# Prompt Engineering
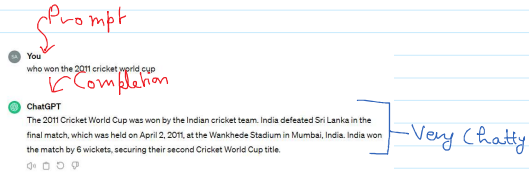
Prompt ⇒ The input that you send into a generative model (text, images, videos etc)
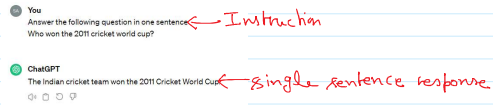
Completion ⇒ Output

Example :
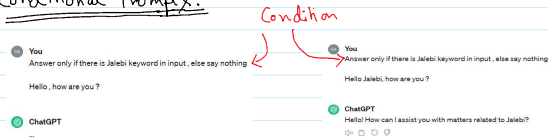
~ Prompt

**You**
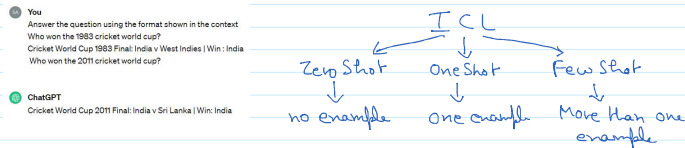who won the 2011 cricket world cup

← Completion

**ChatGPT**
The 2011 Cricket World Cup was won by the Indian cricket team. India defeated Sri Lanka in the
final match, which was held on April 2, 2011, at the Wankhede Stadium in Mumbai, India. India won
the match by 6 wickets, securing their second Cricket World Cup title.

→ Very Chatty

## Prompt Structure :-

Instruction ⇒

**You**
Answer the following question in one sentence    ← Instruction
Who won the 2011 cricket world cup?

**ChatGPT**
The Indian cricket team won the 2011 Cricket World Cup    ← Single sentence response

## Conditional Prompts :-

Condition

**You**
Answer only if there is Jalebi keyword in input, else say nothing    ← Condition

Hello, how are you ?

**ChatGPT**
...

**You**
Answer only if there is Jalebi keyword in input, else say nothing

Hello Jalebi, how are you ?

**ChatGPT**
Hello! How can I assist you with matters related to Jalebi?

## In-Content Learning :-

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

**You**
Answer the question using the format shown in the context
Who won the 1983 cricket world cup?
Cricket World Cup 1983 Final: India v West Indies | Win : India
Who won the 2011 cricket world cup?

**ChatGPT**
Cricket World Cup 2011 Final: India v Sri Lanka | Win: India

```
                    I C L
        ↙           ↓              ↘
   Zero Shot    One Shot      Few Shot
      ↓            ↓              ↓
  no example   one example   More than one
                              example
```
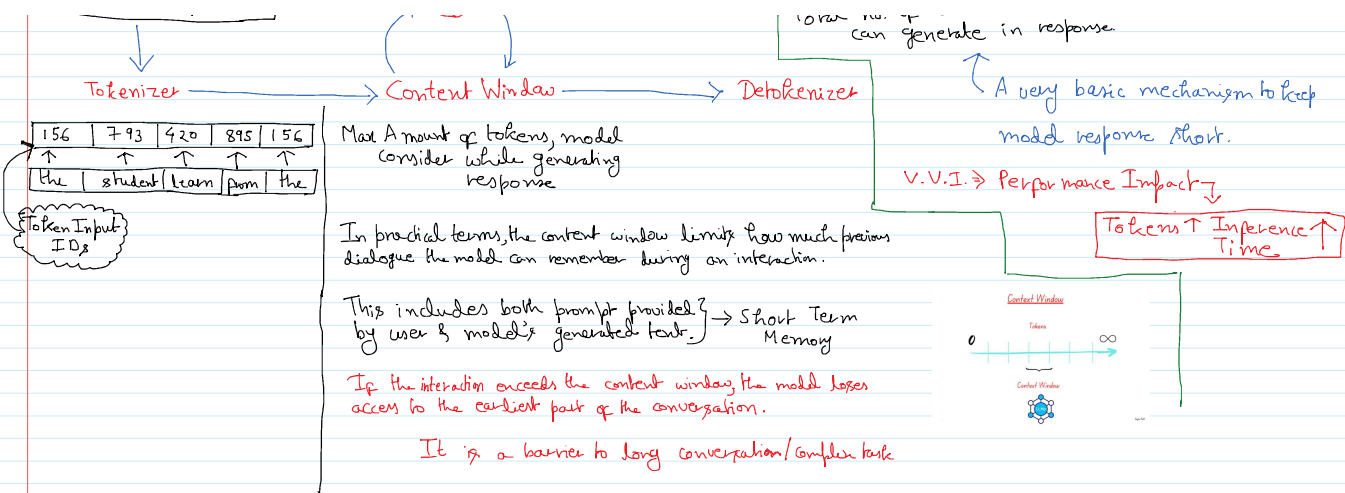
## Tokens :-



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

[OpenAI Tokenizer](OpenAI Tokenizer)

## How Gen AI Works ?

Prompt                    Foundation Model

┌─────────────────────────┐
│ The students learn from the..... │
└─────────────────────────┘
        ↓
    Tokenizer ────────→ Content Window ────────→ Detokenizer

        (Model)

☆ Max Token :-
Total no. of tokens Generative AI
can generate in response

↳ A very basic mechanism to keep

Tokenizer ──────→ Content Window ──────→ Detokenizer

for no. + can generate in response

A very basic mechanism to keep model response short.

| 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. → Short Term Memory

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

It is a barrier to long conversation/complex task

V.V.I. ⇒ Performance Impact

Tokens ↑ Inference ↑ Time

*Context Window*
Tokens
0 ━━━━━━━━━━→ ∞

*Context Window*

Greedy Vs. Random Sampling :-

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

↓

Choose a single token

Greedy                    Random

Choose next token with heighest 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 → Choose token randomly from only top k tokens with heighest probability

| Probability | Token |
|---|---|
| 0.3 | blue |
| 0.4 | limitless |
| 0.2 | Clear |
| ⋮ | ⋮ |

$\longrightarrow$

| Probability | Token |
|---|---|
| 0.4 | limitless |
| 0.3 | blue |
| 0.2 | Clear |
| ⋮ | ⋮ |

$\}$ k = 2

$\boxed{\text{Top } K = 1} \Rightarrow$ Greedy Sampling

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

$$p = 0.7$$

| Probability | Token |
|---|---|
| 0.3 | blue |
| 0.4 | limitless |
| 0.2 | Clear |
| ⋮ | ⋮ |

$\longrightarrow$

| Probability | Token |
|---|---|
| 0.4 | limitless |
| 0.3 | blue |
| 0.2 | Clear |
| ⋮ | ⋮ |

$\}$ p = 0.7

Temperature :-

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

temperature → changes next-token probability distribution → ultimately affect next token prediction.

| Probability | Token |
|---|---|
| 0.3 | blue |
| 0.4 | limitless |
| 0.2 | Clear |
| ⋮ | ⋮ |

Temp ↑   Randomness ↑



blue    limitless    clear

High Probability
Token

| Probability | Token |
|---|---|
| 0.01 | blue |
| 0.5 | limitless |
| 0.02 | Clear |
| ⋮ | ⋮ |

blue    limitless    clear

Temp ↓   Randomness ↓

low temp ⇒ Strong peaks

Low Temp, Low Top-p ⇒ highly focused on narrow range of high probability token.

| Probability | Token |
|---|---|
| 0.02 | blue |
| 0.5 | limitless |
| 0.01 | Clear |
| ⋮ | ⋮ |

→ *sorting* →

| Probability | Token |
|---|---|
| 0.5 | limitless |
| 0.02 | blue |
| 0.01 | Clear |
| ⋮ | ⋮ |

High Temp, High Top-p ⇒ very high randomness.

temp=2, p=1

| Probability | Token |
|---|---|
| 0.3 | blue |
| 0.4 | limitless |
| 0.2 | Clear |
| ⋮ | ⋮ |

→ *sorting* →

| Probability | Token |
|---|---|
| 0.4 | limitless |
| 0.3 | blue |
| 0.2 | Clear |
| ⋮ | ⋮ |

} p = 1

Low Temp, High Top-p ✓

| Probability | Token |
|---|---|
| 0.02 | blue |
| 0.5 | limitless |
| 0.01 | Clear |
| ⋮ | ⋮ |

→

| Probability | Token |
|---|---|
| 0.5 | limitless |
| 0.02 | blue |
| 0.01 | Clear |
| ⋮ | ⋮ |

} high p

## Embedding Vectors :-

Mango → 

| Is it a fruit? | Cost? |
|---|---|
| 1 | 30 |

→ [1, 30]

Banana →

| Is it a fruit? | Cost? |
|---|---|
| 1 | 10 |

→ [1, 10]

Pen →

| Is it a fruit? | Cost? |
|---|---|
| -1 | 25 |

→ [-1, 25]

Is it a fruit?

Angle represents distance

Banana    Mango
1 —

        10   20   30  Cost

-1 —         Pen
Distance more

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

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

$\vec{a} = a_x \hat{i} + a_y \hat{j} + a_z \hat{k}$
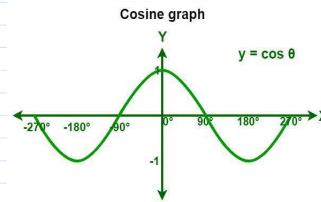
$\vec{b} = b_x \hat{i} + b_y \hat{j} + b_z \hat{k}$

Cosine graph

Pen $\longrightarrow$ | Is it a fruit? | Cost? | $\longrightarrow$ [1, 25]
                      |     −1         |  25   |

$\vec{a} = a_x \hat{i} + a_y \hat{j} + a_z \hat{k}$
$\vec{B} = b_x \hat{i} + b_y \hat{j} + b_z \hat{k}$

**Cosine graph**

$y = \cos \theta$

-270°  -180°  -90°   0°   90°   180°   270°   X

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

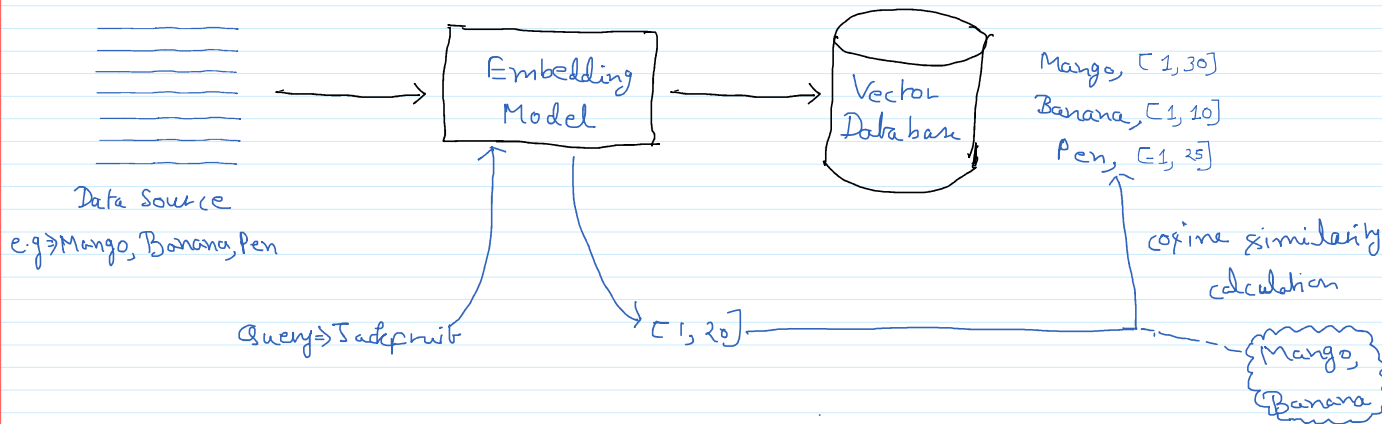$\vec{a} \cdot \vec{b} = a_x b_x + a_y b_y + a_z b_z$

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

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

$$\cos\theta = \frac{a_x b_x + a_y b_y + a_z b_z}{\left(\sqrt{a_x^2 + a_y^2 + a_z^2}\right)\left(\sqrt{b_x^2 + b_y^2 + b_z^2}\right)}$$
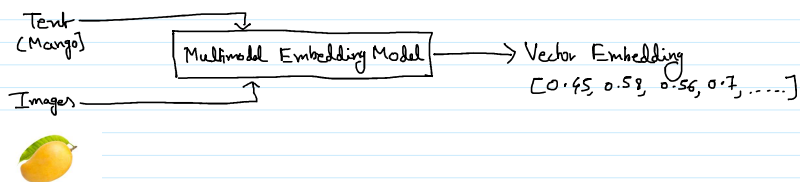
## Semantic Search:-

Semantic Search $\Bigg\{$
→ Traditional keyword matching based output ✗
→ Return output based on meaning & context behind user queries. ✓



Data Source
e.g ⇒ Mango, Banana, Pen

Embedding Model

Vector Database

Mango, [1, 30]
Banana, [1, 10]
Pen, [1, 25]

cosine similarity calculation

Query ⇒ Jackfruit

[1, 20]

{Mango, Banana}

## Multimodal Embedding:

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

Text (Mango) ———→ Multimodal Embedding Model ———→ Vector Embedding
Images ————————↗                                  [0.45, 0.58, 0.56, 0.7, ......]

# RAG :-



TXT

Chunk 1
Chunk 2
Chunk 3
Chunk 4
Chunk n

EMBEDDING

| 0.2 | 0.35 | 0.49 | ......... | 0.53 |
| 0.5 | 0.72 | 0.48 | ......... | 0.57 |
| 0.41 | 0.34 | 0.25 | ......... | 0.34 |
| 0.58 | 0.3 | 0.25 | ......... | 0.34 |
| 0.62 | 0.15 | 0.26 | ......... | 0.7 |

Vector Store

Retrieval

Augmentation

Relevant Chunks + Query

LLM → Response

Generation

Query → Embedding → Query vector → Search relevant information

RAG → Combination of 2 memories

model's own prior knowledge

A search engine

# index-demo Info

## Getting started                                                         Hide

**Step 1. Create an index**

An index is the place where you add your data sources to make them searchable in Amazon Kendra.

⊘ Created!

**Step 2. Add data sources**

Add and sync your data from S3, SharePoint, and other databases to your index.

**Add data sources**

**Step 3. Test & Deploy**

After syncing your data, visit the Search console page to test search and to deploy Amazon Kendra in your search application.
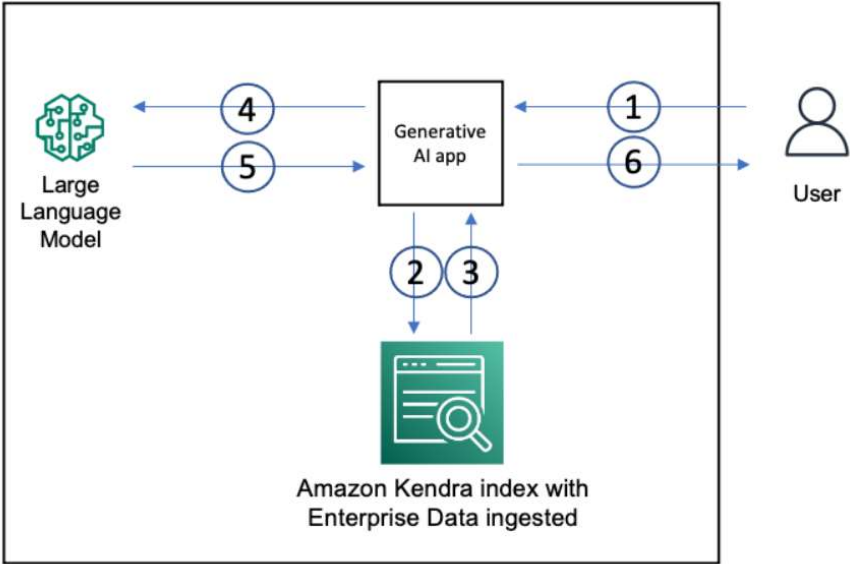
**Search console**

**Add FAQs** - *optional*

Upload FAQs to seamlessly provide curated answers to commonly asked questions in search results.

**Add FAQs**

Large Language Model

④
⑤

Generative AI app

①
⑥

User

②③

Amazon Kendra index with Enterprise Data ingested

# Rag Implementation:-



Knowledge Base

Retrieval

Augmentation

generation