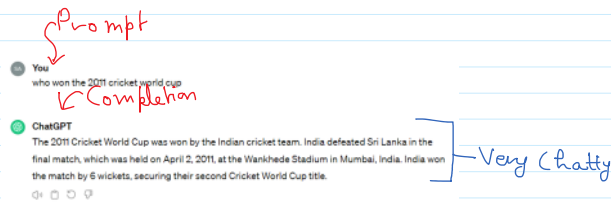


Prompt Engineering

Prompt \Rightarrow The input that you send into a generative model (text, images, videos etc)

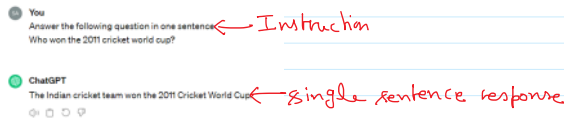
Completion \Rightarrow Output

Example:

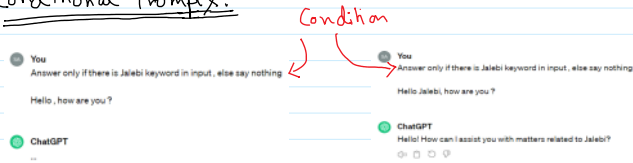


Prompt Structure:-

Instruction \Rightarrow

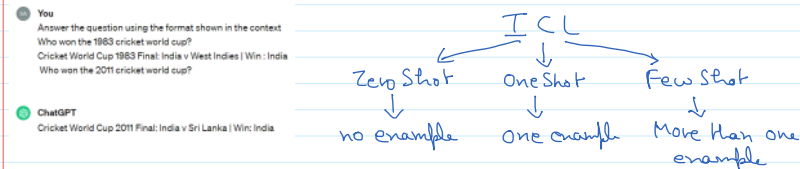


Conditional Prompts:-



In-Context Learning:-

Allows model to perform tasks with a minimal amount of training examples.



Tokens:-



The smallest unit into which text data can be broken down for all AI model to process.

[OpenAI Tokenizer](#)

How GenAI Works?

Foundation Model

Prompt

The students learn from the....

Tokenizer

Content Window

Detokenizer

156	793	420	895	156
the	student	learn	from	the

Token Input IDs

Max Amount of tokens, model consider while generating response

In practical terms, the content window limits how much previous dialogue the model can remember during an interaction.

This includes both prompt provided by user & model's generated text. \rightarrow Short Term Memory

If the interaction exceeds the content window, the model loses access to the earliest part of the conversation.

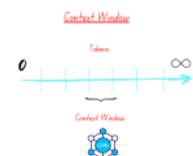
★ Max Token:-

Total no. of tokens Generative AI can generate in response

A very basic mechanism to keep model response short.

V.V.I. \Rightarrow Performance Impact \downarrow

Tokens \uparrow Inference Time \uparrow



If the iteration exceeds the context window the model loses access to the earliest part of the conversation.



It is a barrier to long conversation/complex task

Greedy Vs. Random Sampling:-

Prompt \rightarrow Model \rightarrow Probability Distribution across all tokens in the model's known vocabulary

Choose a single token

Greedy

Random

Choose next token with highest probability

The sky is —

Probability	Token
0.3	blue
0.4	limitless
0.2	clear
:	:

The sky is —

Probability	Token
0.3	blue
0.4	limitless
0.2	clear
:	:

Top P & Top K Random Sampling:-

Top K \rightarrow Choose token randomly from only top k tokens with highest probability

Probability	Token		Probability	Token
0.3	blue		0.4	limitless
0.4	limitless	\rightarrow	0.3	blue
0.2	clear		0.2	clear
:	:		:	:

$k=2$

Top K = 1 \Rightarrow Greedy Sampling

Top P \rightarrow Randomly sampling from the set of tokens whose cumulative probability do not exceed p , starting from highest probability and working down to the lowest.

$p=0.7$

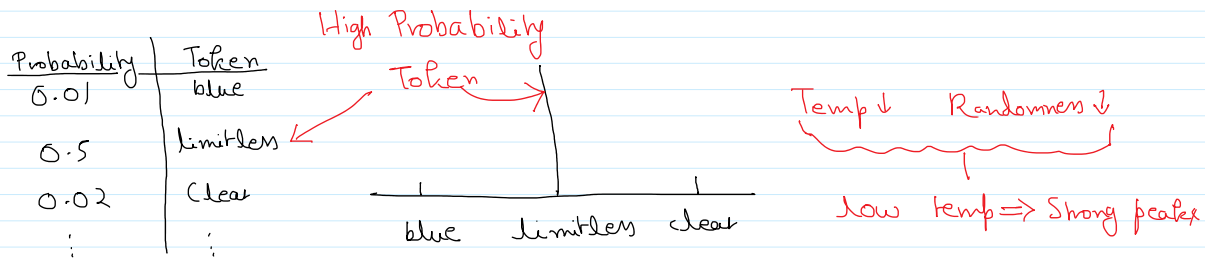
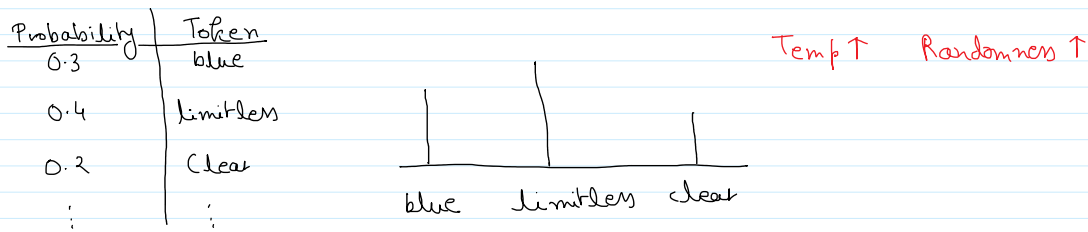
Probability	Token		Probability	Token
0.3	blue		0.4	limitless
0.4	limitless	\rightarrow	0.3	blue
0.2	clear		0.2	clear
:	:		:	:

$p=0.7$

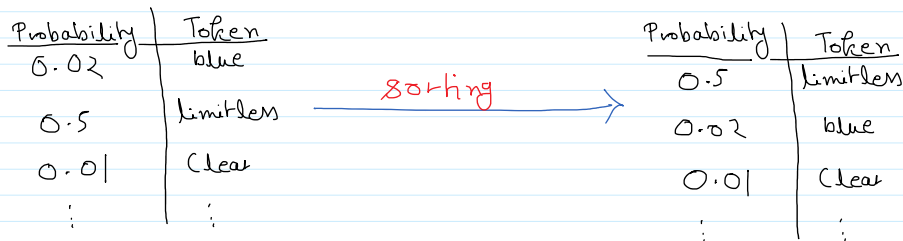
Temperature:-

top-p / top-k \rightarrow affects next token prediction after probability distribution is generated

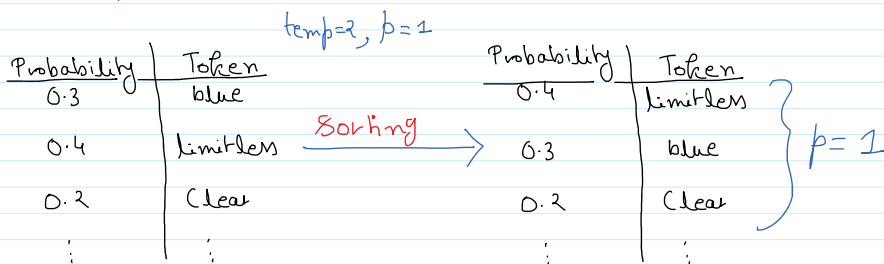
temperature \rightarrow changes next-token probability distribution \rightarrow ultimately affect next token prediction.



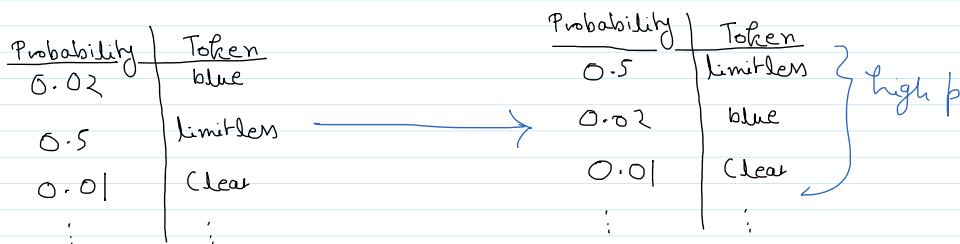
Low Temp, low Top-p \Rightarrow highly focused on narrow range of high probability tokens.



High Temp, High Top-p \Rightarrow very high randomness.



Low Temp, High Top-p ✓



Embedding Vectors :-

Is it a fruit?



Banana
Angle represents distance

Mango \rightarrow

Is it a fruit?	Cos?
1	30

 $\rightarrow [1, 30]$

Banana \rightarrow

Is it a fruit?	Cos?
1	10

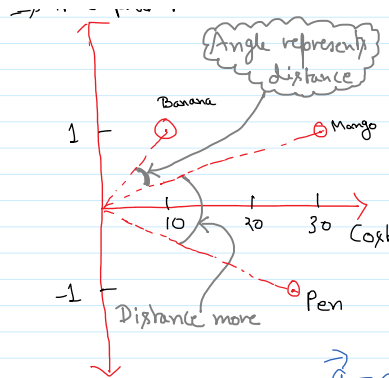
 $\rightarrow [1, 10]$

Pen \rightarrow

Is it a fruit?	Cos?
-1	25

 $\rightarrow [1, 25]$

Embedding are numerical, vectorized representation of any type, including text, images, audio clips, videos etc.



$\vec{a} \cdot \vec{b} = |\vec{a}| |\vec{b}| \cos \theta$ ← Another Equation

$\cos \theta = \frac{\vec{a} \cdot \vec{b}}{|\vec{a}| |\vec{b}|}$

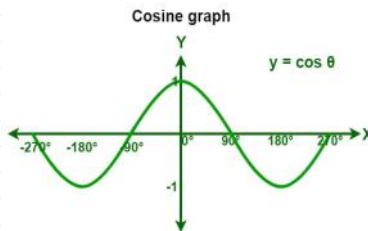
$\vec{a} = a_x \hat{i} + a_y \hat{j} + a_z \hat{k}$

$\vec{b} = b_x \hat{i} + b_y \hat{j} + b_z \hat{k}$

$\vec{a} \cdot \vec{b} = a_x b_x + a_y b_y + a_z b_z$

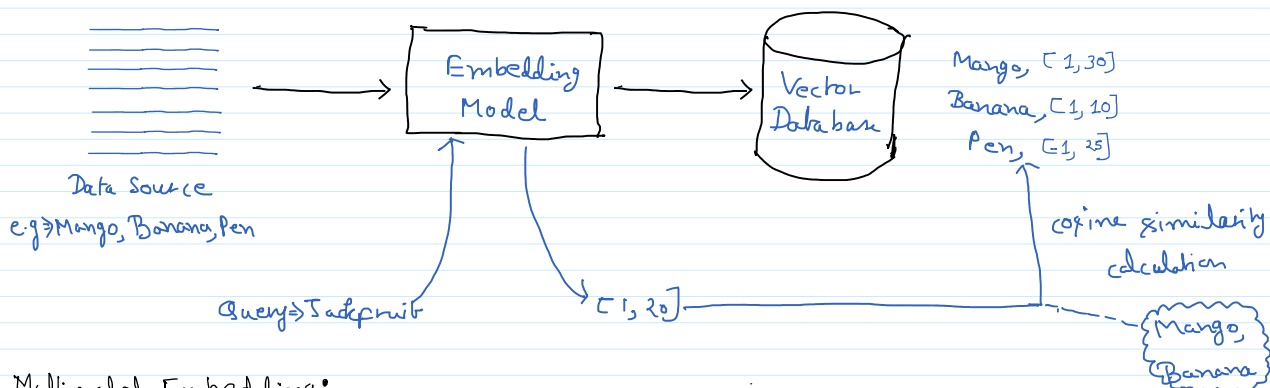
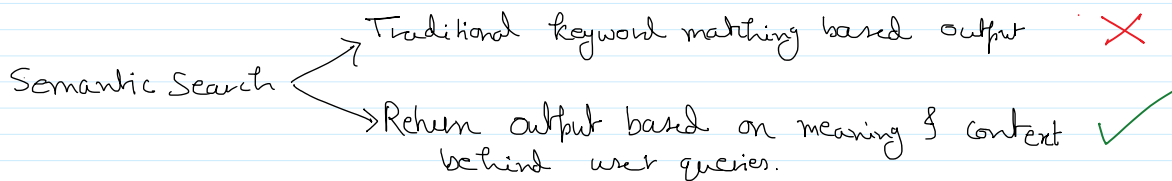
$|\vec{a}| = \sqrt{a_x^2 + a_y^2 + a_z^2}$

$|\vec{b}| = \sqrt{b_x^2 + b_y^2 + b_z^2}$



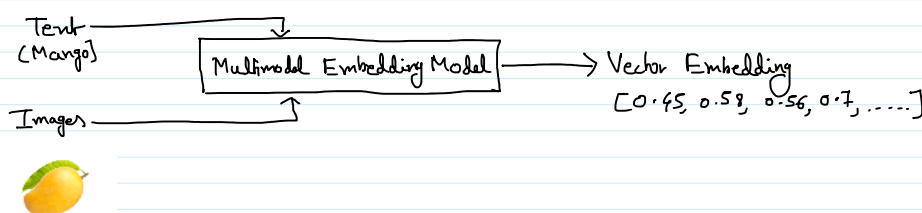
$\cos \theta = \frac{a_x b_x + a_y b_y + a_z b_z}{(\sqrt{a_x^2 + a_y^2 + a_z^2})(\sqrt{b_x^2 + b_y^2 + b_z^2})}$

Semantic Search:-

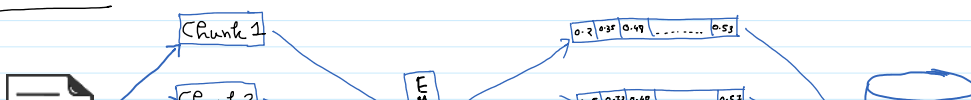


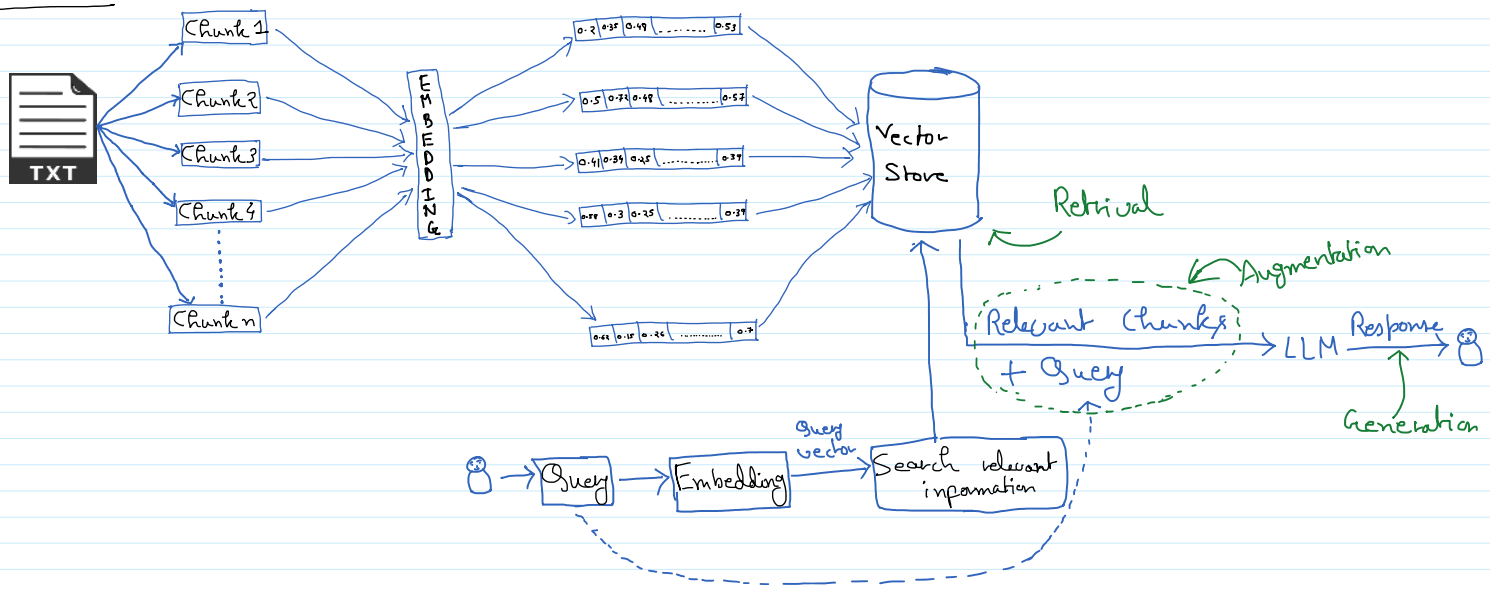
Multi modal Embedding:-

\rightarrow Integration of Text, Images & various other datatypes in single vector space.



RAG:-





RAG \rightarrow Combination of 2 memories

model's own
prior knowledge

A search engine