic or meaning

- Depends on conceptual similarity.

Semantic Search: It is of two types:-

- (i) Dense retrieval: It is a retrieval method that uses dense vector embedding to find documents that are semantically similar to a user query.

→ Why it is required? - Because it understands meaning but not just exact words. It handles synonyms. Works well for long, natural language queries & also captures contexts & relationships between terms.

Limitation -

- Answer must exist in corpus(db).

- DR cannot retrieve correct answers if the infⁿ isn't present in index-text).

- Weak at exact matches.

- (DR may miss precise keywords, phrases, names or codes).

- Domain Sensitivity
(Performance drops when embeddings are applied to domains different from their training data).
- Chunking Dependent
(Retrieval quality heavily depends on how documents are split into chunks).

13-12-25

Ex: Consider a small document & process using Dense retrieval concept.

Doc - Semantic search helps computers understand meaning. Dense retrieval uses vector embeddings. Text is broken into chunks for efficient search.

Step 1: Break the document into chunks.

Semantic search system always splits long text into small pieces called chunks, because embedding models work best with short & focused text. The chunks are -

Semantic search helps

chunk 1

Dense retrieval uses

chunk 2

Text is broken..

Step 2: Convert each chunks into embedding vectors, i.e., get to vector and measure

C1: [0.21 | 0.55 | 0.13 | ...]

C2: [0.89 | 0.1 | 0.66 | ...]

C3: [0.32 | 0.69 | 0.05 | ...]

Step 3: Process a query

let Q: What uses embeddings?

Step 4: Convert a query into an embedding form

Q: [0.26 | 0.42 | 0.51 | ...]

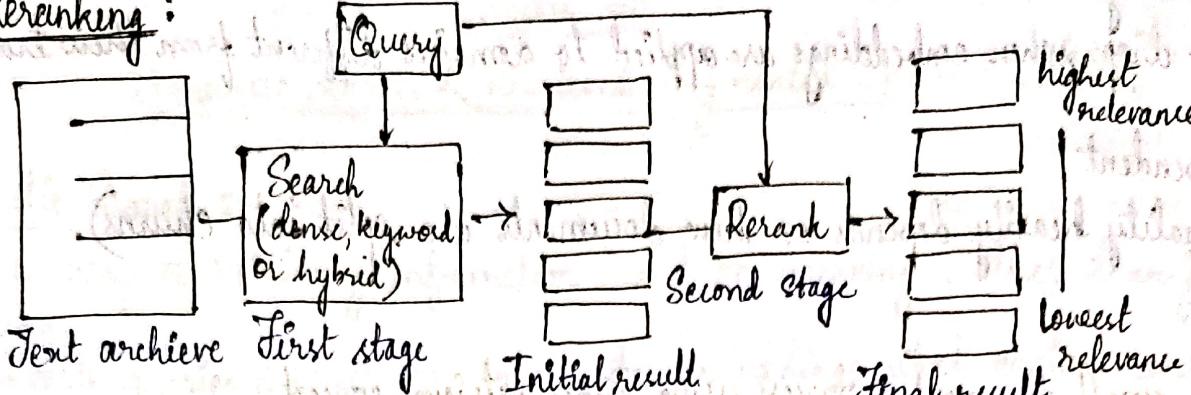
Step 5: Dense retrieval by using/finding nearest vectors

$$(Q, C1) = 0.21$$

$$(Q, C2) = 0.91$$

$$(Q, C3) = 0.18$$

(ii) Reranking :



- In Semantic Search system, Reranking is the two-stage refinement step that takes the initial set of search results & reorders them based on deeper semantic relevance.
- Reranking is needed because the first pass dense retrieval is fast but only approximates relevance. However, it can miss subtle intent, handle synonyms imperfectly or be fooled by high frequency words.
- A second stage is more expensive model that looks at the query & each candidate together, capture contextual meaning, the context & the domain specific meaning, & reorder the results so that the top few are truly the most useful.

15-12-25

⇒ Retrieval Evaluation Metrics :

i) It measures how many of top k results, returned by a system are actually relevant.

$$\text{Precision at } k = \frac{\text{No. of relevant doc in top } k}{k}$$

ii) Average Precision (AP) - Avg of Precision value is computed at each rank where a relevant document appears.

$$AP = \frac{1}{R} \sum_{i=1}^N P(i) \cdot rel(i)$$

iii) Mean AP (MAP) - It is the arithmetic mean of AP scores across a set of queries. It gives a single query independent measure of a ranking system's overall performance.

$$MAP = \frac{1}{Q} \sum_{j=1}^Q AP_j$$

iv) Normalized Discounted Cumulative Gain (nDCG) - It measures how well a rank list of items matches an ideal ordering, taking into account both relevance & position.

$$nDCG = \frac{DCG}{IDCG}$$

where

$$DCG = \sum_{i=1}^n \frac{\text{relevance}_i}{\log_2(i+1)}, IDCG = \sum_{i=1}^n \frac{\text{Grade}_i}{\log_2(i+1)}$$

here, most doc, relevance = 1
relevance = 0

most doc, Grade is either binary (i.e 1/0) or can consider 3pt/5pt rating scale.

Q) An user issues a query & the retrieval system returns 5 documents, ordering from top to bottom by mentioning relevant or not relevant.

Rank	Doc	Relevant?	Relevant Grade
1	A	Rel	1
2	B	Not Rel	0
3	C	Rel	1
4	D	Not Rel	0
5	E	Rel.	1

(sorted)

$$\text{Step 1: Precision at } k=3 \text{ (tot)} = \frac{2}{3}$$

$$\text{Step 2: AP} = \frac{1}{3} (1(\gamma_1) + 2(\gamma_3) + 3(\gamma_5)) = \frac{1}{3} (1 + \frac{1}{3} + \frac{3}{5}) = \frac{84}{45} = 0.75.$$

$$\text{Step 3: MAP} = 0.75$$

$$\text{Step 4: DCG} = \sum_{i=1}^n \frac{\text{relevance}_i}{\log_2(i+1)} = \frac{1}{\log_2(1+1)} + \frac{0}{\log_2(2+1)} + \frac{1}{\log_2(3+1)} + \frac{0}{\log_2(4+1)} + \frac{1}{\log_2(5+1)}$$

$$1 + 0 + 0.5 + 0 + 0.38 = 1.88.$$

$$\text{IDCG} = \sum_{i=1}^n \frac{\text{Grade}_i}{\log_2(i+1)} = \frac{1}{\log_2(1+1)} + \frac{1}{\log_2(2+1)} + \frac{1}{\log_2(3+1)} + \frac{0}{\log_2(4+1)} + \frac{0}{\log_2(5+1)}$$

$$1 + 0.63 + 0.5 + 0 + 0 = 2.13. \quad nDCG = \frac{1.88}{2.13} = 0.88.$$

Q) An IR system produces the following ranking in answer to queries q_1 & q_2 . The underlined documents are the ones relevant to the user.

RG, q, R	1	<u>q₁</u>	<u>q₂</u>
1	1	<u>F</u>	
0	2	<u>L</u>	<u>G</u>
1	3	<u>G</u>	<u>D</u>
0	4	<u>F</u>	<u>E</u>
0	5	<u>D</u>	<u>L</u>
1	6	<u>E</u>	<u>I</u>
1	7	<u>B</u>	<u>H</u>
0	8	<u>H</u>	<u>C</u>
0	9	<u>I</u>	<u>B</u>
10	C	A	

i) Precision at k

$$P_{q_1} = \frac{2}{5}, \quad P_{q_2} = \frac{3}{4}$$

$$\text{ii) AP}_1 = \frac{1}{5} (1(\gamma_1) + 2(\gamma_2) + 3(\gamma_3) + 4(\gamma_4) + 5(\gamma_5)) = 0.65$$

$$\text{AP}_2 = \frac{1}{4} (1(\gamma_1) + 2(\gamma_2) + 3(\gamma_3) + 4(\gamma_4)) = 0.79$$

$$\text{iii) MAP} = \frac{0.65 + 0.79}{2} = 0.72$$

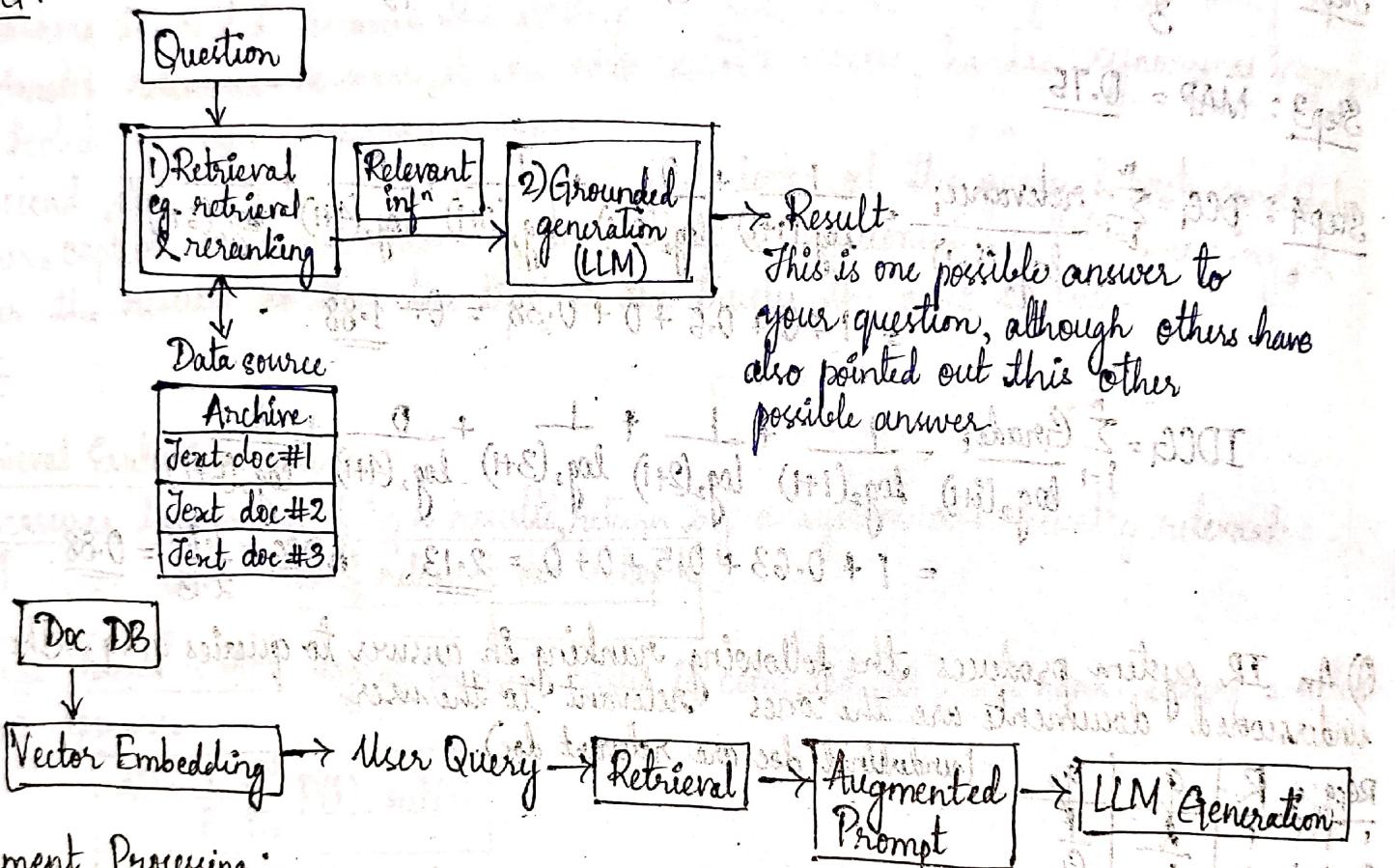
$$\text{iv) DCG}_1 = \frac{1}{\log_2(1+1)} + \frac{0}{\log_2(2+1)} + \frac{1}{\log_2(3+1)} + 0 + 0 + \frac{1}{\log_2(6+1)} + \frac{1}{\log_2(7+1)} + \frac{1}{\log_2(8+1)} + 0 + 0$$

18-12-25 \Rightarrow Retrieval-Augmented Generation (RAG):

It is a framework that enhances large language models by retrieving semantically relevant document chunks from an external knowledge store & injecting them into a prompt, allowing the model to generate fact-grounded & upto date responses without retraining.

19-12-25 Q) Given a vocabulary size of 40,000 & an embedding dimension of 768, calculate the no. of parameters in the embedding layer. Multiplying it we will get no. of parameters.

\Rightarrow RAG:



\Rightarrow Document Processing:

Step 1: Chunking

Step 2: Embedding Vectors

Step 3: User Query

Step 4: Similarity Computation

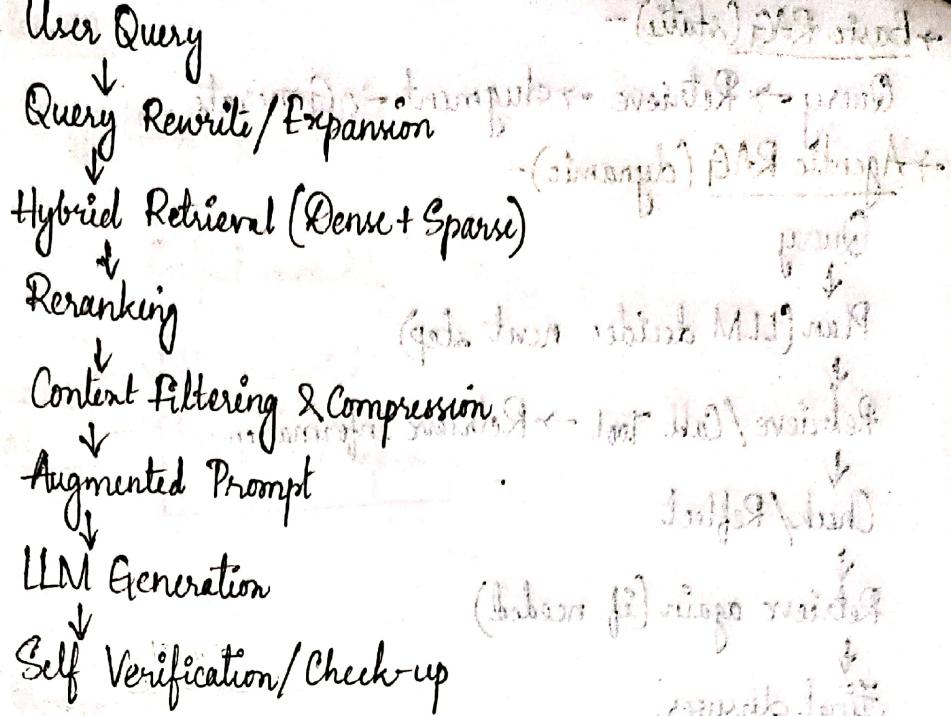
Step 5: Retrieval Result

Step 6: Content integration (or Prompt construction)

Step 7: Generate Answer

Advanced RAG Techniques:

- Before retrieval
- During retrieval
- After retrieval
- During generation



Query Rewriting: LLM rewrites the query into a clearer, retrieval-friendly form.

→ Multi-query retrieval: Extension of query retrieval. The next improvement we can introduce is to extend the query rewriting to be able to search multiple queries if more than one is needed to answer a specific question.

Generate multiple variations of the query & retrieve documents for all.

Q₁: What is dense retrieval?

Q₂: How does dense retrieval work?

Q₃: Dense retrieval

21-12-25 → Multi-hop RAG:

Eg: User Question - "Who are the largest car manufacturers in 2023? Do they each make EVs or not?"

To answer this, the system must first search for
Who are the largest car manufacturers in 2023.

Then after it gets this information (the result being Toyota, Volkswagen, & Hyundai), it should ask follow-up questions:

Step 2, Query 1 - "Toyota car makes EVs."

Query 2 - Volkswagen

Query routing:

It simply does redirection

It is process of classifying a user query & directing it to the most suitable retrieval strategy or data source.

Agentic RAG:

A language model acts as an agent, that means LLM plans to what it will do, then decide what to process, ...

→ Basic RAG (static) -

Query → Retrieve → Augment → Generate

→ Agentic RAG (dynamic) -

Query

↓
Plan (LLM decides next step)

↓
Retrieve / Call Tool → Retrieve information

↓
Check / Reflect

↓
Retrieve again (if needed)

↓
Final Answer

→ RAG Evaluation :

This is to check how far the answer is appropriate to query.

RAG System

└ Retrieval Evaluation (Precision at k, AP, MAP, nDCG)

↳ Basic evaluation strategy

└ Generative Evaluation ()

↳ It evaluates results along four axes :- Fluency : whether

Perceived utility

Citation recall

Citation precision