

Learn and Launch

- Generative AI Application -



1

Understand Generative AI

2

Word Embedding

3

Prompt Engineering for LLM and Image Model

4

Build a Generative AI Application

AGENDA

Lecture

Lab



Ayyanar Jeyakrishnan



bit.ly/awsaj



ABOUT ME

I started my career as humble Hardware and Networking engineer in 2005 in HCL Infosystem.

As a Multi Cloud Architect, I have around 18+ years of IT Experience, 50+ Industry Certification, I am deeply passionate about designing and implementing scalable and efficient cloud solutions also Machine Learning and Data Platform for enterprises.

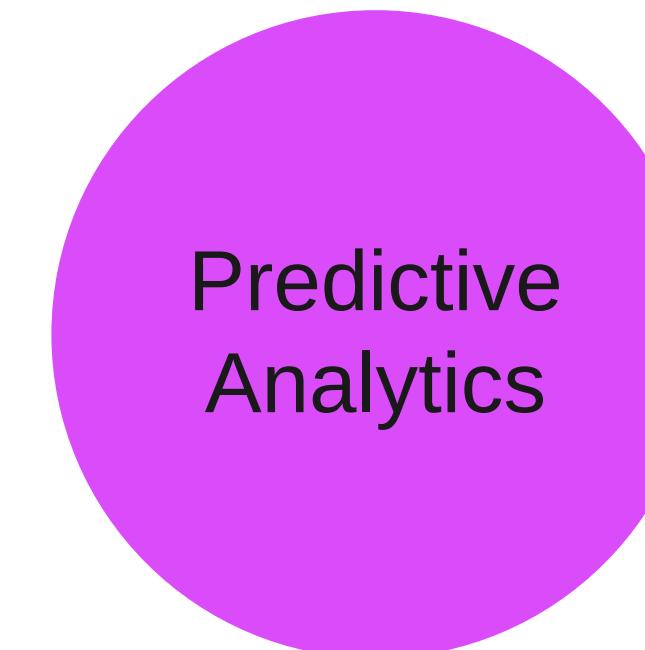
Currently, I am working as VP - Principal Engineer in Wells Fargo.

I always love to learn - Unlearn - Relearn, Motivated to Share knowledge with peers, community and learn from them.

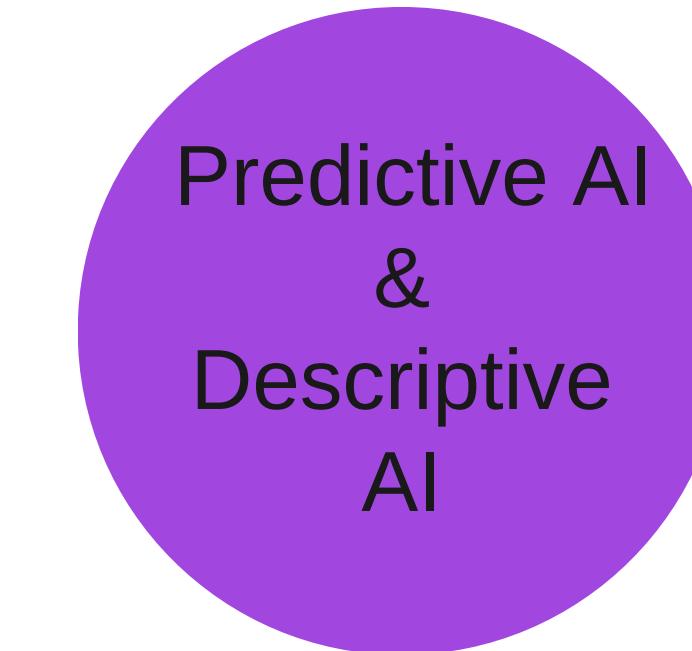
Path to Gen AI



Charts and
Graphs - Using
BI



Linear
Regression
Logistic
Regression
Clustering



Word
Embeddings
Word2Vec

RNN, LSTM,
GRU
CNN, YOLO

Path to Gen AI



Charts and
Graphs - Using
BI



Linear
Regression
Logistic
Regression
Clustering



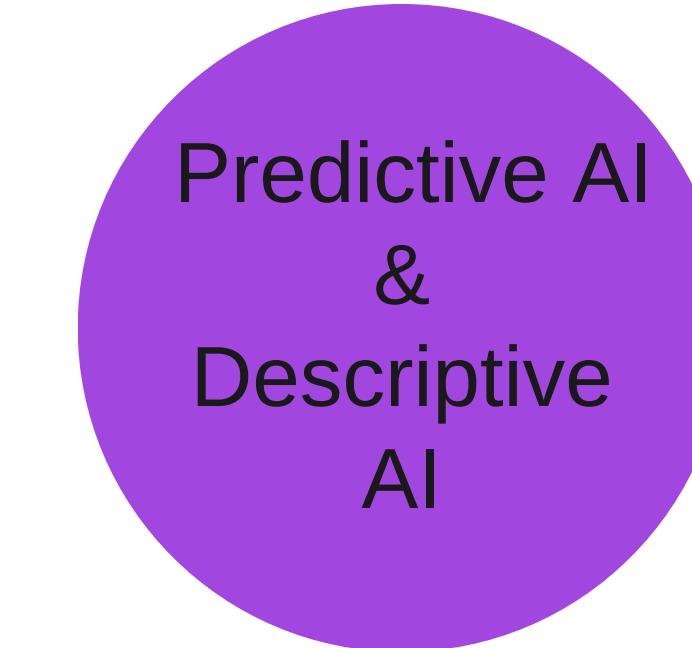
Path to Gen AI



Charts and
Graphs - Using
BI



Linear
Regression
Logistic
Regression
Clustering



Word
Embeddings
Word2Vec

RNN, LSTM,
GRU
CNN, YOLO



Path to Gen AI

Descriptive Analytics

Charts and
Graphs - Using
BI



Predictive
Analytics

Linear
Regression
Logistic
Regression
Clustering



Predictive AI
&
Descriptive
AI



Word
Embeddings

Word2Vec

RNN, LSTM,
GRU
CNN, YOLO



Generative
AI

GAN (Gen and Des)

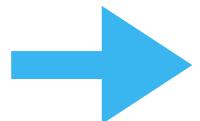
Transformer Models

Falcon, MPT, Dolly, GPT-3.5, Bloom, FLAN-T5, BART and BERT, LLaMA-2

Stable Diffusion

Let us Understand Basics on How Computer Learn

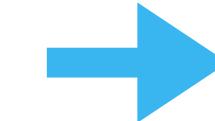
Dog
Cat
Fruit
Banana



3
5
18
19

Word2Vec

[-0.29418945, 0.3930664, 0.76464844, -0.722167...
[0.72265625, -0.30395508, 0.54296875, 0.904296...
[-1.1162109, -0.6743164, -0.111694336, -0.2362...
[0.87353516, -0.5004883, 0.734375, -1.1494141,...
[-2.2753906, -0.21533203, -1.0986328, 0.472656...]



cosine_similarity

Semantic Search

Dog		
Cat		
	Fruit	
		Banana

Tokenization

Embedding

KMeans

Clustering

Dog	
Cat	
	Fruit
	Banana

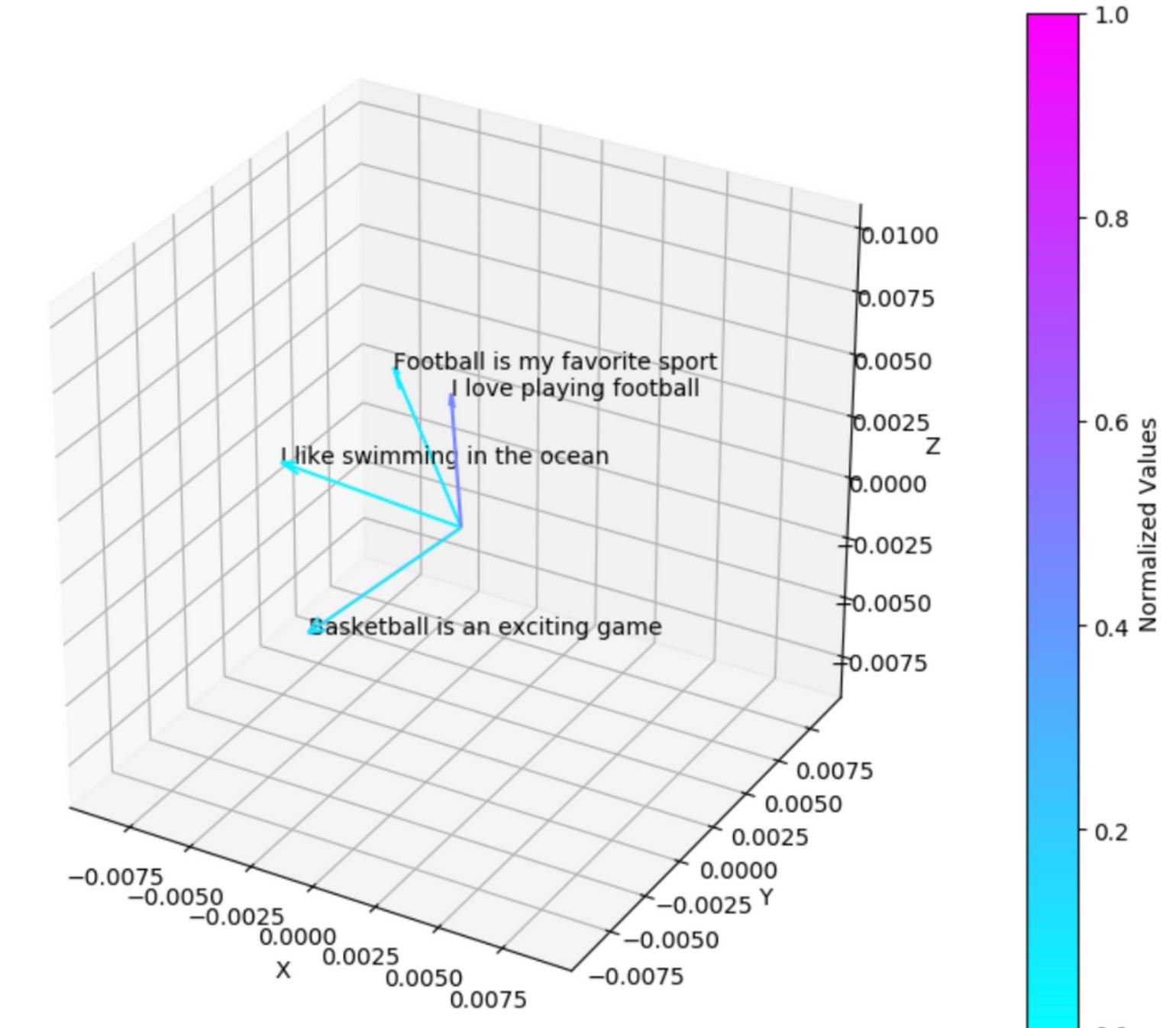
Lab 1 - Word Embedding

```
corpus = [  
    "I love playing football",  
    "Football is my favorite sport",  
    "I enjoy watching football matches",  
    "Soccer is popular worldwide"  
]
```

Vocabulary size: 15

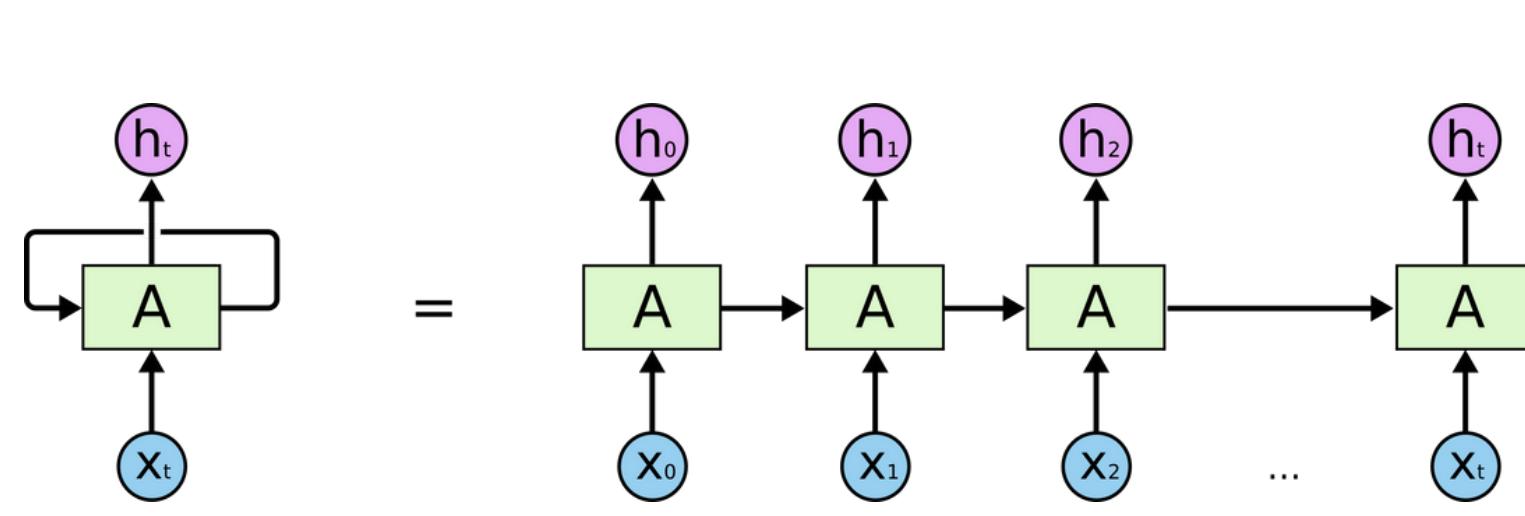
Word Vocabulary:
is
football
I
worldwide
popular
Soccer
matches
watching
enjoy
sport
favorite
my
Football
playing
love

Similar words to 'football':
love: 0.16072481870651245
playing: 0.13725270330905914
worldwide: 0.06797593832015991
enjoy: 0.03363519534468651
Soccer: 0.009391184896230698
my: 0.008315940387547016
popular: 0.004503006115555763
Football: -0.003644464071840048
is: -0.010839183814823627
I: -0.023671666160225868

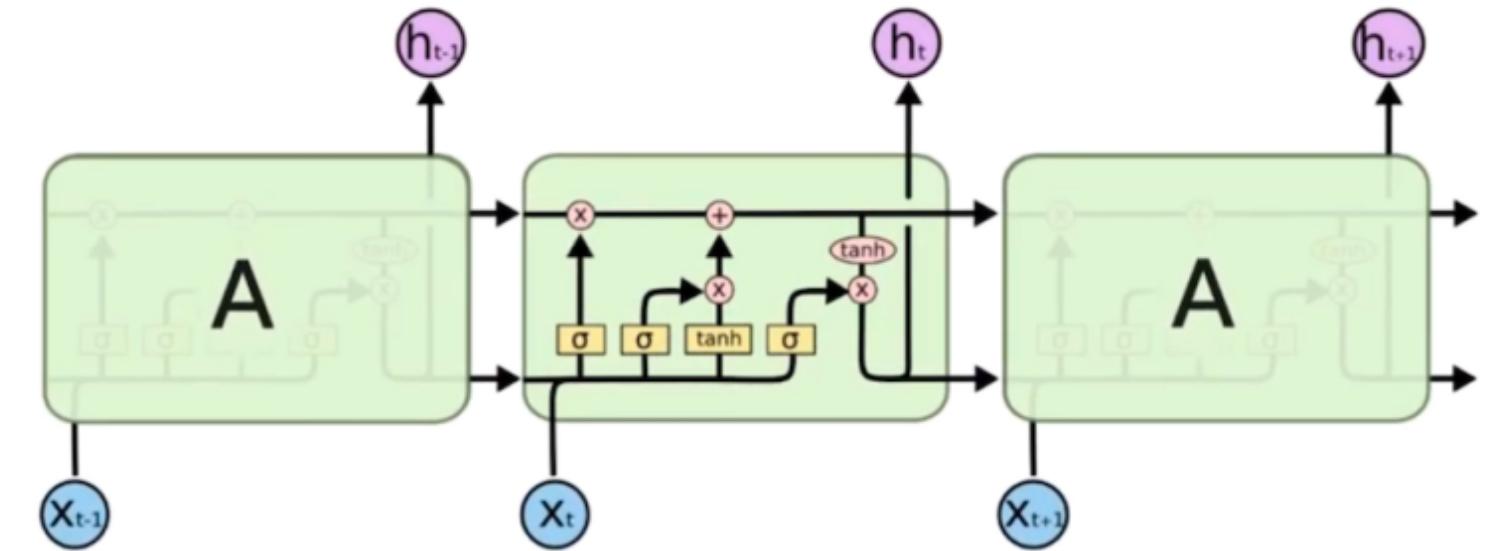


Feed Forward Network - Sequential Models

RNN



LSTM Networks



The repeating module in an LSTM contains four interacting layers.

the challenge

- It's slow to train.
- Long sequences lead to vanishing gradient or the problem of long-term dependencies. In simple terms, its memory is not that strong when it comes to remembering old connections.
- Vanishing gradient
- Slow training

Transformers - Attention Is All You Need - 2017

The Transformer model architecture Advantage.

- “Multi Headed Attention” enabled models to scale the understanding of relationships between words and help to give context.
- Efficiently use Parallel computing.
- Transfer Learning and Adaptability

Scalability

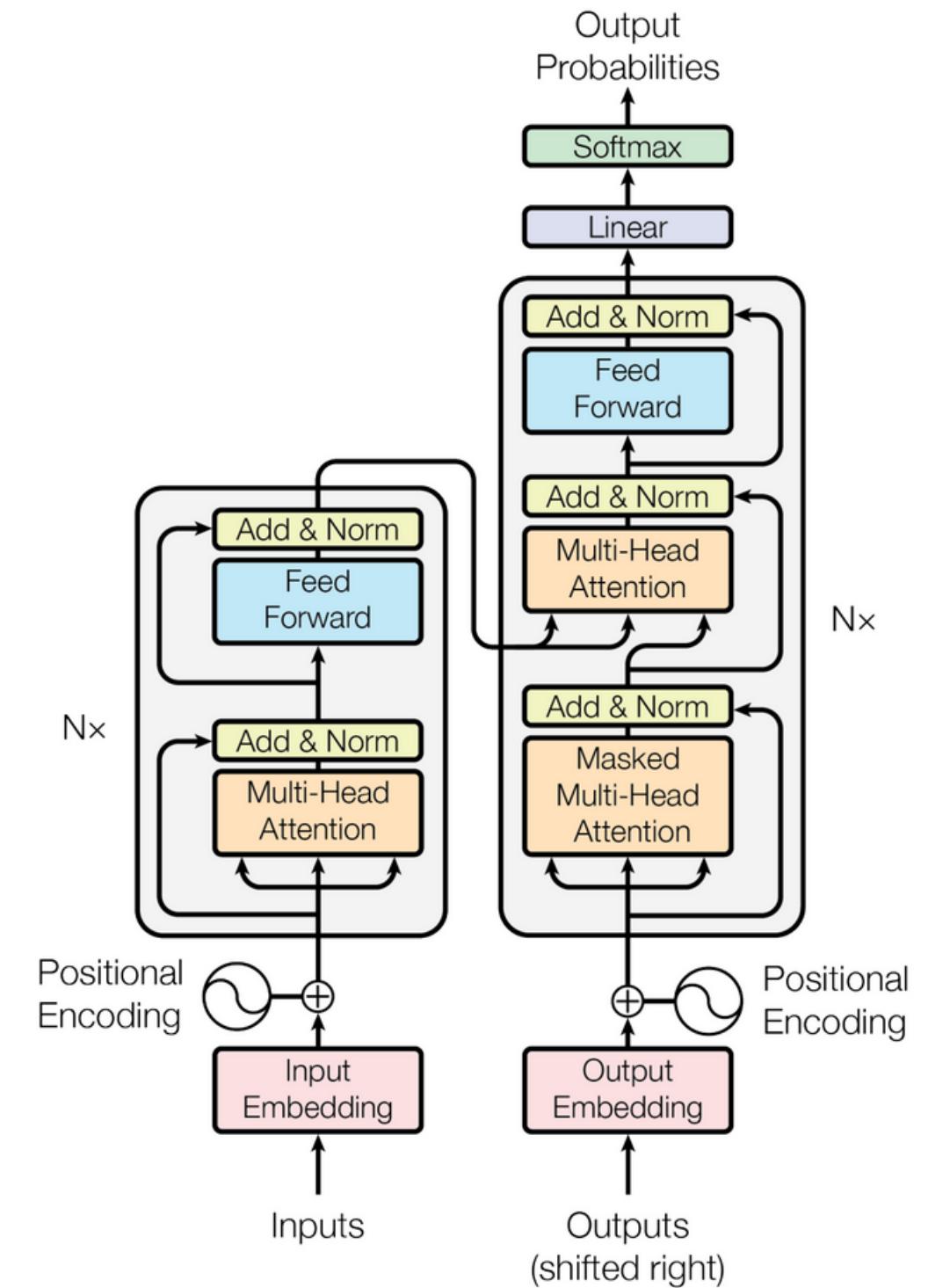
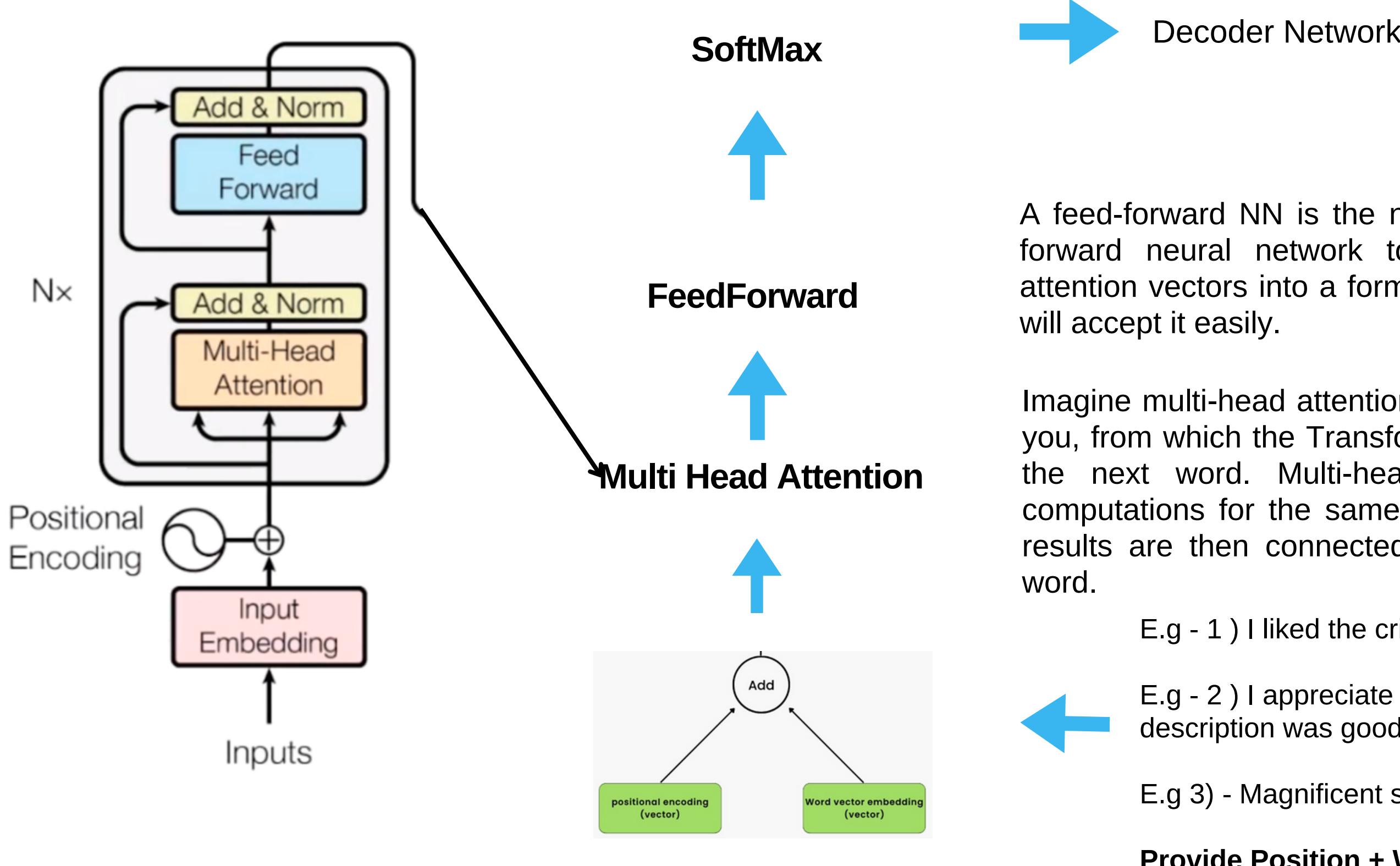


Figure 1: The Transformer - model architecture.

Transformer - Encoder - Intuitive Understanding

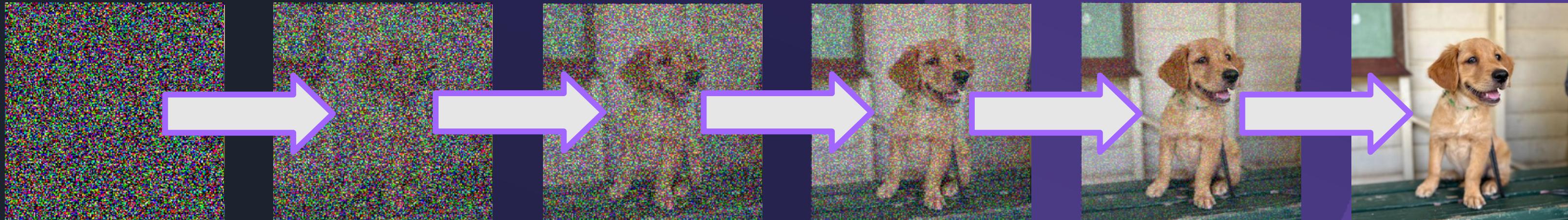


Stable Diffusion

Diffusion models.



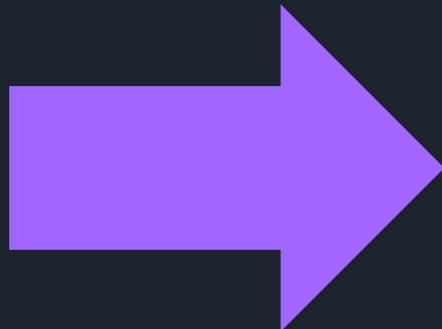
Adds noise and learns how to work backwards to the original image.



Trained model works from random noise to generated an image.

Stable Diffusion Examples

Diffusion models.



the puppy”



“Marigold the puppy on a space station”

Foundation models available on SageMaker JumpStart for self-managed access

Publicly available

stability.ai



Models

Text2Image
Upscaling

Tasks

Generate
photo-realistic
images from
text input

Improve quality
of generated
images

Features

Fine-tuning on
SD 2.1 model

Models

AlexaTM
20B

Tasks

Machine
translation

Question
answering

Summarization

Annotation

Data generation

Models

Flan T-5 models
(8 variants)

DistilGPT2, GPT2

Bloom models
(3 variants)

Tasks

Machine
translation

Question
answering

Summarization

Annotation

Data generation

Proprietary models

co:here



AI21 labs

Models

Cohere
generate-med

Tasks

Text generation

Information
extraction

Question
answering

Summarization

Models

Lyra-Fr
10B

Tasks

Text generation

Keyword
extraction

Information
extraction

Question
answering

Summarization

Sentiment
analysis

Classification

?

Models

Jurassic-1
Grande 17B

Tasks

Text generation

Long-form
generation

Summarization

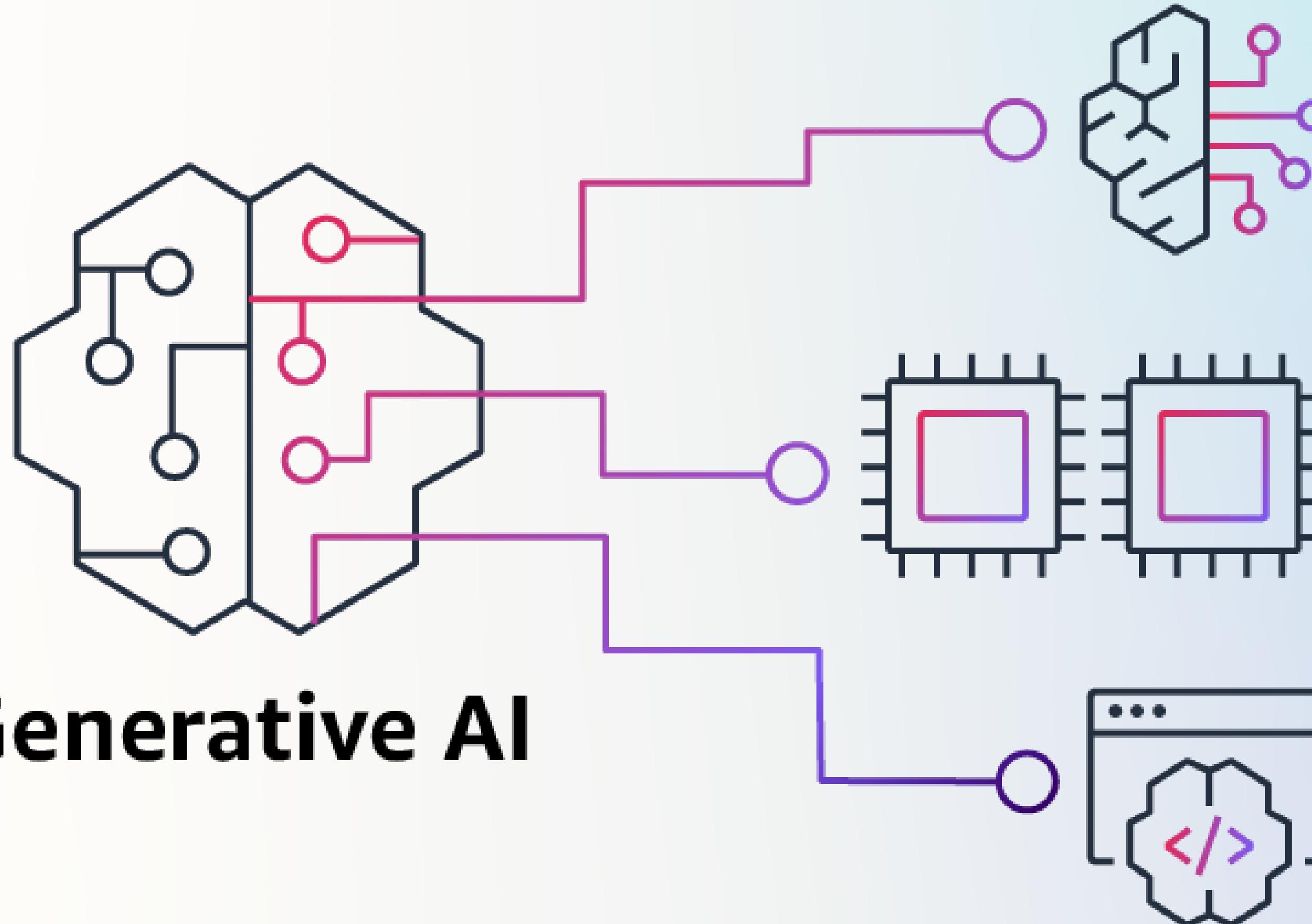
Paraphrasing
Chat

Information
extraction

Question
answering

Classification

Generative AI



Amazon Bedrock

**Amazon EC2 Trn1n
and Amazon EC2 Inf2**

Amazon CodeWhisperer

Connect with Large Language Model Easily

- Prompting -

The screenshot shows the Sagemaker Foundation Model - Playground Console interface. On the left, the **Prompt** section contains a text input field with the placeholder "Be yourself; everyone else is already taken. Rewrite the sentence above:" and a "Generate text" button. Below it, the **Output** section displays the generated text: "A higher temperature value (e.g., 1.5) leads to more diverse and creative text, while a lower value (e.g., 0.5) results in more focused and deterministic text." On the right, the **Example** section includes settings for "Max length" (200), "Temperature" (0.7), "Top P" (1), and "Stop Sequence" (##). Four purple circles with numbers 1 through 4 are overlaid on the interface:

- 1**: Points to the "Context Window" label in the Prompt section.
- 2**: Points to the "Maximum Output Token" label in the Example section.
- 3**: Points to the "Temperature" slider in the Example section.
- 4**: Points to the "Top P" slider in the Example section.

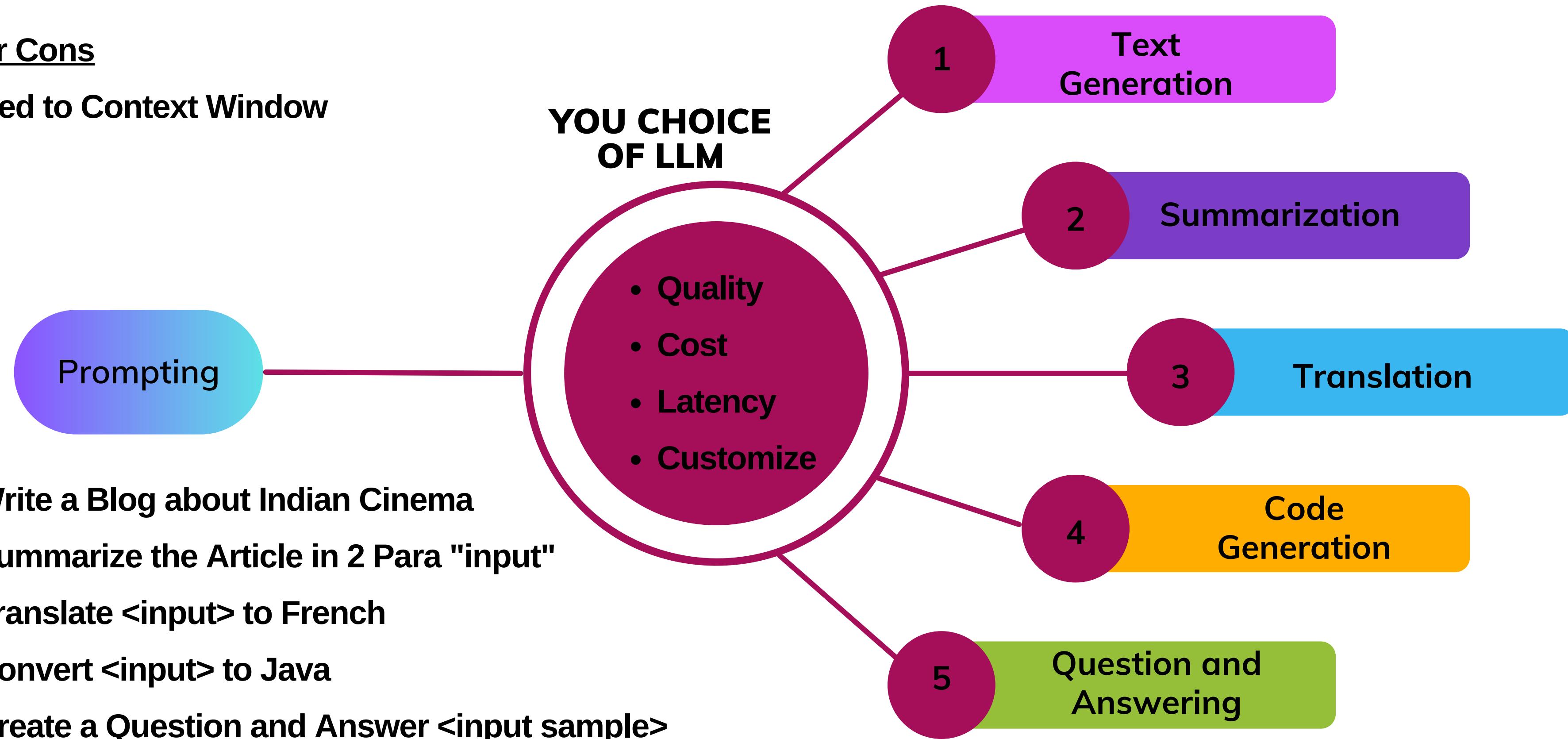
- If the generated text is too narrow in scope and lacks diversity, consider increasing the probability threshold (p).
- If the generated text is too diverse and includes irrelevant words, consider decreasing the probability threshold (p).

Sagemaker Foundation Model - Playground Console

LLM - Prompting

Major Cons

Limited to Context Window



LLM - Prompting

Zero Shot - Task described, but demonstrations not given



Few-shot: Task described and random demonstrations provided



Chain of Thoughts Prompting

Standard prompting

Input:

Q: Roger has 5 tennis balls. He buys 2 more cans of tennis balls. Each can has 3 tennis balls. How many tennis balls does he have now?

A: The answer is 11.

...

Q: The cafeteria had 23 apples. If they used 20 to make lunch and bought 6 more, how many apples do they have?

A:

Model output:

The answer is 50. 

Input:

Q: Roger has 5 tennis balls. He buys 2 more cans of tennis balls. Each can has 3 tennis balls. How many tennis balls does he have now?

A: Roger started with 5 balls. 2 cans of 3 tennis balls each is 6 tennis balls. $5 + 6 = 11$. The answer is 11.

...

Q: The cafeteria had 23 apples. If they used 20 to make lunch and bought 6 more, how many apples do they have?

A:

Model output:

The cafeteria had 23 apples originally. They used 20 to make lunch. So they had $23 - 20 = 3$. They bought 6 more apples, so they have $3 + 6 = 9$. The answer is 9. 

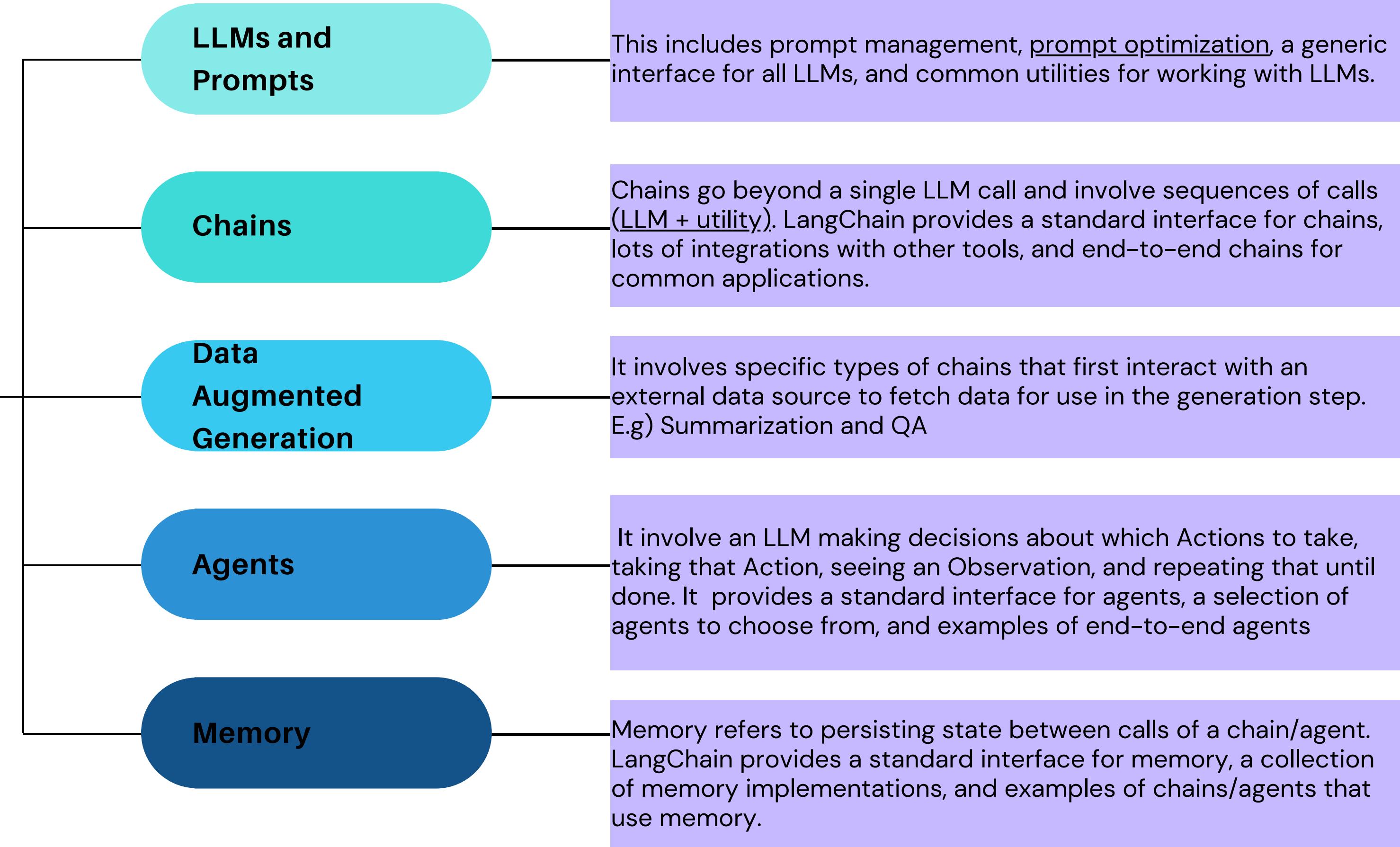
Connect with Large Language via API

- Persona - Developers - Prompting

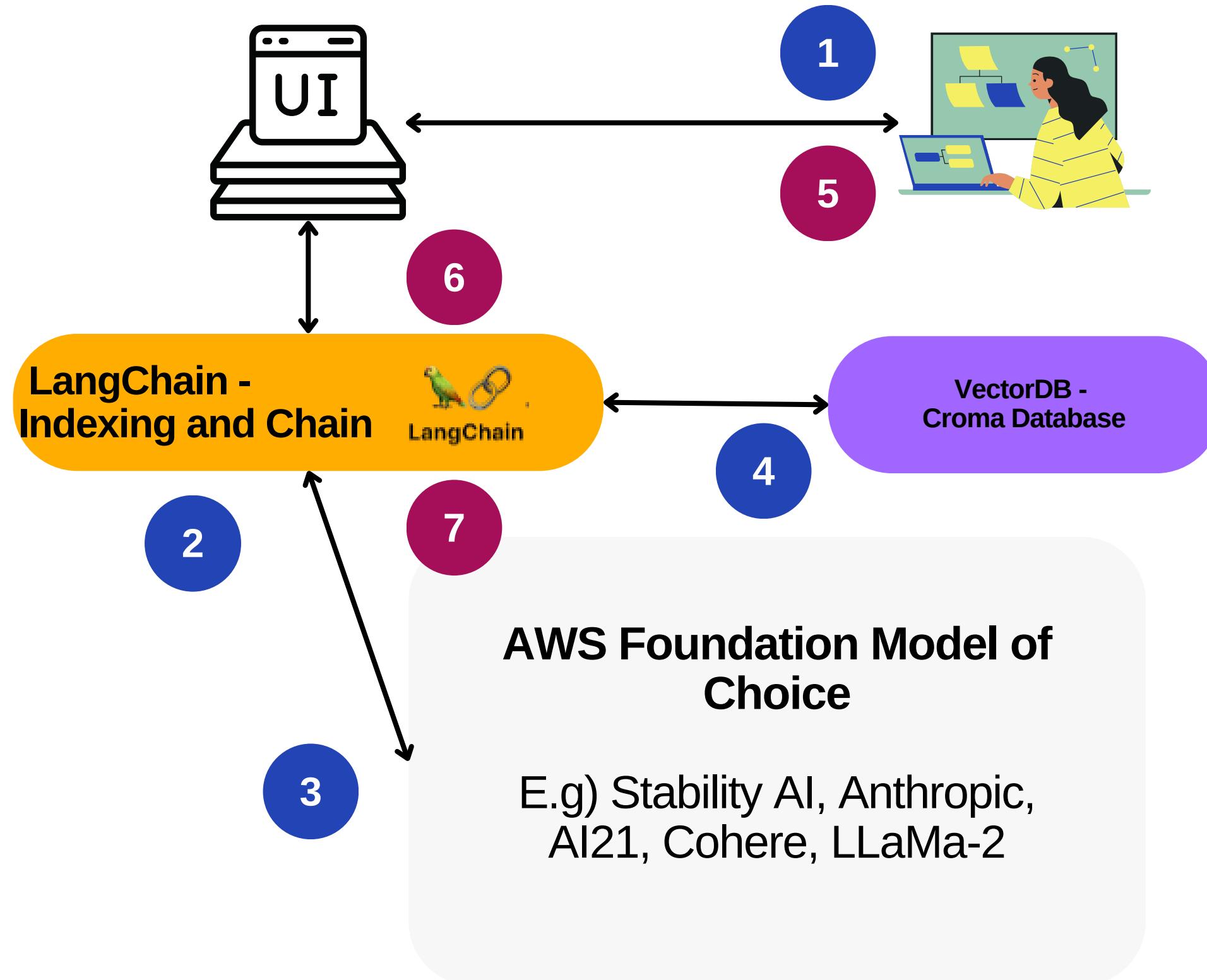
```
# Generate the content
import cohere
co = cohere.Client(api_key) # This is your trial API key
response = co.generate(
    model='command',
    prompt='Write a LinkedIn post about starting a career in tech:',
    max_tokens=300,
    temperature=0.9,
    k=0,
    stop_sequences=[],
    return_likelihoods='NONE')
print('Prediction: {}'.format(response.generations[0].text))
```

Sagemaker Foundation Models - Connect via API

LLM - Langchain

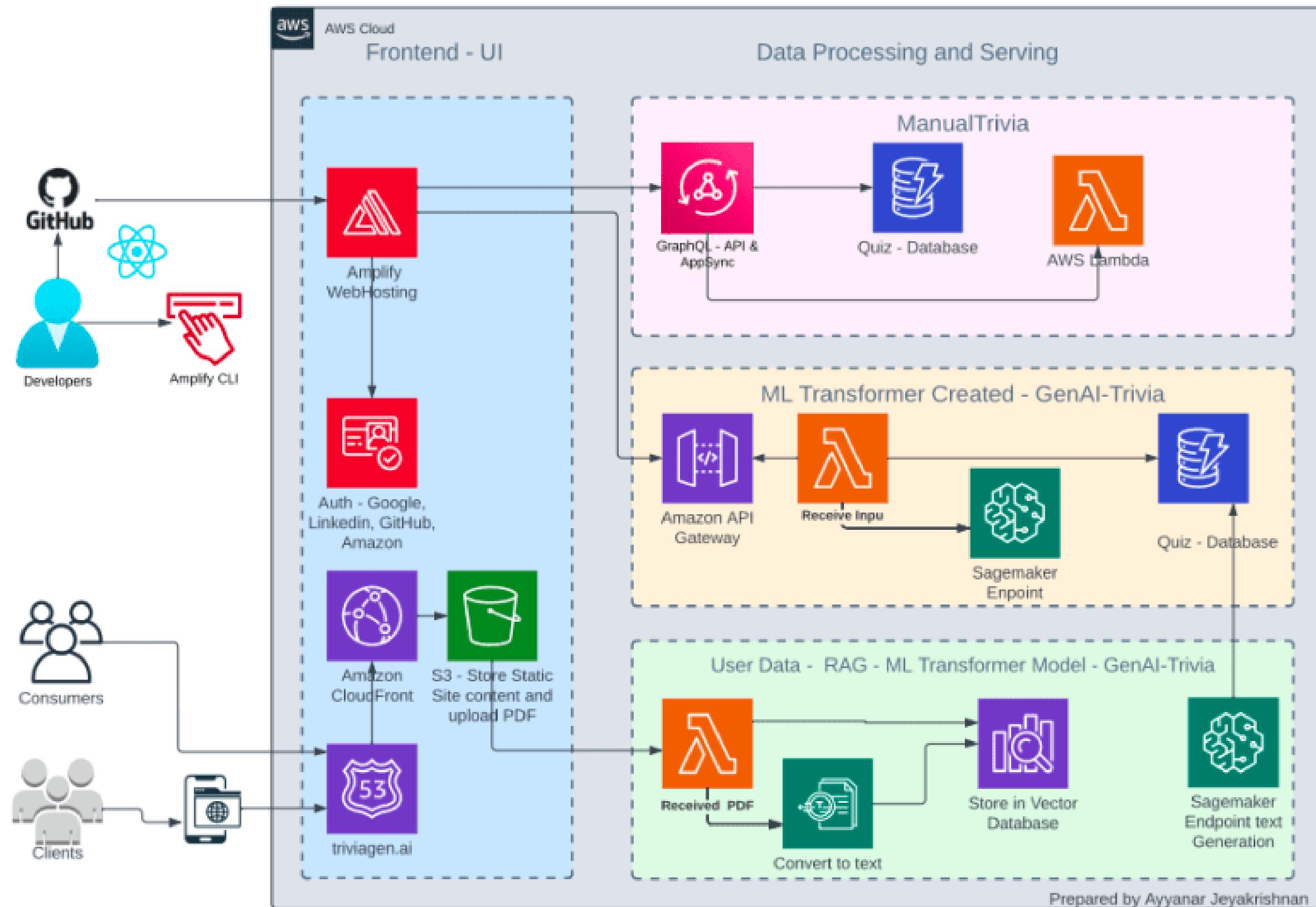


Small Application - Student Study App - P



1. Student upload PDF file via UI
2. Application built on Langchain read the PDF file as chunk and
3. Use Embed Model from Cohere to covert word to Vector
4. Upload Vector to Vector DB - Croma Database
5. Student Provide the Prompt from UI
6. Langchain use the Prompt Template and read the Input and forward the request to LLM and get the context / Semantic search from Vector DB
7. Use the Cohere API and return the result

Sample Usecase - Fun Project



<https://dev.to/jayyanar/triviagenai-m98>



bit.ly/feedback-awscommunity

Thank you!

Ayyanar Jeyakrishnan

QUESTIONS

WELCOME