

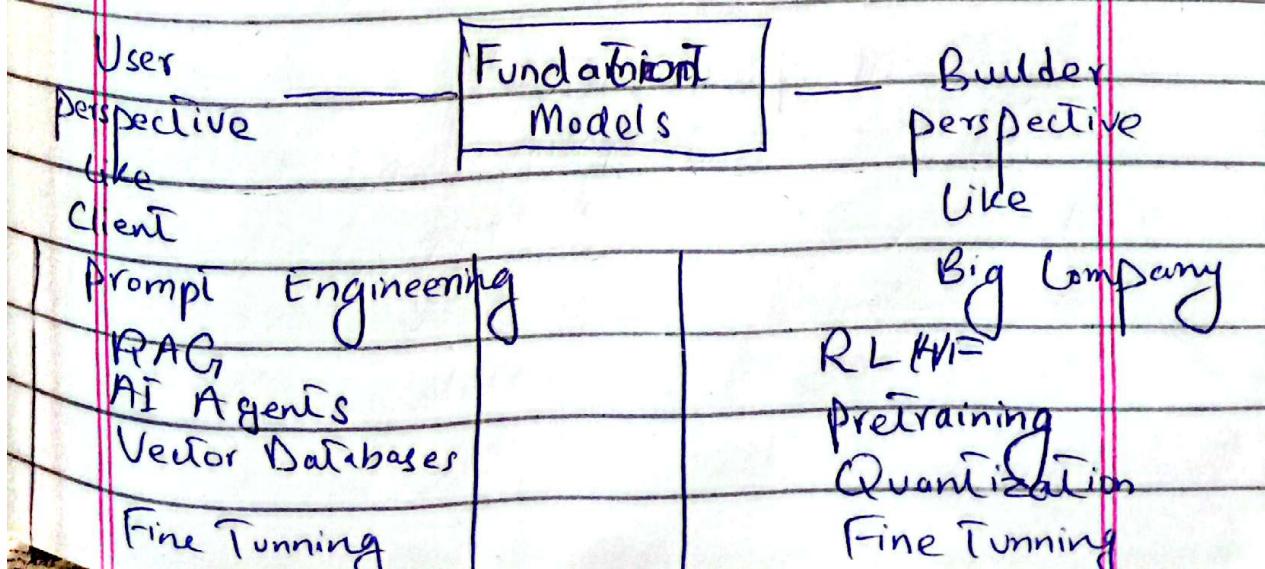
Generative AI

- What is GenAI?
 - type of AI that creates new content, text, images, code

Generative AI

[mimicking human creativity]

 - GenAI Impact Areas:-
 - i) Customer Support
 - ii) Content Creation (Human creativity)
 - iii) Education (Change) ← Teacher
 - iv) Software Development - Is GenAI successful?
 - Does it solve real world problems?
- Yes
- Is it useful on a daily basis?
 - Is it creating new jobs. (Yes) in AI
 - How To learn in (everyday advancement)
 - ~~Using~~ Using foundation models



- The Builder's Perspective:-
 - i) Transformer Architecture.
 - ii) Transformer Types:-
 - Encoder Only (BERT)
 - Decoder Only (GPT)
 - Encoder and Decoder based (T5)
 - iii) PreTraining
 - Training Objectives
 - Tokenization strategies
 - Handling Challenges
 - iv) Optimization
 - Training optimization
 - Model compression
 - Optimizing Inference
 - v) Fine Tuning:
 - Task Specific Tuning
 - Instruction Tuning
 - Continual Pretraining
 - (vi) Evaluation
 - (vii) Deployment

The User's Perspective

(How To build application using LLMs)

i) Building Basic LLM Apps

Open Source vs Closed Source LLMs

fetching api using Hugging Face

Using LLM APIs

LangChain

Hugging Face

Ollama

ii) Prompt Engineering (How To use and generate)

RAG

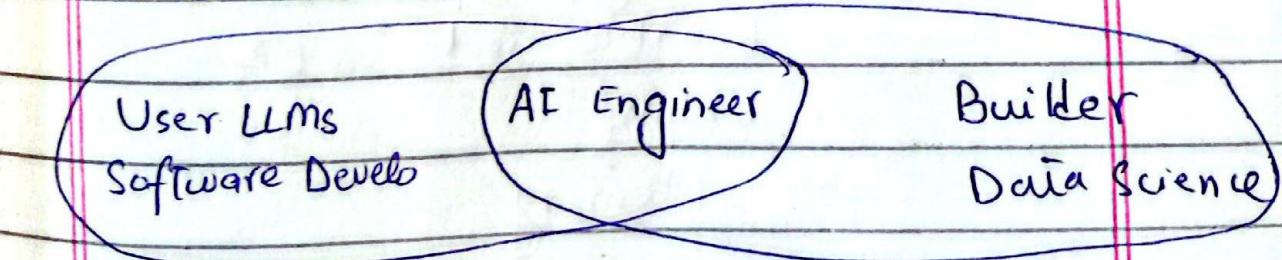
Fine Tuning

Agents (Like ChatGPT but do extra work)

LLMOps

Miscellaneous

- Do we need To learn both



The answer is Yes.

LangChain

→ User's perspective

→ first step because it

integrate with all other Tool's

LangChain is an open-source framework that helps in building LLM's based application. It provides end-to-end tools to developers build complex AI applications such as chatbots, question-answering systems, RAGs, support all LLMs. Support all major CreAI usecases.

PlayList Videos

① Fundamental :-

- i) What is LangChain
- ii) Models
- iii) Prompts
- iv) Parsing Output
- v) Runnable and LCEL
- vi) Chains
- vii) Memory

② RAG

- i) Documented Leaders
- ii) Text Splitters
- iii) Embeddings
- iv) Vector Stores
- v) Retrievers

via Building a RAG application

③ Agents

i) Tools and Toolkits

ii) Tool Calling

iii) Building on AI Agent

Why do we need Langchain
reader Pdf

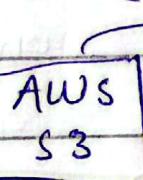
option chat or as LLM

web 2014 - 2023

- ESWC 2023

- User

PDF upload



store
split into parts



Page → Embed

Page → Embed

Page → Embed

all embed

database emb
relate emb

database

emb
Question

User query

Semantic Search

Page related
System Query
+ Query

Brain
NLG + CA IESI

Final Output

- Challenge 1 - Building Brain
[query understanding + relevant text semantic]

So we just add LLM's which have these capacity.

- Challenge 2 - LLM's Weight
It is very heavy. On Loading Server it take many computation. So many companies faces these solution and get solution.

Now you can not load on server already loaded on google server you can use api.

- Challenge 3 - Organization of these components.

OpenAI code is for slicing API
it's like this key (OpenAI) is embedding
or perform model so it's an embed L-WG change
the code is as follows
- 6% code
- 6% Langchain is 5%

Langchain has the built-in
functionalities like component (LLM)
change models, the code is as follows

Benefits

i) Concepts of Chains

→ Big tasks, big components, small tasks.

→ first task in pipeline will perform

perform task p → next task in chain

→ output task → ...

→ (i) input task

ii) Model Augmented Development

→ model of the task process

→ change model

iii) Ecosystem Computed:

→ ecosystem of big & longchain
excel, pdf, etc., load document → Doc Loader
avail models, doc embeddings, etc.,
etc., in database - OS

iv) Memory and state handling

→ linear seq of conversation

→ mention of LR in conversation

→ E in related conv?

- What can you Build?

i) Conversation Chatbots

chatbot → like call center

ii) AI knowledge Assistant

also trained on your data.

iii) AI Agents

also give some suggest
also do work for you
book your tickets like a
person or agent talk like
human work like human.

(iv) Work flow automation

(v) Summarization & Research helpers

- Alternative

Llamaindex.ai

Haystack.

Video #2 Components of LangChain

1. Models.
2. Prompts.
3. Chains.
4. Memory
5. Index.
6. Agents.

1. Models

are core interface which interact AI.

NLP task \Rightarrow ChatBot

Challenge 1 : Natural Language Understanding

Challenge 2 : Reply \Rightarrow Context aware text generation

Problem \rightarrow Challenges solved by LLM's

Challenge 1 :- Huge size of LLM's

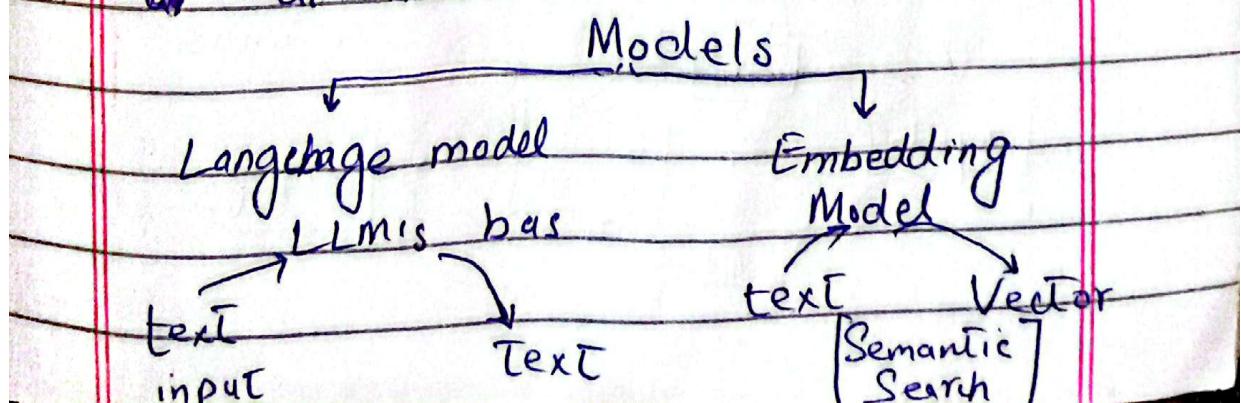
which small company drive on server

Solution :- Builder company make API

around LLM's.

Challenge 1 :- Implementation:-

Langchain make interface that everyone communicate with every ~~ai~~ company ai model



2. Prompts:-

$LLM \rightarrow \text{input} = \text{prompt}$

Chatgpt: What is langchain, = prompt

Prompt is very very important

So prompt Engineering comes

Types:- i) Dynamic and Reusable (Summarize topic)
ii) Role-based

Hi you are experienced {profession}

Tell me about {topic}

iii) Few Shot

3. Chains:-

To build pipeline

LLM application \rightarrow pipeline

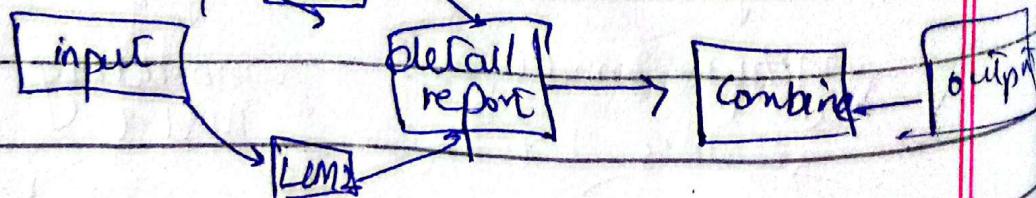
input \rightarrow LLM \rightarrow translated \rightarrow LLM \rightarrow output

Use chain which do ~~one~~ one input other is output

Complex pipeline

parallel
chain

Conditional
chains.



④ Indexes:-

Indexes connect your application to external knowledge - such as PDFs, website or databases.

Components of Indexes:-

- Document Loader
- Text Splitter
- Vector storage (Database)
- Retrievals

ChatGPT has no excess to external sources of your company So, it cannot answer you about your company privacy.

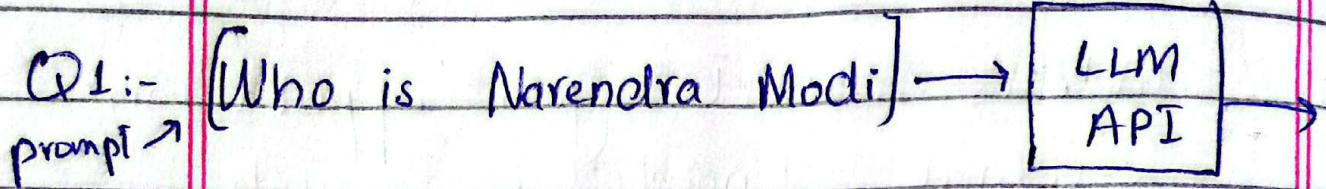
So we build LLM's and external knowledge To these LLM's, so it trained on these data also.

"Indexes is the way To build LLM based application which has excess to external knowledge source."

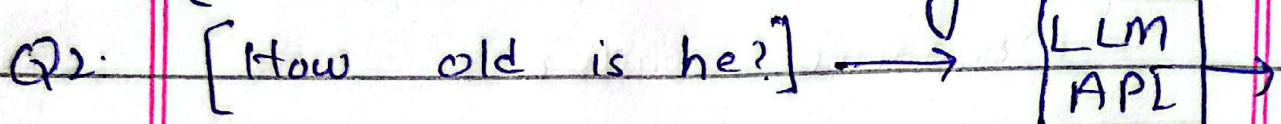
⑤ Memory

LLM's API calls are stateless
 - Each call, LLM API forgets previous's stateless
 - One call previous's stateless independent request
 - Own understanding

compon = 6' Langchain & C solve problem
This is the main problem. . Model



Narendra Modi is the Indian
Prime Minister Since May 2014



As an AI, I don't have access

To personal data about individuals
unless it has been shared with me

i) Conversation Buffer Memory:-

→ stores interaction history
→ uses to call next API

→ uses to call next API

ii) Conversation Buffer Window Memory

→ stores last N interactions

→ uses to call API

iii) Summarize-Based Memory:-

Summary of older chat

iv) Custom Memory:-

AI Agents

LLM property → NLU → Text generation

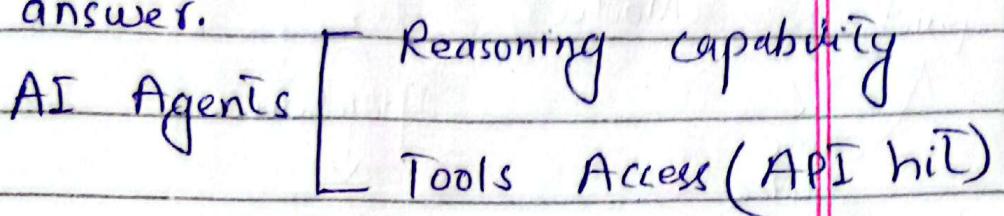
Chatbot

Chatbot :- \rightarrow AI with super power

AI Agents = \rightarrow AI with super power

(Chatbot with super power)

When you ask question, AI agents hit on API and do answer.



Tools => calculator, Weather API

Q: Can you multiply 3 with temperature of Lahore?

First he check an API who give him temperature.

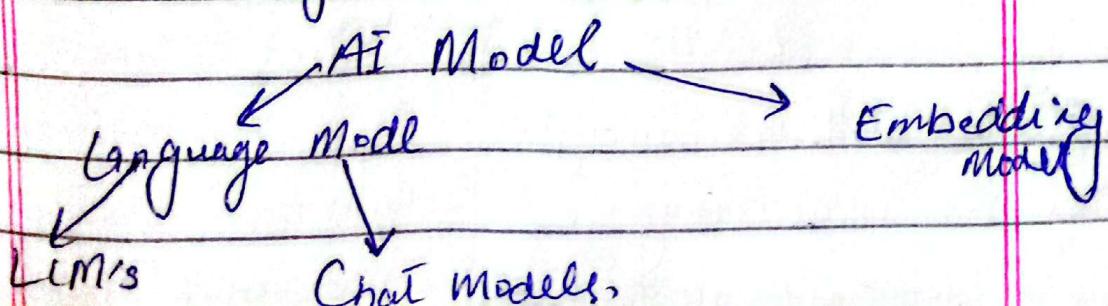
then he search calculator and he call calculator.

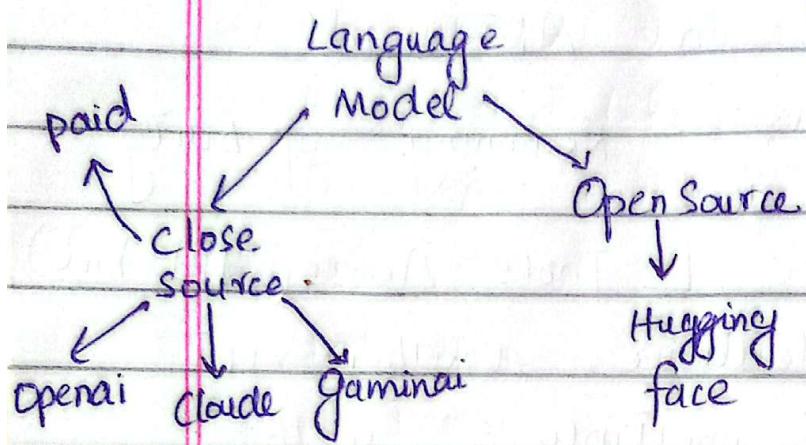
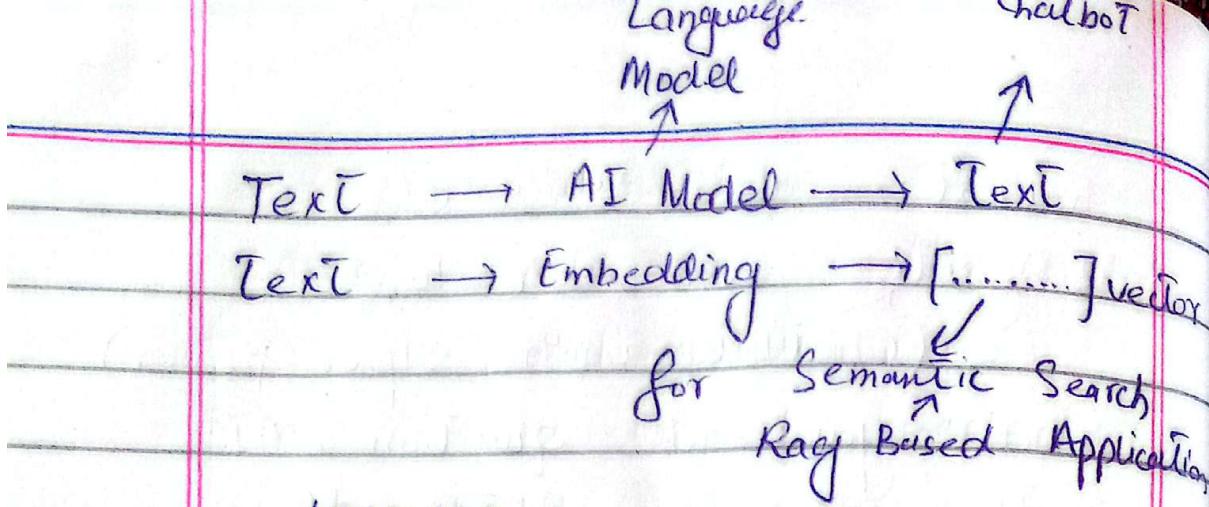
$$= 25^{\circ}\text{C} \times 3 = 75$$

Video # 3

Model Component

"Model component in Langchain is nothing but an interface which helps you connect with different kinds of AI Model."





Language Model.

Types :-

i) LLM's

ii) Chat Model

LLM's :- General purpose model that used in raw-text-generation. These are older model.

Chat Model :- Language Model that are specialized for conversational task. They take a sequence of message as input and return chat messages as output

LLM's

Chat Models

No built-in memory

No understanding of

user and assistant rule

Support structured

Understanding system

User and assistant rule

Use Cases:-

- Text generation, summarization, translation creative writing, code generation
- Conversational AI, chatbots, virtual Assistant customer support, AI tutor.

Setup:-

LLM's Models

- i) Openai
 - ii) Anthropic
 - iii) Google
 - iv) Hugging face → Opensource
- Openai - Langchain

from langchain-openai import OpenAI

from dotenv import load_dotenv file

* dotenv load your api key from .env

~~dot~~ load_dotenv()

llm = OpenAI(model="gpt-3.5-turbo-", insta)

temperature=(0 - 2) => range, max-complexity

result = llm.invoke("This is your prompt") words

invoke is everywhere in Langchain

for every ai model to give

prompt of ai model and get answer

print(result).

"string input → string output"

llm

Chat Models

same code as Clm code.

- from langchain-openai import ChatOpenAI
- from dotenv import load_dotenv
load_dotenv()

model = ChatOpenAI(model='gpt-4')
result = model.invoke("Something")

print(result)

Anthropic and ChatGPT
generative ai is same as
code this openai

Open Source AI Model

are freely available AI
model than be downloaded,
modified, fine-tuned and
deployed without restriction.

"Closed Source Mode is on the
server, on builder company. These
company gives API to user and
get back many. These model
running on server"

"Closed Source Mode is on the
internet, not on company server. It is

uploaded on internet. you can just download them.

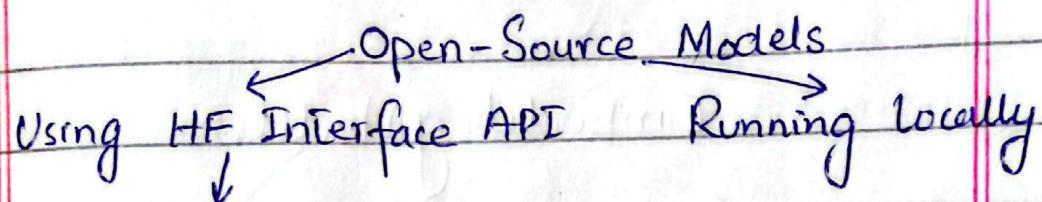
" You have control on it".

" You have data privacy. of your company."

Where To find?

On Huggingface.

Largest repository of open source LLM's Models.



API like OpenAI but free. you can use thousand of model.

- Disadvantage:-

High Hardware Required

Step Complexity

Lack of RLHF

Limited Multimodel Abilities

Code is same as other.

$\text{Um} = \text{HuggingFaceEndpoint}(\text{rapid_id} = \text{tinyLlama} / \text{tinyLlama-1.1B},$

$\text{task} = \text{'text-generation'}$

$\text{model} = \text{ChatHuggingFace}(\text{Um} = \text{Um})$

$\text{result_model.invoke(...)}$

Model:

```
Locally HuggingFace downloaded  
import HuggingFacePipeline  
lm = HuggingFacePipeline(  
    model_id = ' ',  
    task = 'text-generation',  
    pipeline_kwarg = dict(  
        temperature = 1.5,  
        max_new_tokens = 100  
)
```

0 = same
0.1 = some changes
0.7 = more changes
1.5 = completely changes!

Embedding Models.

```
from openai import OpenAIEmbed  
from dotenv import load_dotenv  
model = OpenAIEmbedding(model='text-embedding-avata')
```

result = model.embed_query('What is Campus X')
print(result)

1536 - for small embedded model

3072 - for large model,

you can also do multiline query

```
document = ['...', ',
```

```
'.....',
```

```
'.....; ]
```

~~result = embedding.embed_document(document)~~

Video # 4: What are Prompts:-

Prompts are input instruction

give to model to guide its output.

Prompts → Text-based sound
Multimodel (image, pdf)

- Static prompts.

```
from langchain_opai import ChatOpenAI
```

```
from dotenv import load_dotenv
```

```
import streamlit as st
```

```
load_dotenv()
```

```
st.title("CHATBOT HELPER")
```

~~ver~~ input = st.text_input("Enter your prompt")

```
model = ChatOpenAI(model='gpt-4')
```

~~if~~ st.button("Summarize"):

```
result = model.invoke(user_input)
```

```
st.write(result.content)
```

~~for i in range user_input if prompt static~~

~~then if user input == user_out put - if~~

~~if - else if result == user_out put~~

use `research` paper web =
give `style` - & give `research` paper =
1. `Umig` & `text` file
2. `result`
- `template` col

Dynamic Prompt

"Please summarize the research paper titled "`{paper-input}` with the following specification:-

Explanation-style : `{style-input}`

`input` < user `for` `style` `input` ↪
`prompt` = `style` `input` `for` `Umig`
↳ `Umig` & `text` file
↳ `style` `input` < `user` or `key` values

Prompt-Template

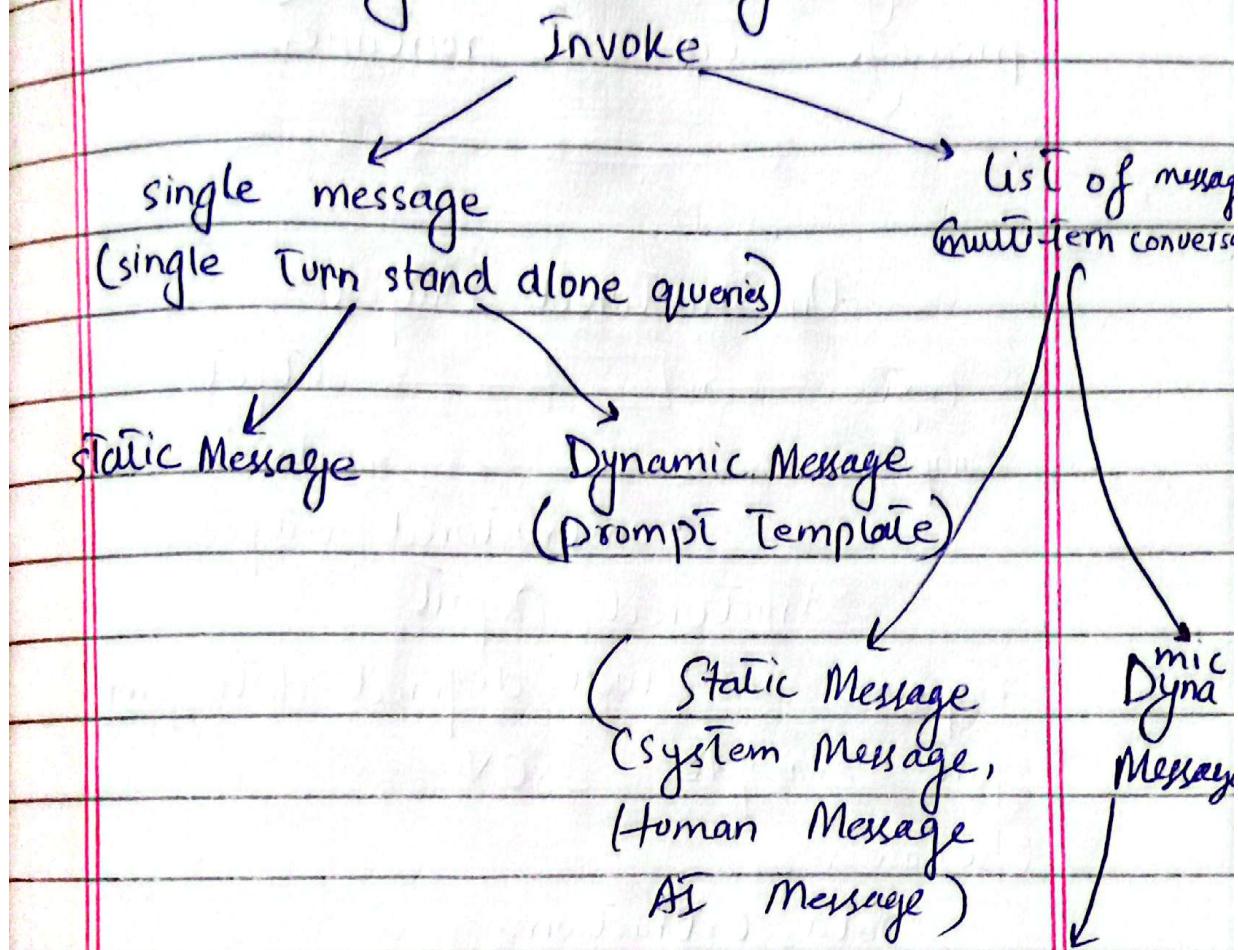
A prompt template in LangChain in a structured way to create a prompt dynamically by inserting variable into a pre-defined template. Instead of hardcoding template allow you to define placeholders that can be filled in runtime in different path.

Why use prompt template
over f-string

i) Default validation

ii) Reusable

iii) ChaiChain - Ecosystem.



PromptTemplate :-

Chai PromptTemplate

PromptTemplate is used for
single query.

Chai PromptTemplate.

While ChaiPromptTemplate
is used for multiple query using
messages.

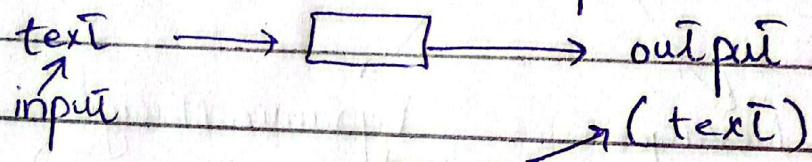
Message Placeholder:-

in Langchain used

inside a chatprompt template

To dynamically insert chat-history or a list of message to at runtime.

Un Structured Output



unstructured output

Structured Output

response in well-defined data format
for example JSON.

Use cases:-

Data Extraction:-

API Building :-

Agents:-

Ways to get Structured Output

LIMs

Can

Can't

like OpenAI

CPT Model.

Langchain has

Output parsers.

with structured output

with structured output parameter

model invoke K waalk hum

ye duty hain our sathe main
data format

TypedDict Pydantic JSON

Typed Dict. → Simple

is a way to define dict in python.

Class Person(TypedDict) which the type
name: str → of the key.
age: int

Annotated:-

is the type of dict
which is used to give description

to the TypedDict parameter

Class Person(TypedDict):-

name: Annotated[str, "Name of Author"]

age: Annotated[int, "Age of Author"]

model: Annotated[dict, "Model description"]

Pydantic: Simple:-

Pydantic is a data validation and
data parsing library in python. It
ensures that the data you work with
is correct, structured and type-safe.

Default, Optional, Coerce, Building Field Function

```
from peewee import BaseModel  
class Student(BaseModel):  
    name: str = 'nitish'  
    age: Optional[int] = None  
    email: EmailStr  
    cgpa: Field(gt=0, lt=10, descrip="")  
new_student = {'age': 32, 'email': 'ab'}
```

JSON-schema

{

"title": "Student",

"description": "",

"type": "object",

"properties": {

},

"required": ["name"]

}

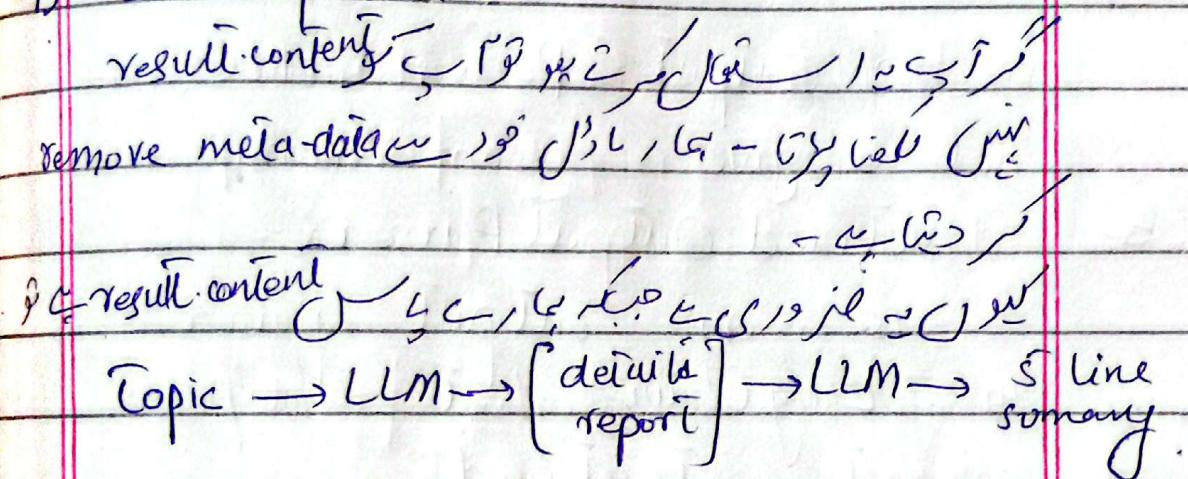
to integrate language (python)
use file '16 JSON-schema'

OpenSource Model
with structural output
Object with 2 meta-data
- Configuration parameter

Output parsers:-

Output parsers help to convert raw LLM responses into structured format like JSON, CSV, Pydantic models and more. They ensure consistency, validation and ease of use in application.

-1) StrOutputParser:-



```
parser = StrOutputParser()
```

```
chain = template1 | model | parser | template2 |
```

model | parser

```
chain.invoke({ "topic": "Black Hole" })
```

-2) JSONOutputParser

```
Parser = JSONOutputParser()
```

template = Prompt | Template

template = 'Give me the name, age and city of person, \n {
 } format instruction',

input-variables = [],
partial-variables = {format instruction} :
parser.get_format_instruction
prompt = template1.format()
result = model.invoke(Prompt)
final-result = parser.parse(result-content)
print(final-result)

JSON Output Parser does
not give json-schema.

↳ Structured Output Parser:-

is an output parser
in LangChain that helps
extract structured json data
from LLM responses based on
predefined field schema.

using `pydantic` schema

["fail", "description": "Response schema"]

Disadvantage:-

- Precise, strict data validation

4)

Pydantic Output Parser
 is a structured output parser
 that uses pydantic models to
 enforce schema validation when
 processing LLM responses.
 data validation]- Pydantic
 schemas

Seamless integration = Works well
 with other Langchain components.

Easy Validation = Uses Pydantic's built-in
 validation to catch incorrect or
 missing data.

Video # 7

Chain

Get input from user

- Simple | Sequential | prompt | model | parser

↳ 1. get it pipeline & help (Qn)

↳ output b' step w. ↳ 2. 3. 4. steps ↳ ↳
 - automatically get input b' step

- Sequential :-

↳ 1. get - Qn 2. / call model, ↳ 3. get ↳

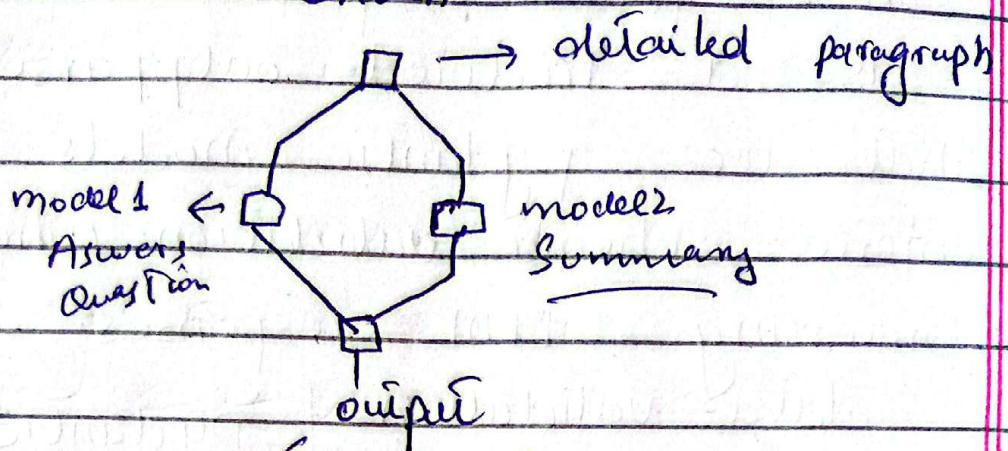
- " " " prompt ↳

- ↳ simple ↳

↳ " " " - details ↳ text ↳ ↳ user use ↳

- ↳ ↳ ↳ ↳ ↳ summary ↳ 5-point ↳ text

parallel Chain



Runnable parallel \leftarrow i.e. \leftarrow in parallel chain
 \leftarrow \leftarrow \leftarrow \leftarrow \leftarrow

parallel_chain = RunnableParallel({

'notes': prompt1 | model1 | parser,

'Ans': prompt2 | model2 | parser,

)

merge_chain = prompt3 | model1 | parser

chain = parallel-chain | merge-chain

chain.invoke ({'text': text}).

generate notes or |||, notes \leftarrow g. prompt1 \leftarrow b3 |

||| b "generate prompt1" \rightarrow Answer, ||| b |

||| b | merge d \leftarrow (j) Runnable parallel
jed b " > ans (j) prompt \leftarrow j \leftarrow g. output

quieres ||| notes i f. already or
input (j) j \leftarrow b, \leftarrow input variable

\leftarrow j \leftarrow last chain \leftarrow ||| b2 (j),
 \leftarrow b " \leftarrow j (j) output \leftarrow prompt.

Conditional Chains

فیلم = کارہر جسے بنا کر بنیں (conditions)

- "E chain = فیلم کو اپنے chain میں chain کا

- 052 / (E chain) RunnableChain سے لے کر 051
کو اپنے message اور feedback کا

1116" E chain message اور feedback negative

- 6" E chain message اور positive

= RunnableChain(

(,chain),

(,chain);

RunnableLambda()

- ڈیم / (یہ فیلم کو اپنے chain میں کو Runnable کی

concept by runnable کو دیکھ کر یہ کو chains

کے میں لے لے لے

- اسی طرز کو اپنے Runnable

2022 → ChatGPT realize

→ OpenAI API → build LLM
application.

AI Agent,

ویسے PDF Readers ChatBot, AI Assistant

پر، اسے رہا کر لے لیں اسی LangChain کو

کو اپنے help کو developers

کے لئے LLM's کی Company کو identify کی

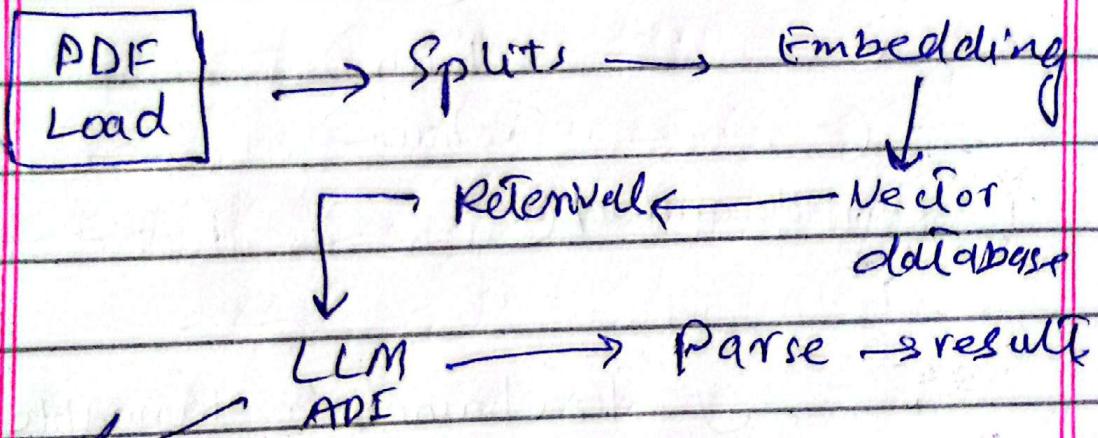
لے لے لے

API LLM's company's framework
problem we can't solve

i languagechain's problems
of us & application's to identify
only i'm interested API

For example: task, we can get 101

PDF Reader



gets only one component.

but we can't task b'c' t'i' j'k' l'

in pdf we can't find the component
parse by component PDF Loader
(i) developer -> in parser
- it's not simple application complex

it's 16' component in AI Engine

it's in application w/o d.o.t

which is languagechain's problem

the application's of notice

U.S. prompt $\omega = (\mu, \omega_0, \nu) \mapsto \omega / \nu^2$

أفضل فيلم will be called the best film

وہ میں کے میں built-in function فری دیکھو
اکثر call of um, prompt پر کھلے گئے

- LLM (chain) / If result,
- ~~chain~~ (um, prompt)

ویرایشی Components لیست می شوند

longchain f (at task simplech¹ g) = (g¹)

application RACI and related to the chain complex

user triggers task 6's retrieval

• At the same page (the document) will go to query

Um C merge of C++ and Java
is called as call of

Diagnostic team helped call it increased function with

19) query \hookrightarrow Tries Out C.R. function w/
 \Rightarrow JDBC Chain w/ get / & retrieval
 ... (or a moment by result ...)

- CCR (automated regulation) $(U_m = U_m)$
- PAA chain (retention, ~~efficiency~~)

Epoxy chain length اینچوں گلی

task = نکار و کار کیوں کیا جائے گا

chain' $\cup_{i=1}^n c^i \cup_{j=1}^{m-i} b_j$, reverse \rightarrow b_1

Simple Sequential Chain

llm $\xrightarrow{\text{call}}$ joke $\xrightarrow{\text{llm call}}$ Explanation.

chain so + i team (langchain)

our idea = of AI in code -> ->
-> gives no problem

(AI Engineers) intent b'j'

with our help

-> plug in our component
so chain can solve problem

→ database huge

→ new learner

so chain can do
task or how to define it

call function chain

in which each component

i.e. parser -> is component

common method (standardize)

parse & call of parser

call of prompt template -> (is called JIN)

then -> in JIN format

JIN / predict > TLL API call

->

design & build which is

components & method abstract

JIN / invoke of TLL API call

& standard components

-> (like chain ps)

(Runnable)

[Task Specific
Runnables]

- These are the components like
 - Prompt Template
 - Um, parsers
 - standards

↓
Runnables

[Runnable]
Primitives

For connecting task specific components

Prompt \leftrightarrow LM
input $\xrightarrow{\text{output}}$ input

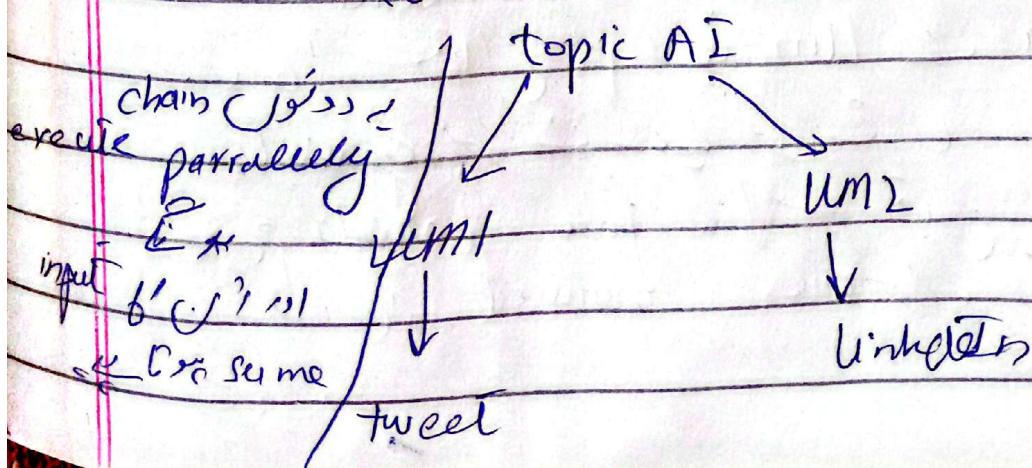
- Runnable Sequence
- Runnable Parallel
- Runnable Lambda

Runnable Sequential

↳ Runnable objects \rightarrow help to make chain = $\rightarrow - \rightarrow - \rightarrow$ connect one-by-one
input $\xrightarrow{\text{process}}$ output $\xrightarrow{\text{process}}$ $\xrightarrow{\text{process}}$ $\xrightarrow{\text{process}}$
 $\xrightarrow{\text{process}}$ $\xrightarrow{\text{process}}$

↳ Runnable Sequential (prompt, Um, parser)

Runnable Parallel:-



Runnable Parallel (

'tweet' : RunnableSequence(prompt, mode
parser)

'linkedInPost' : RunnableSequential(prompt, mode)

),

- Runnable Passthrough

is a special Runnable primitive that simply returns the input as output without modifying it.

↳ let's say b' has input b' ↳ i ↳ b'' ↳ c ↳ no output

For example

topic : AI

joke x

↓
explanation.

my requirement

joke + explanation

RunnablePassthrough

prompt → Un → parser

prompt → Un →

Runnable Lambda

extra. un Un is if file 16' or

Un -> Un → generate output

as file 16 from Lambda p. Un

or file named Un

RPC {

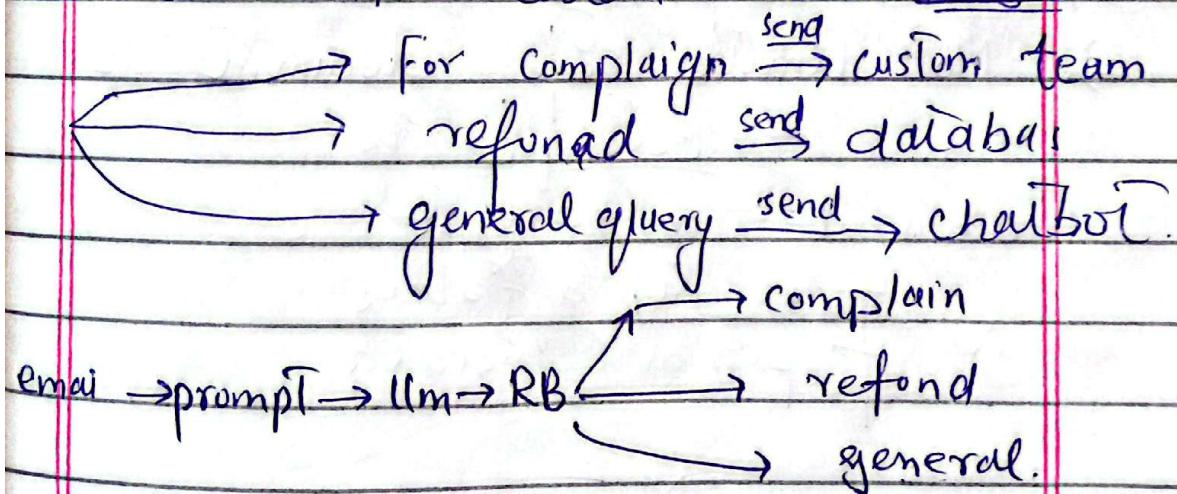
'cardai': Runnable Lambda (Lambda x: len(x.split())

(if no condition
match,
default runnable
is used)

Runnable Branch

(you define set)
of conditions) ← for conditional chain

if-else, 6' universe's LangChain
is a control flow component
in LangChain that allow you
to conditionally route data
to different chains or runnable
based on custom data. e.g.



anyone branch will triggered.
execute branch (لأى فرع، يُ執نّي)

LCEL (1)

RunnableBranch = team of LangChain

sequence (prompt, model, parser), it's like

RunnableSequence (prompt, model, parser)

RunnableSequence (prompt, model, parser)

RunnableSequence (prompt, model, parser)

= prompt | model | parser

Document Loader

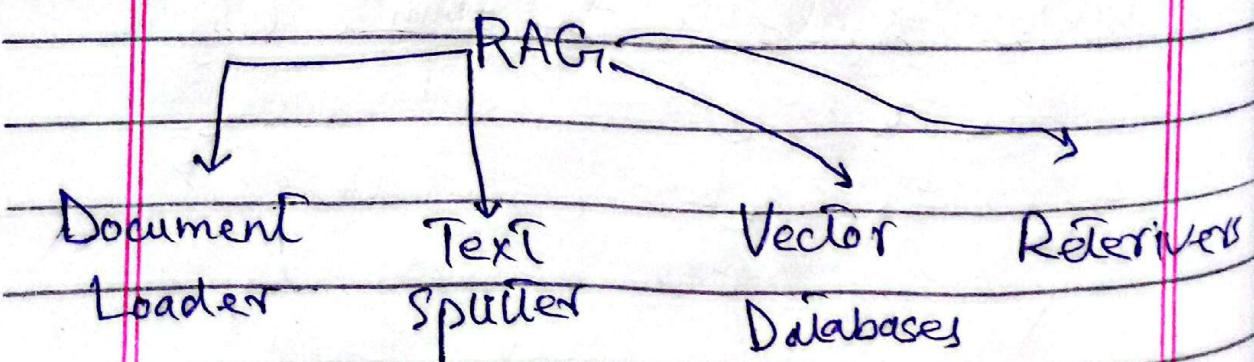
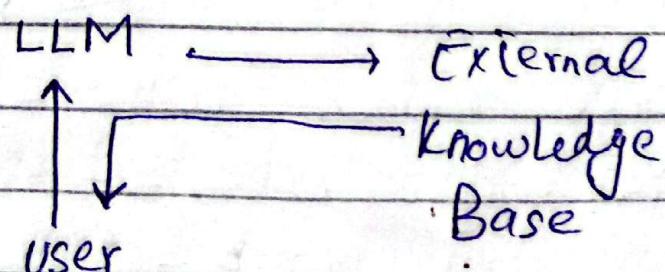
- RAG:-

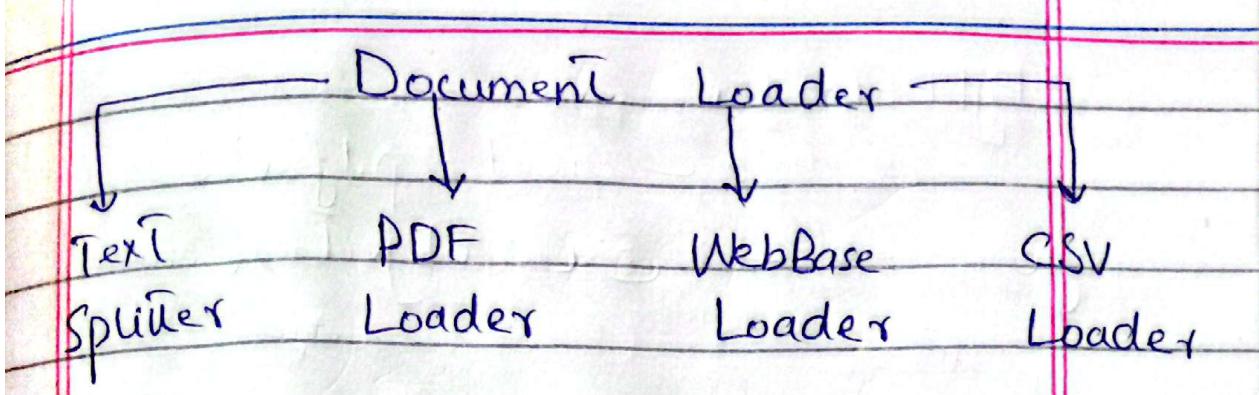
is a technique that combines information retrieval with language generation where a model retrieves relevant documents from a knowledge base and then uses them as context to generate accurate and grounded responses.

- 1). Use of up-to-date information
- 2). Better privacy
- 3). No limit of document size

Example

Chatbot → Chatgpt
 X current affairs
 X personal data





Document Loader are components in LangChain used to load data from various sources into a standardized format (usually as Document object), which can then be used for chunking, embedding, retrieval, and generation.

Document(

```

page-content="The actual text content",
metadata[{source: "file-name", ...}]
)
  
```

① Text Loader

is a simple and commonly used document Loader in LangChain that reads plain text files (*.txt) and converts them into Langchain document.

- chat logs . scraped text .
- transcript . code snippets
- output in [a]
- visi

PyPDFLoader

→ $\text{page} \rightarrow \text{list}$ → Document object

→ $\text{Document} \rightarrow \text{metadata}$

Document.page-content = "Text from page"

Document(page-content = page 2, metadata[page=1, ...])

]

Limitation:-

pages is not in PyPDF loader
"6" is in document as part of pdf
- PDF text simple.

PyPDFLoader pdf not in pdf loader

Use Cases

Simple, clean PDF's

PDF's with tables/columns

Scanned/Image PDF's

Need layout/img Data

Want best structure extraction

Recommended Loader

PyPDF Loader

PDFPlumber Loader

Unstructured PDF/AjaxPDF Loader

PYMPDF Loader

Unstructured PDFLoader

Directory Loader

use the `File` to load a directory
 directory \rightarrow \rightarrow file (or a list)
 \rightarrow txt \leftarrow files of pdf \leftarrow files
 files of directory \rightarrow file load
 \rightarrow file \leftarrow file \leftarrow file
 \rightarrow pdf file is loaded
 \rightarrow loader_cls \leftarrow file pdf
 - use `Loader_cls` \leftarrow file pdf

Directory Loader()

path = ''

glob = '* .pdf'

loader_cls = 'PyPDF'

you also
give

*.txt

Text Loader

)

Load()

- Eager Loading (loads everything)
- Return a list of document objects
- Loads all documents immediately into memory
- Best when:-
 - The document is small
 - You want everything loaded upfront

Lazy_Load()

- Lazy Loading (loads on demand).
- Return a generator document.
- Document are not all loaded at once: they're fetched once at a time as needed.

Best when :-

You are loading with large document

WebBase Loader :-

is a document loader in Langchain used to load and extract text content from web pages (URL). It uses BeautifulSoup under the hood to parse HTML and extract visible text.

- When To Use:-

For blogs, news article, or public websites when the content is primarily text-based and static.

- Limitations:-

Don't handle JavaScript-heavy pages well (Use SeleniumURLLoader for this).

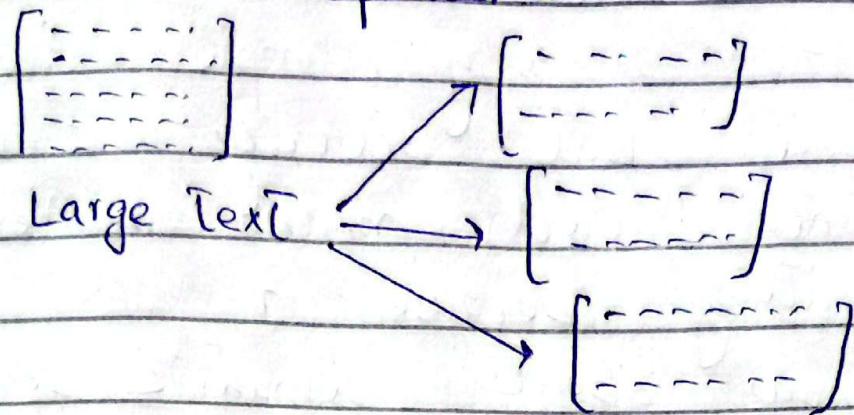
WebBase Loader:

url → page ← website ↗ ← request
- (User)

CSV Loader.

→ URL → path of CSV file,
→ CSV document (1 row per)

Text Splitter:-



is a process of breaking large chunks of text (like Articles, pdf's) into smaller and manageable pieces that LLM can handle effectively.

- Overcoming model limitation:-

Many language model and language model has maximum input size constraints. Splitting allows us to process document that would otherwise exceed this limit.

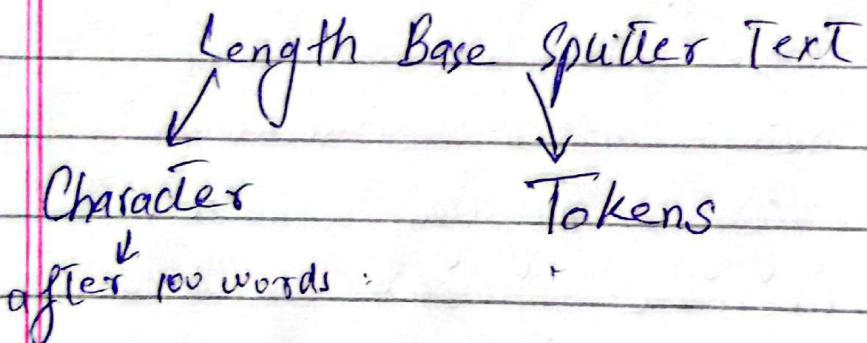
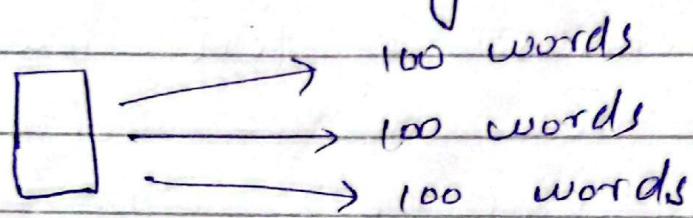
- Downstreaming Task:- Text splitting improves every LLM powered task.

- i) Embedding
- ii) Semantic Search
- iii) Summarization.

- Optimizing computational Task.

Length Base Text Splitting

It is very simple and fast but actually it is not used much in text splitting because it does not know about semantic search even not know grammar.



Text Structured Base

Recursive Character Text Splitter

\n, \n, ', -, "

para, line, space, character

Allowed

chunk-size=10

↳ chunk-size 10^{10} base of para m

↳ chunk-size 10^{10} base of line n in para m

↳ chunk-size 10^{10} base of character c in line n in para m

character, base of space etc.,

→ base of

Benefits:

→ para → tryout many
→ split sentence

Document-Structured Based

- ~ for piece of code
- ~ for markdown.

"\n class", "\n def" } first step

for class : for function

\n\n, \n, , , } second step
para, line, space

Semantic Meaning Based

Semantic = text splitting based

sentences → text-splitting based

s1 s2 s3 s4 s5
0.9 0.8 0.85 0.1
embedding model.

semantic windows
sliding window

semantic

text splitting

std::vector<float> semantic-score()

text-splitting → 1. 1. 1. 1. 1.

semantic → 1. 1. 1. 1. 1.
embedding models

Vector Stores

store just or in vectors from Embedding vectors v_1, v_2, \dots, v_n in \mathbb{R}^d
- just or in databases functional
semantic search as like v_1 first v_2 \dots
 v_n \neq v_1

A vector store is a system designed to store and retrieve data represented as numerical vectors.

- Storage :- Ensure that vectors and their associated metadata are retained, whether in-memory and on-disk.
- Similarity Search:- Help retrieve the vectors most similar to query vector.
- Indexing:- Provide a data structure or method that enables fast searching on high-dimensional vectors (ie approximate nearest neighbour lookups)
- CRUD Operations

Use Cases:-

- i) Semantic Search , RACI
- ii) Recommendation System
- iv) Image / Multimedia Search.
- A vector database is effectively a vector store with extra database features (e.g. clustering, scaling, security, metadata filtering and durability)

Vector Stores in LangChain:-

- Supported Stores :- Langchain integrate with multiple vector stores (FAISS, Chroma ie).

Common Interface

Metadata Handling:-

Most vector stores allows you to attach metadata, i.e Timestamp, author.

Chroma:-

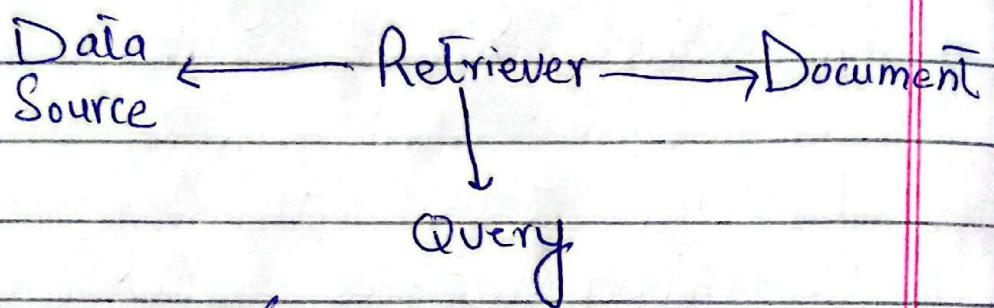
Chroma is a lightweight, open-source vector database that is especially friendly for local development and small to medium scale production need.

Video no # 13

(Retrievers)

A retriever is a component in Langchain that fetches relevant documents from a data source in response to a user's query.

All retrievers are runnables.



- ~~use the retriever to query~~
- ~~if you have 1 data source in your app~~
- ~~it will do the searching or it~~

— Wikipedia Retriever

is a retriever that queries the wikipedia api to fetch relevant content for a given query.

= How it works:-

- You give it a query (e.g. "Albert Einstein")
- It sends the query to wikipedia api

• It retrieves the most relevant articles.

• It returns them as Langchain Document object.

$\text{ret} = \text{WikipediaRetriever}(\text{top-k-result}=2, \text{LangChain})$

- Vector Store Retriever:-

is the most common type of retriever that let you search and fetch document from vector store based on semantic similarity using vector embedding.

= How it works.

- You store your document in vector store (FAISS, CHROMA, Weaviate)

- Each document is converted into dense vector using embedding model.

- When the user enter query:-

 - It also turn into a vector.

 - retriever compare the query vector with the stored vector.

 - It receives the Top-K-most similar ones.

→ Minimal Marginal Relevance

"How can we pick results that are not only relevant to the query but also different from each other."

MMR is an information retrieval algorithm designed to reduce redundancy in the retrieved results while maintaining high relevance to the query.

→ Why MMR Retriever:-

In regular similarity search, you may get documents that are:

- All very similar to each other
- Repeating the same info.
- Lacking diverse perspective.

MMR retriever avoids that by . Picking the same relevant document first

- Then picking the next most similar and least similar to already selected docs: and so on.

This helps especially in RAG pipelines where:

- You want your context window to contain diverse but still relevant information.

- Especially useful when documents are semantically overlapping

- Multi-Query Retriever:-

"Sometimes a single query might not capture all the ways information is phrased in your document."

Query :-

"How can I stay healthy."

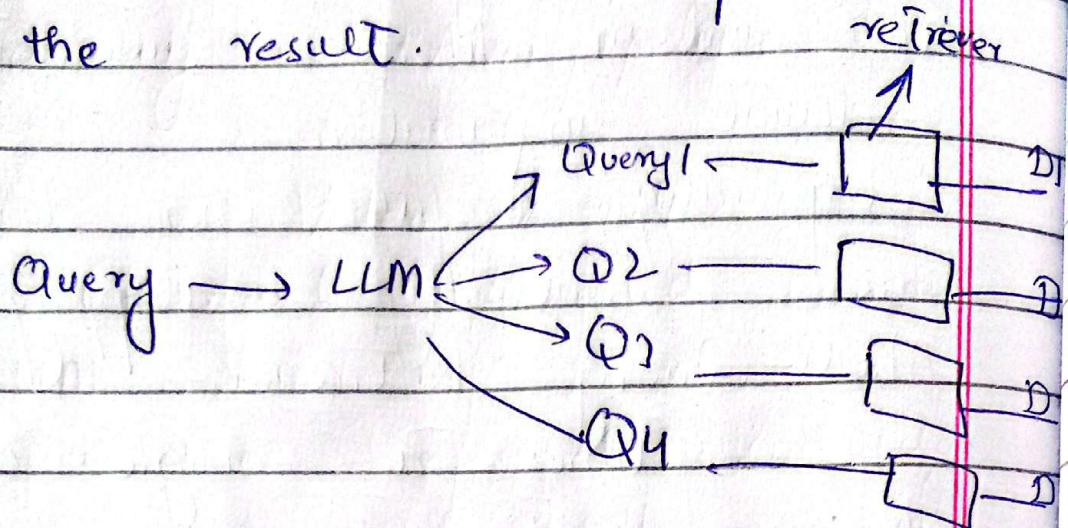
Could means :-

- What should I eat.
- How often should I exercise.
- How can I manage stress.

A simple similarity search might miss documents that talk about those things

- Take your original query.
- Use an LLM to generate multiple semantically different versions of that query.

- perform retrieval for each subquery.
- combines and deduplicates the result.



— Contextual Compression Retriever is an advanced retriever that improves retrieval quality by compressing document after retrieval - Keeping only the relevant content based on the query.

? 'Query'

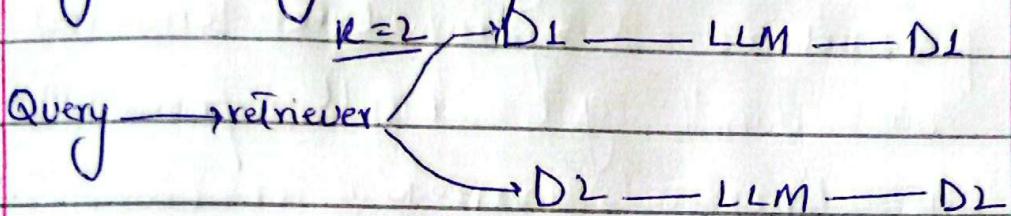
'What is photosynthesis?'

Retrieval Document (by traditional retrieval)

"The grand Canyon is famous natural site. photosynthesis is how plants convert light into energy. Many tourists visit every year."

Problem :-

- The retriever returns the entire paragraph.
- Only one sentence actually relevant to the query.
- The rest is irrelevant noise that waste context window any may confuse the LLM.



Contextual Compression Retriever does:-

- Returns only the relevant part "photosynthesis is how plants convert..."

= How it works:-

- Base Retriever (e.g. FAISS) retrieves N-Documents
- A compressor (LLM) is applied to each document.
- The compressor keeps the part only relevant to query.
- Irrelevant content is discarded.

When To USE:-

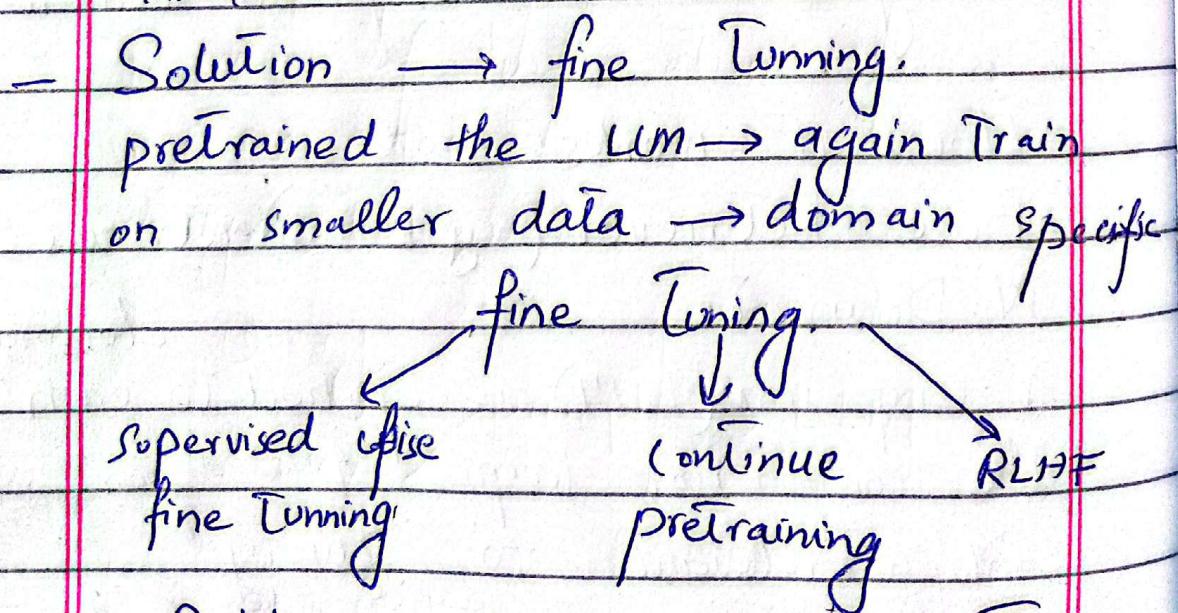
Your documents are long and contain mixed information.

- You want to reduce context length for LMs.
- You need to improve answer accuracy in RAG pipelines.

Video no #14

(Retrieval Augmented Generation)

- LM store information in parameters → parametric knowledge
- You can't get result from LM on these scenarios.
 - i) Private data,
 - ii) Recent Information
 - iii) Hallucination.



- = Problems Again on Fine Tuning
- i) Training on large data computation expensive
 - ii) Technical expertise.

Solution:- Again

In context Learning:-

Solve Task by seeing examples in prompts without updating its weights.

Prompt Example:-

Please sentiment these:

I love my school → positive.

This app crashes a lot → negative

I hate movies →

few shot prompting

It is an emergent property.

PLM (3.3) → examples (prompt) → "I

LMG model (Lc) → (S) context

- 6" → improve

RAG is to make a language model smarter by giving it extra information at the time you ask the question.

- Understanding the RAG:-

i) Indexing

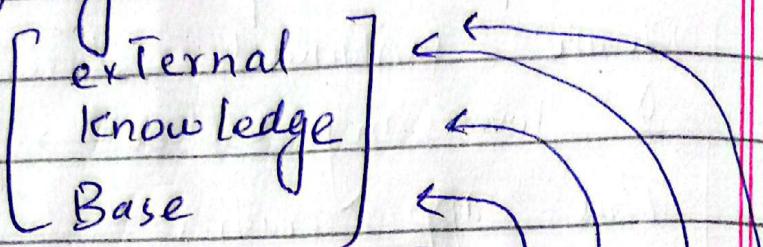
ii) Retrieval

iii) Augmentation

(iv) Generation

— Indexing :-

(is a process of preparing your knowledge base so that it efficiently searched at query time.

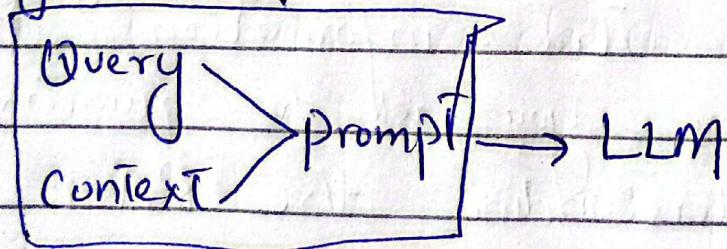


- i) Document Ingusion.
- ii) Text Chunking
- iii) Embedding Generation
- iv) Storage in Vector store

— Retrieval

finding the most relevant pieces of information from vector store.

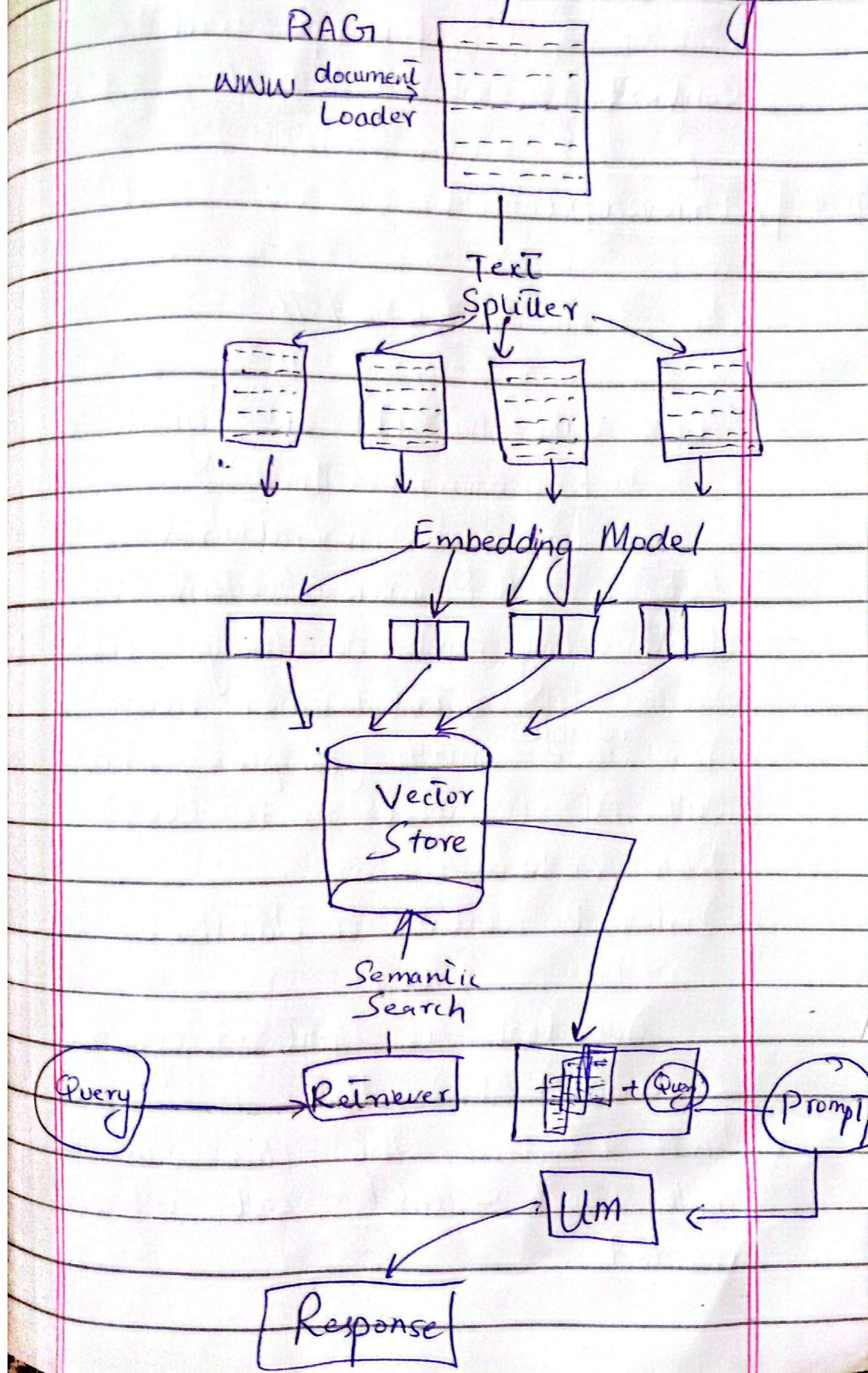
— Augmentation:-

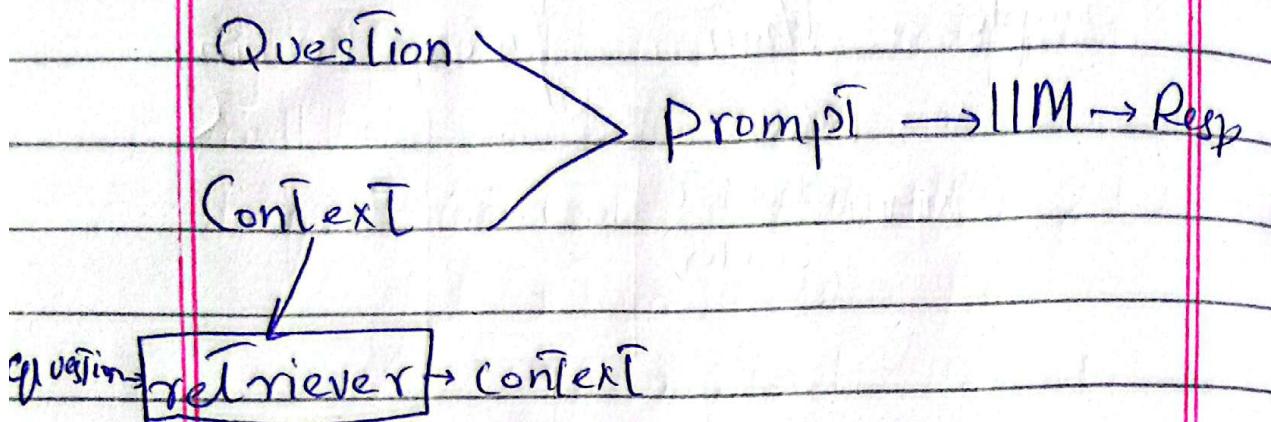


— Generation

Prompt → LLM → Response

Now RAG solve all problem 157
and it is cheaper and
simple than fine tuning.



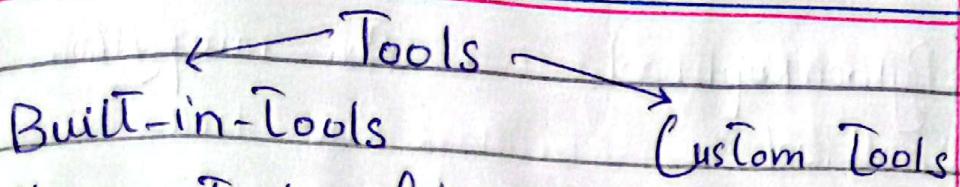


Video no # 16

Tools

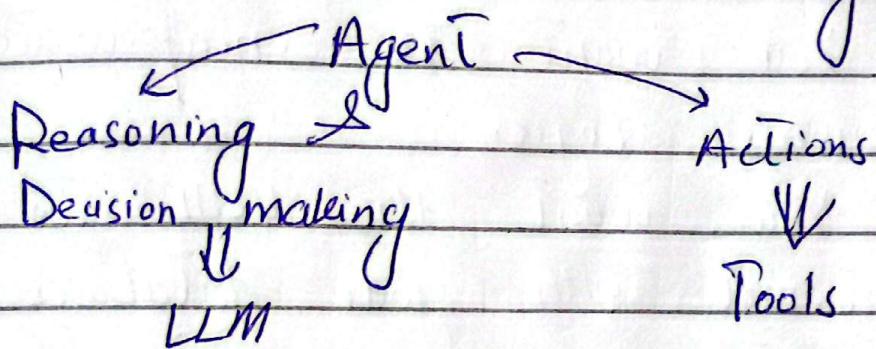
- LLM (like ChatGPT) are created at:-
 i) Reasoning (Think)
 ii) Language generation (speak)
 that is all power of LLM only.
- But they can't do things like
 • Access live data (weather, news),
 • Do ~~reliable~~ math (complex math)
 • Call APIs (tweet on tweeter)
 • Run code
 • Interact with a database

A tool is just a python function (Or API) that is packaged in a way the LLM can understand and call when needed.



How Tools fits into the Agent ecosystem.

"An AI Agent is an LLM-powered system that can autonomously think, decide and take actions using external tools or API's to achieve a goal."



Built-in-Tools:-

A built-in Tool is a Tool that LangChain already provides for you — it's pre-built, production-ready, and requires minimal or no setup.

You don't have to write the function logic yourself, — you just import and use it.

- DuckDuckGo Search Run
- Wikipedia Query Run

These tools are also runnable.

- Python Repl Tool Run raw python code
- Shell Tool Run shell commands.
- RequestsGet Tool Make HTTP get requests

Custom Tools

A custom Tool is a tool that you define yourself.
Use them when:

- You want to call your own API's
- You want to encapsulate business logic.
- You want the LLM to interact with your databases, product or app.

from langchain.core.tools import Tool

- create function
- doc string
- type hint
- add decorator.

a tool make this a special function that can communicate with LLM.

of types of LLM's & tool types
json tools will be converted to tool in
the schema

Ways To create custom Tool:-

(i) using a tool decorator (ii) using StructuredTool
 (iii) Using BaselTool Class.

- "A structured Tool in Langchain is a special type of Tool where the input to the Tool follows a structured schema, typically defined using a Pydantic model."

- "BaseTool is the abstract base class for all Tools in Langchain. It defines the core structure and interface that any Tool must follow, whether it's a simple one-liner or a fully customized function." All other Tool types like aTool, StructuredTool are built on top of base Tool.

Toolkits:-

A toolkits is just a collection of related Tools that serve a common purpose -

best benefit.

package together for convenience
and reusability

- i) GoogleDriveCreateFileTool
- ii) GoogleDriveSearchTool

GoogleDrive Toolkit.

Video no #17

Tool Calling

- Tool Binding

- ~~Tool binding~~ -
- ~~Tool binding~~ -
- ~~Tool binding~~ -
- ~~Tool binding~~ -

- Tool Calling:-

is the process where
Hm decides during a conversation
or task that it needs to
use a specific tool - and
generate a structured output
with

- the name of the tool
- and the arguments to call
it with.

- * The LLM does not actually run the tool. It just suggest the tool and the input arguments. The actual execution is handle by you or Langchain.
- Tool Execution:-
- Tool Execution is the step where the actual python function is run using the input arguments that the LLM suggests during Tool culling.

Video no # 20

AI - Agents in Langchain problem:-

Agenda - LS-Aug-2025
 i) Problem solving of tool with LLM models.
 ii) FAQ's Chat Assistant
 iii) House price prediction
 iv) Pakistan house price prediction updating.

- Cible, web, trip
- goal high level
- call tools external or source API or
- maintain agent context or

$$\text{Agent} = \text{LLM} + \text{Tools}$$

Reasoning

API, database, Search Tools

Rethinking
↑
↓ Planning

AI - Agents:-

- i) ~~Ex.~~ Goal-driven (v) Adaptive
- ii) Autonomous planning
- iii) Tool using
- iv) Context Aware
(Maintains memory across steps).

Steps:- IIM + Tool

create-react-agent, Agent Executor
hub.

prompt = hub.pull("hewchase 17/react")

- Agent = create-react-agent (IIM, tools, prompt)

→ is the main guy → for reasoning

= Agent-Executor (agent, tools, verbose)

→ work actually. or execute
the agent.