

Getting Started

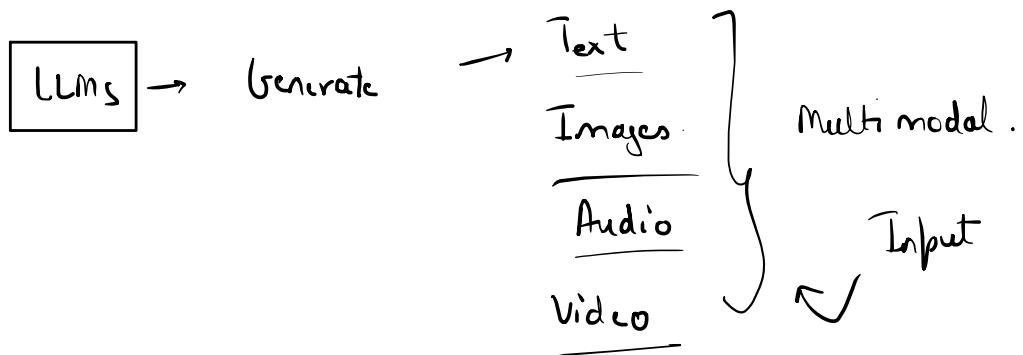


Image captioning , Video → Transcript

Audio → Video

2) Memory → Context aware

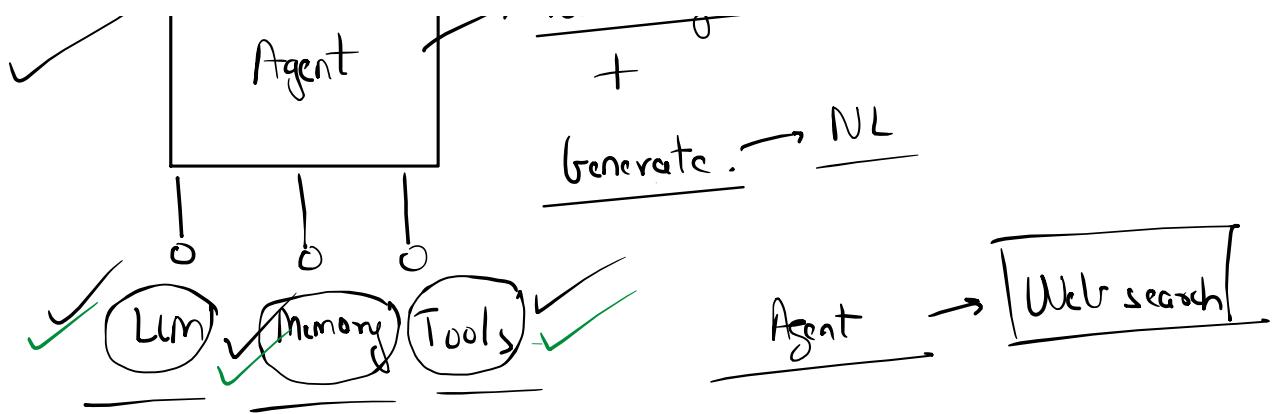
↓
RAG → Documents → QA

3) Reasoning → Cognitive ability → Reactive
↓
Proactive

Complex task → Multiple simple tasks

4) Tools → Outer world. → Do task / take action.

✓
Agent → Reasoning +



Agno → Python Framework

↓
Declarative

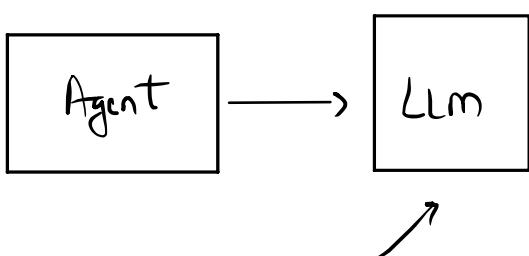
- i) brain → LLM
- a) Memory → Context Aware
- 3) Tools.

Brain

Agent :

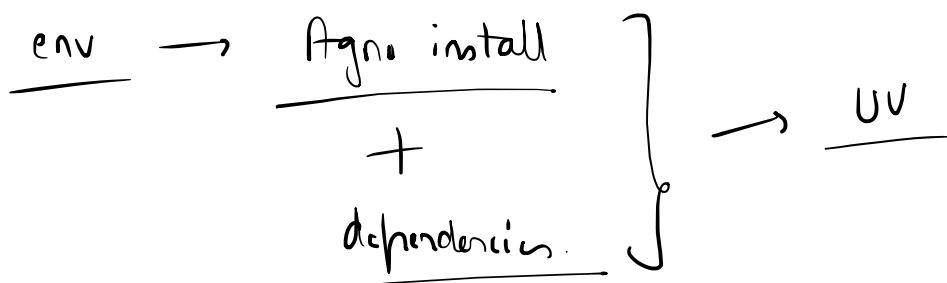
Brain → LLM

- 1) cognitive ability
- 2) generate new data / text



Respond

500 word report → ben AI



Memory

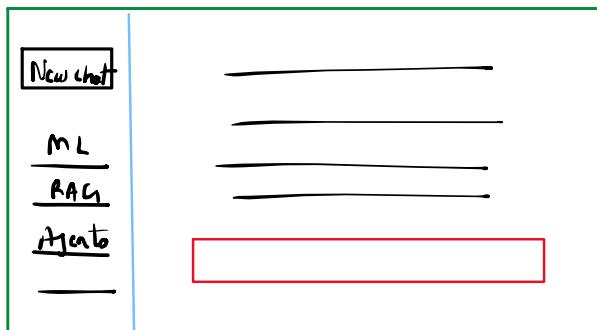
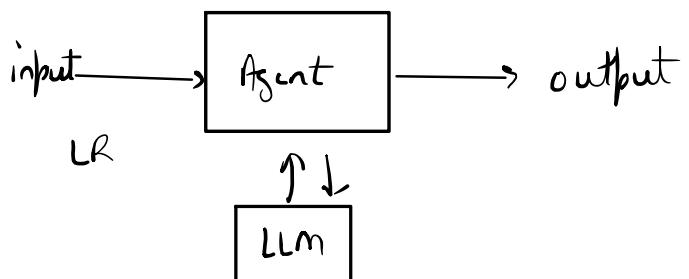
Memory → chat history / conversational history

Agent → Session ← Session-id

100 word summary

topic ?

Agents → Per execution basis



Separate conversations

↓
Session

unique
Session id

A single user can hold multiple conversations.

Session → Chats continue, pause, read. → memory
ML

Memory



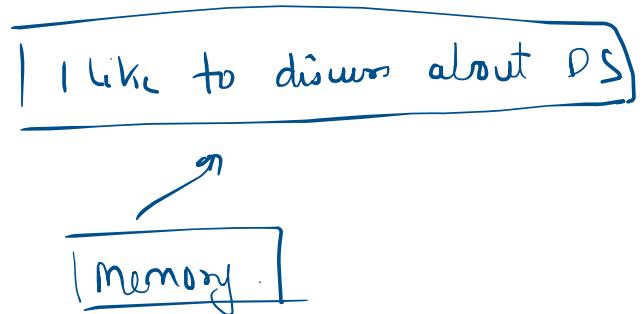
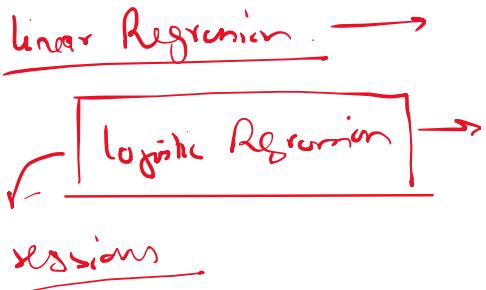
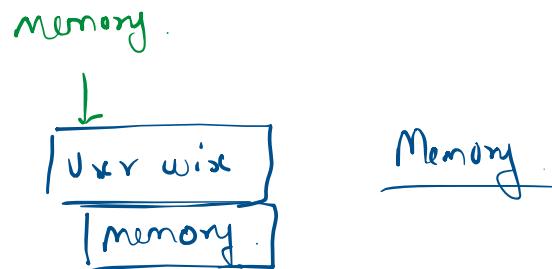
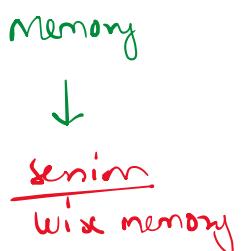
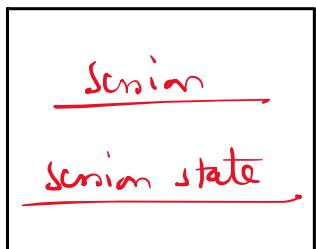
Short term

Memory

long term

Memory

Session



1) conversation / chat history → session .

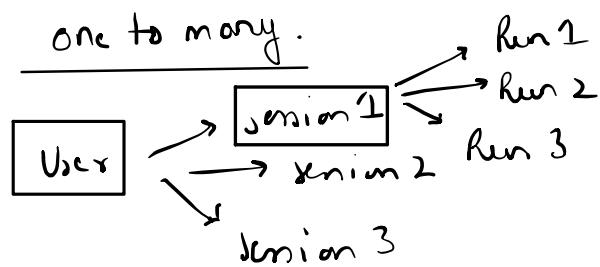
a) session related extra info → session state

3) User related memory → Memory / User memory

User & Session

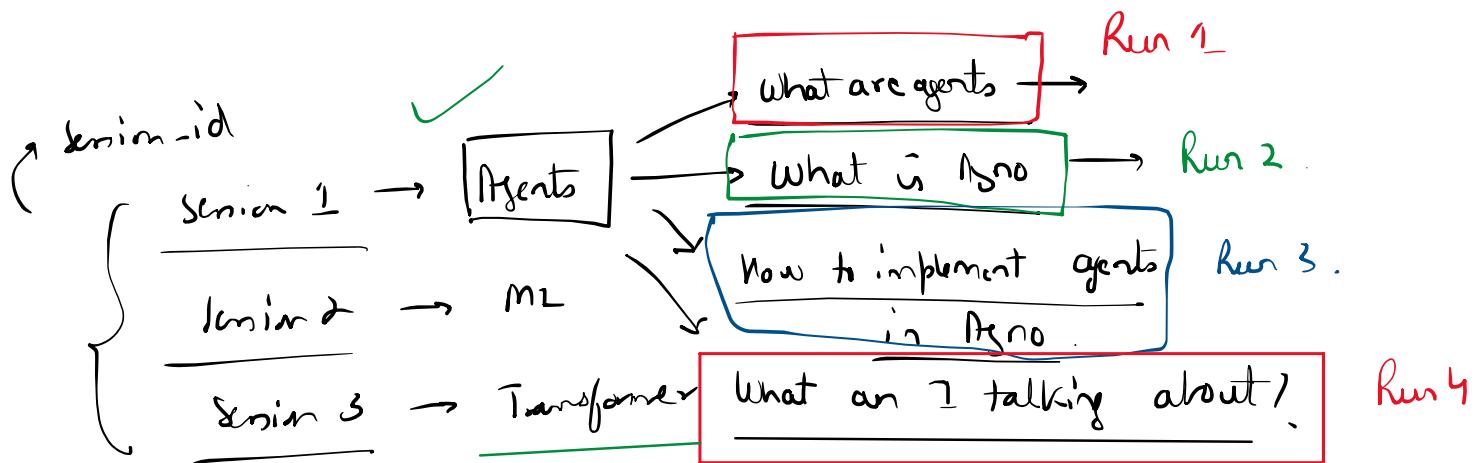
User → User

Session → one single conversation



Runs → Execution
of the agent

Input → Agent → Output



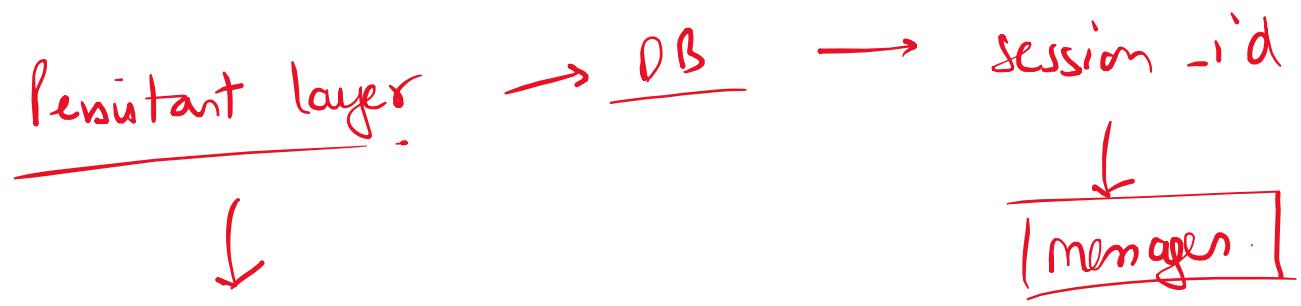
users → unique_id.

↓
sessions → session_id.

↓
runs → run_id.

session_id.

Storage



Storage → Backend

Memory → Frontend

Session State

→ dict Key : value

Working memory

Senior state = { username : Minanshu,
age : 34, } ✓

Information manage across runs → to do list,
shopping list,

The diagram illustrates the relationships between various database operations:

- RUD** (Read, Update, Delete) is at the top, connected by a horizontal arrow pointing right to **delete**.
- Add items** is on the left, connected by a vertical arrow pointing down to **Remove items**.
- Remove items** is on the left, connected by a curved arrow pointing right to **Update**.
- Update** is positioned below **Remove items**, with an arrow pointing from it to **delete**.
- delete** is on the right, receiving arrows from both **RUD** and **Update**.
- A checkmark symbol (\checkmark) is placed near the **Add items** label.

Agno

What is my name?

Senior state

↓

Name

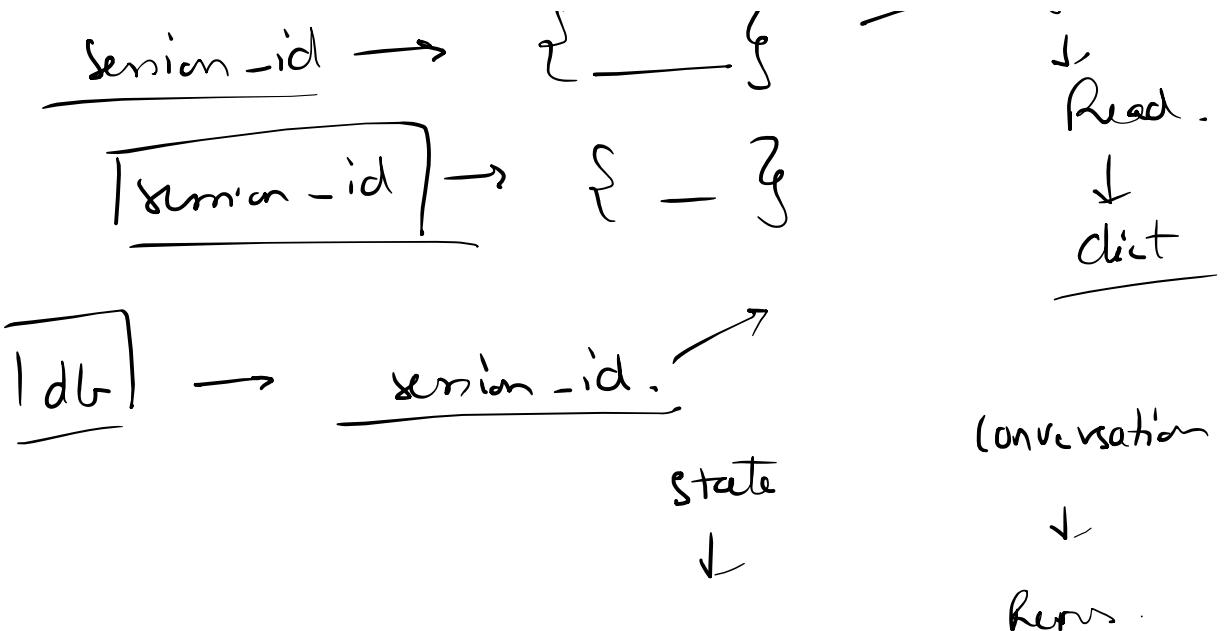
↓

Himanshu

Add task → Record Video

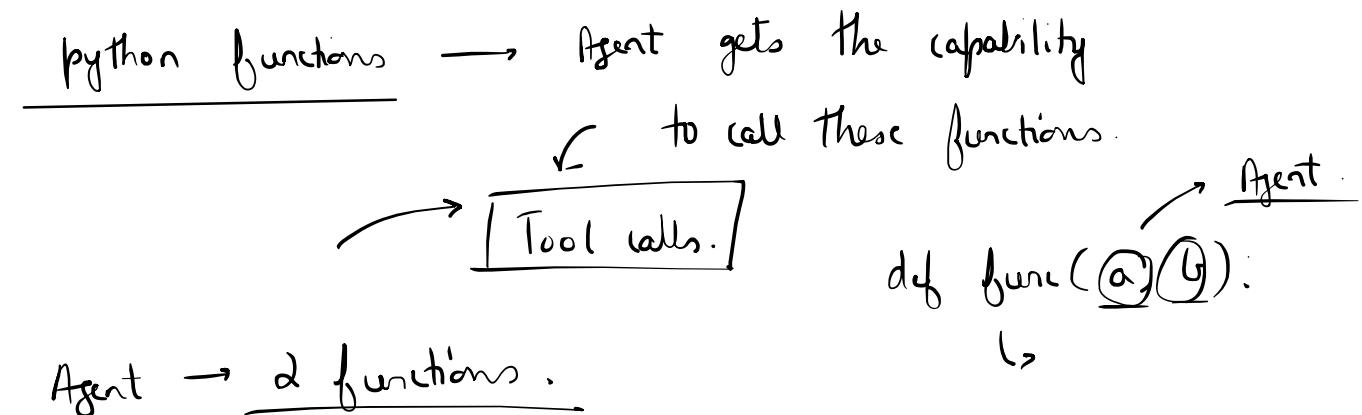
Recording completed. All done.

The diagram illustrates a mapping from a session identifier (session_id) to a collection of agents. The session_id is shown as an input arrow pointing to a brace that encloses a horizontal line. This line then points to an output arrow labeled "agent". Below the brace, there is a downward-pointing arrow and the word "agents", indicating that the session_id maps to a set of agents.



Tools

To do actions / perform tasks.

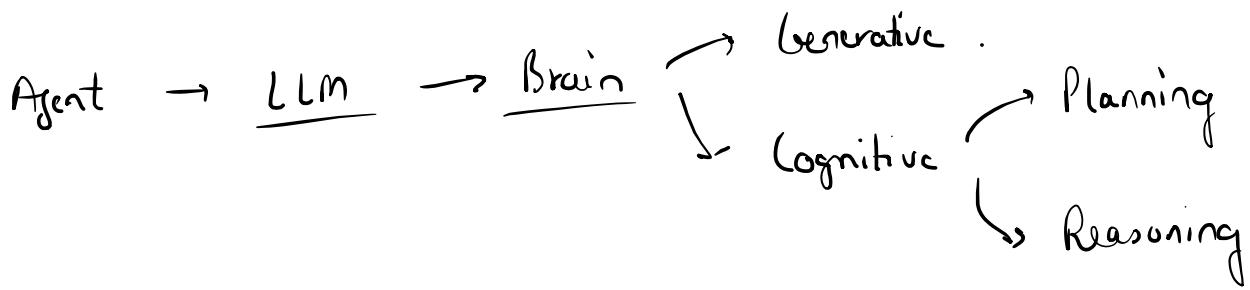


1) Agent decides on which tool to call based on
The user's query.

2) What input values / arguments to pass to the tool.
Argument type.

def add_numbers (a, b):
 return a+b

"3", "4"
add_numbers (3,4)



Planning → Complex goal → Sub tasks

Reasoning → Sub task → achieve

for a given query -

✓) which tool to call?

✓ 2) If tool is selected, what should be its inputs?

tools → docstring → LLM → tools access
tools → task

type hinting

```
def add_item (semi-n-state: dict,  
                item: str)  $\rightarrow$  str
```

@ tool → langchain

"This tool adds items to our list"

Tools = [tools'1 , tool2 , tool3] → custom tools
↓ python functions

Instructions = You have access to 3 tools.

↓

1) Tool 1 →

System

instructions

3) $\overline{100} \mid 2 \rightarrow$

2) Tool 2 →

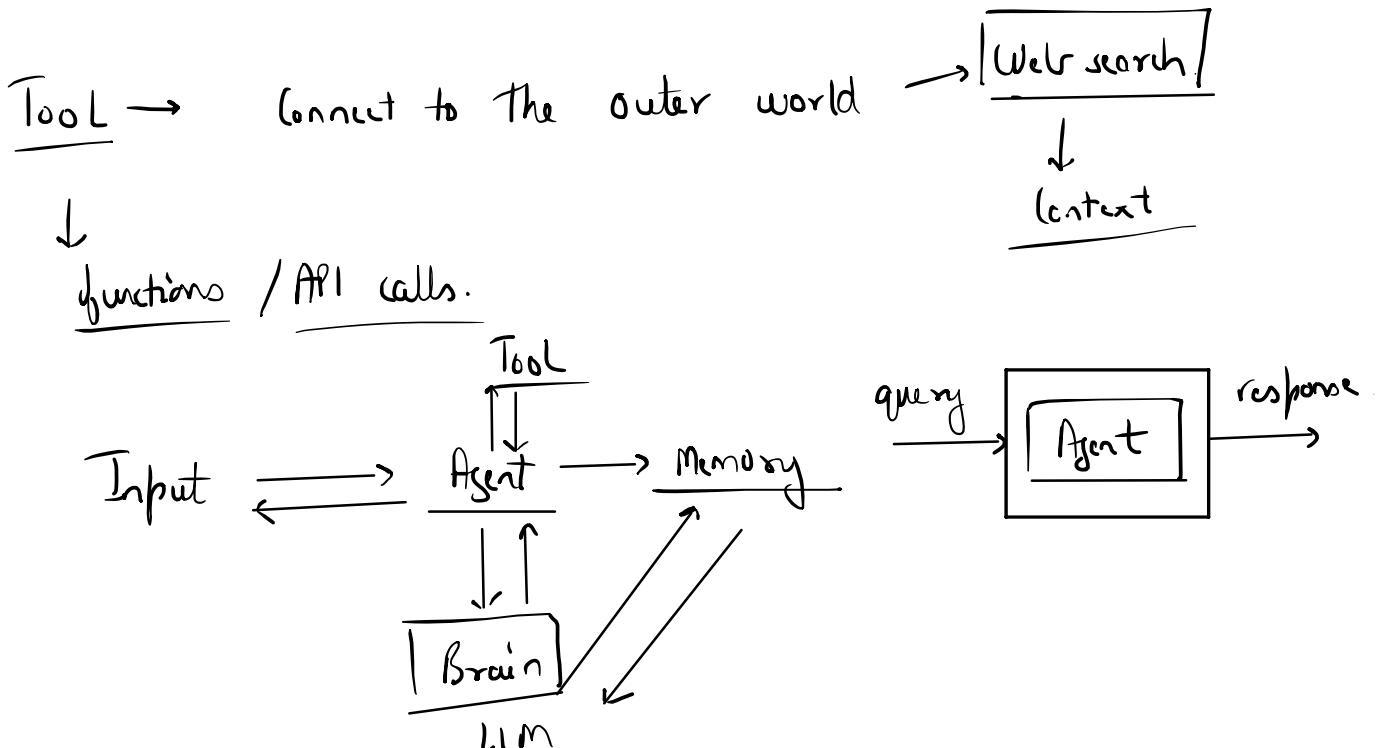
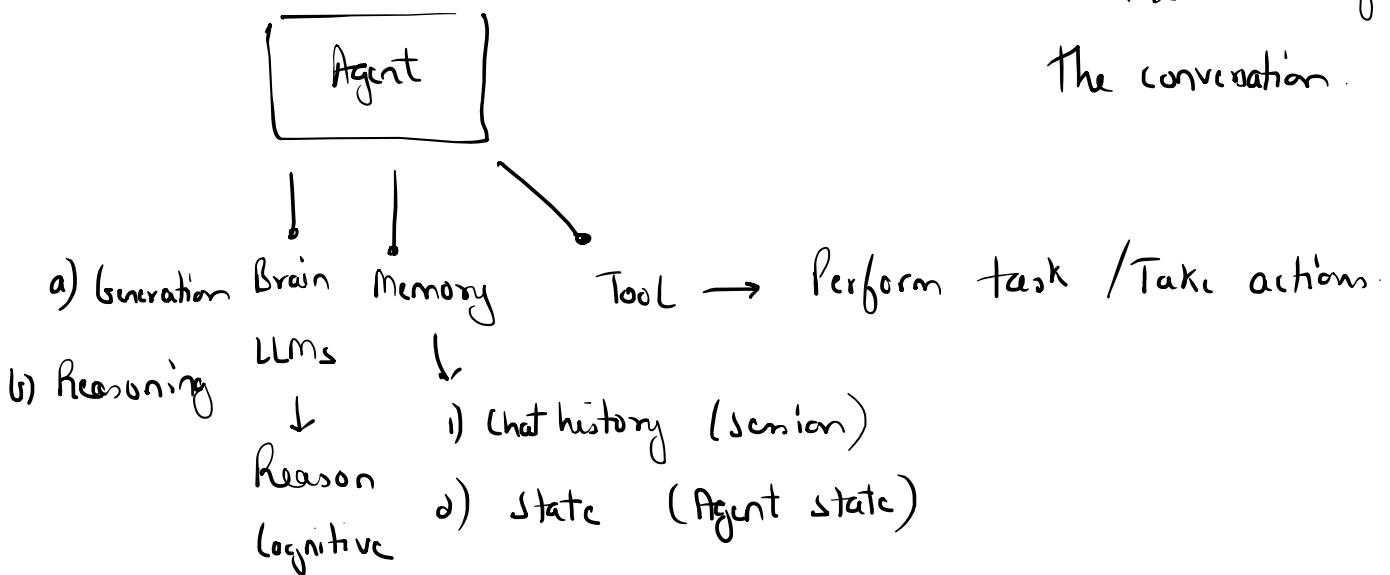
3) $T_{00}(z) \rightarrow$

knowledge → Agentic RAG

Components of Agent

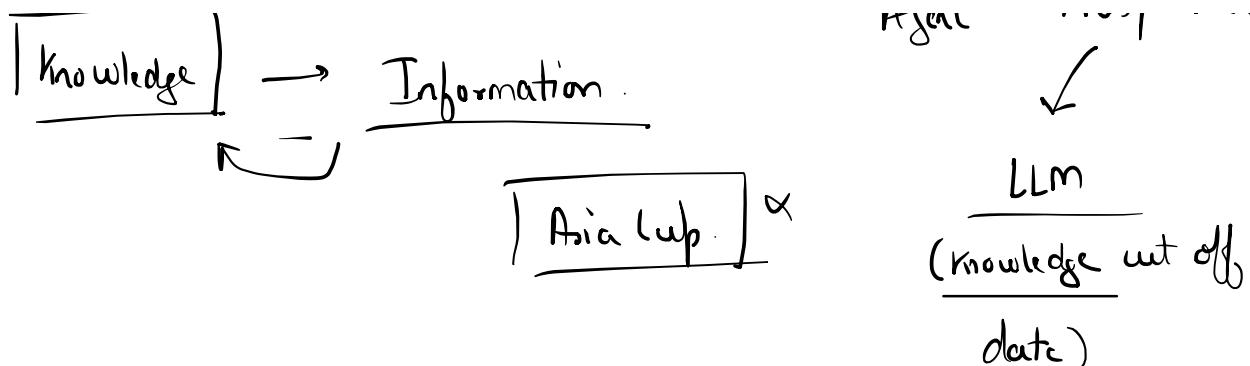
Memory → conversational history

Maintain the context of
the conversation.



knowledge → Information

Agent → Response.
↙



I don't have info about this event.

Hallucinate → $\boxed{\text{Asia Cup 2025}}$ was won by $\boxed{\text{Australia}}$

Agent → Superpower Informational $\boxed{\text{superpower}}$
 ↓
 $\boxed{\text{LTM}}$ → Recent events Up to date info
 ↴ knowledge $\boxed{\text{Knowledge}}$

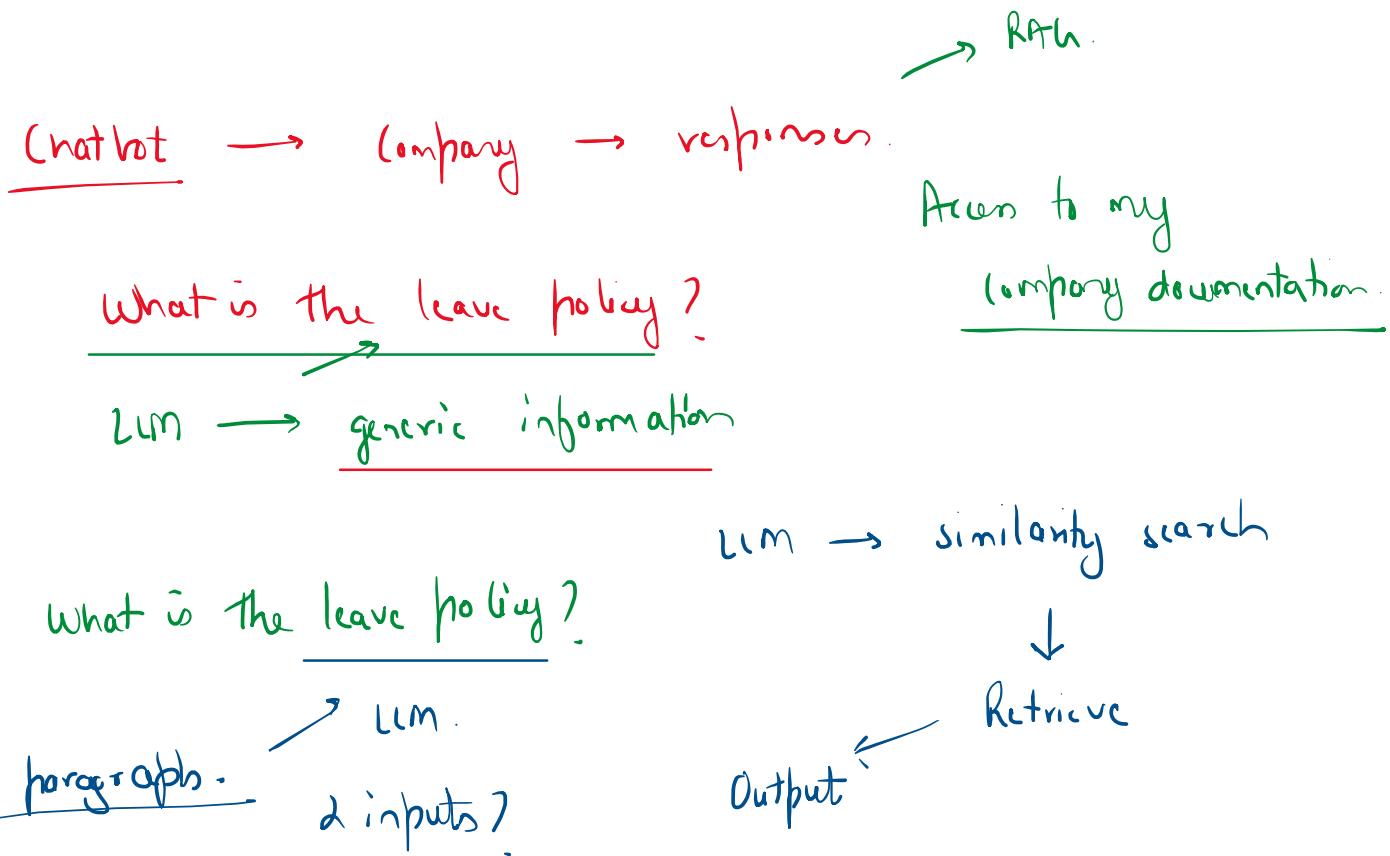
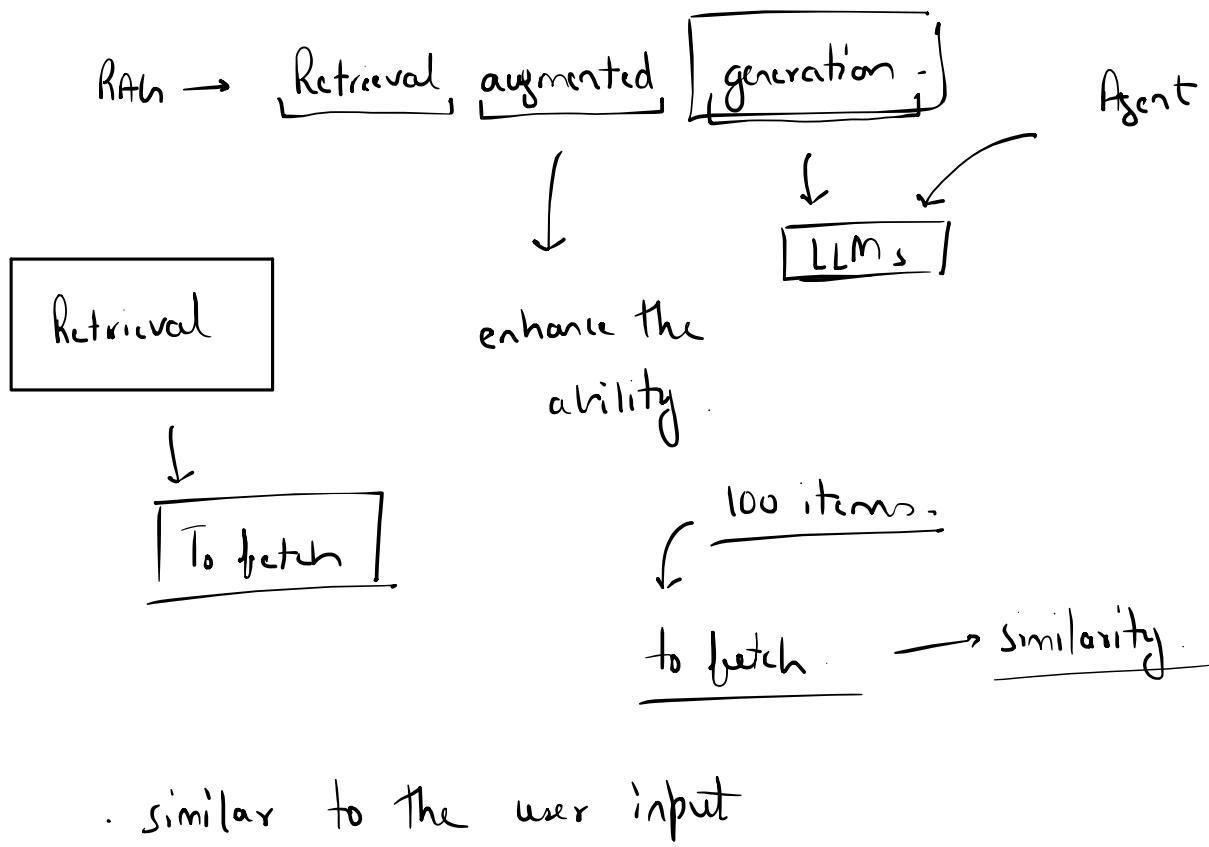
Customer Support + Agent → $\boxed{\text{Info}}$

Status

Knowledge → $\boxed{\text{RAG}}$ → Retrieval Augmented Generation
Agentic RAG

What is RAG

Knowledge → Agentic RAG.

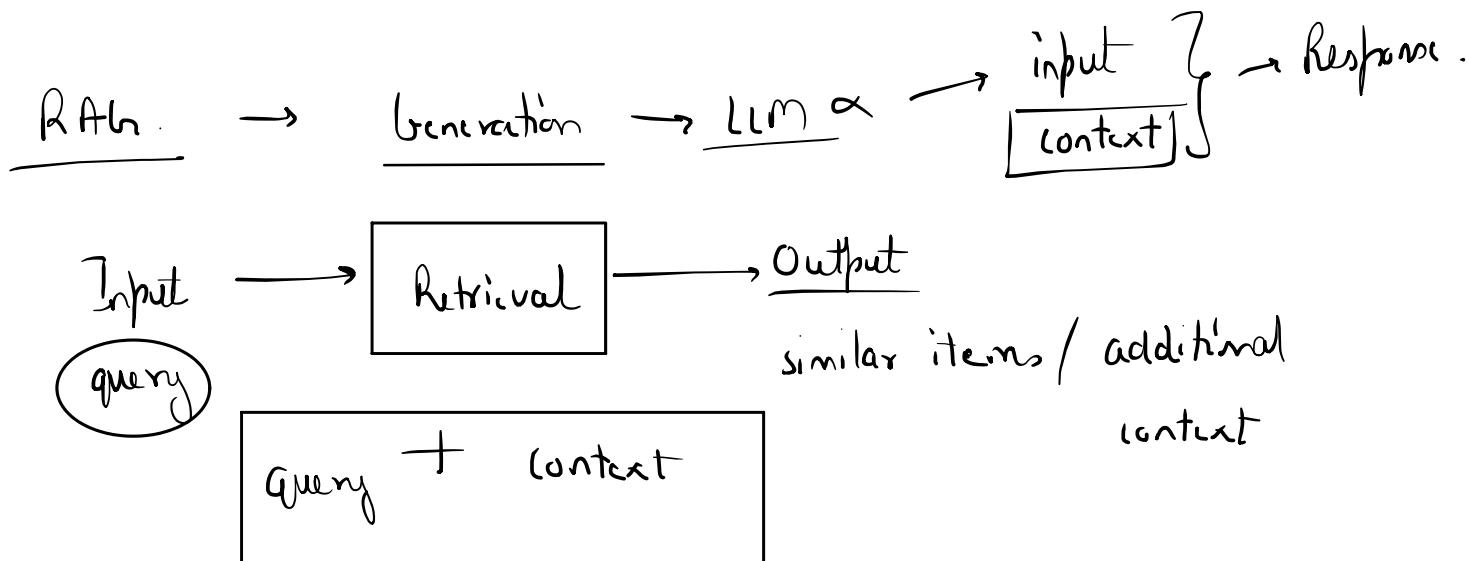


input query: What is the leave policy?

additional context: 3 paragraphs

response → The company policy regarding leaves is -
↓ ,)
specific for my organisation.

Retrieval, augmented generation, → RAGh.
↓
more specific.



RAG components: (Work in sequential order)

↓
Pipeline

1) Readers / Loaders

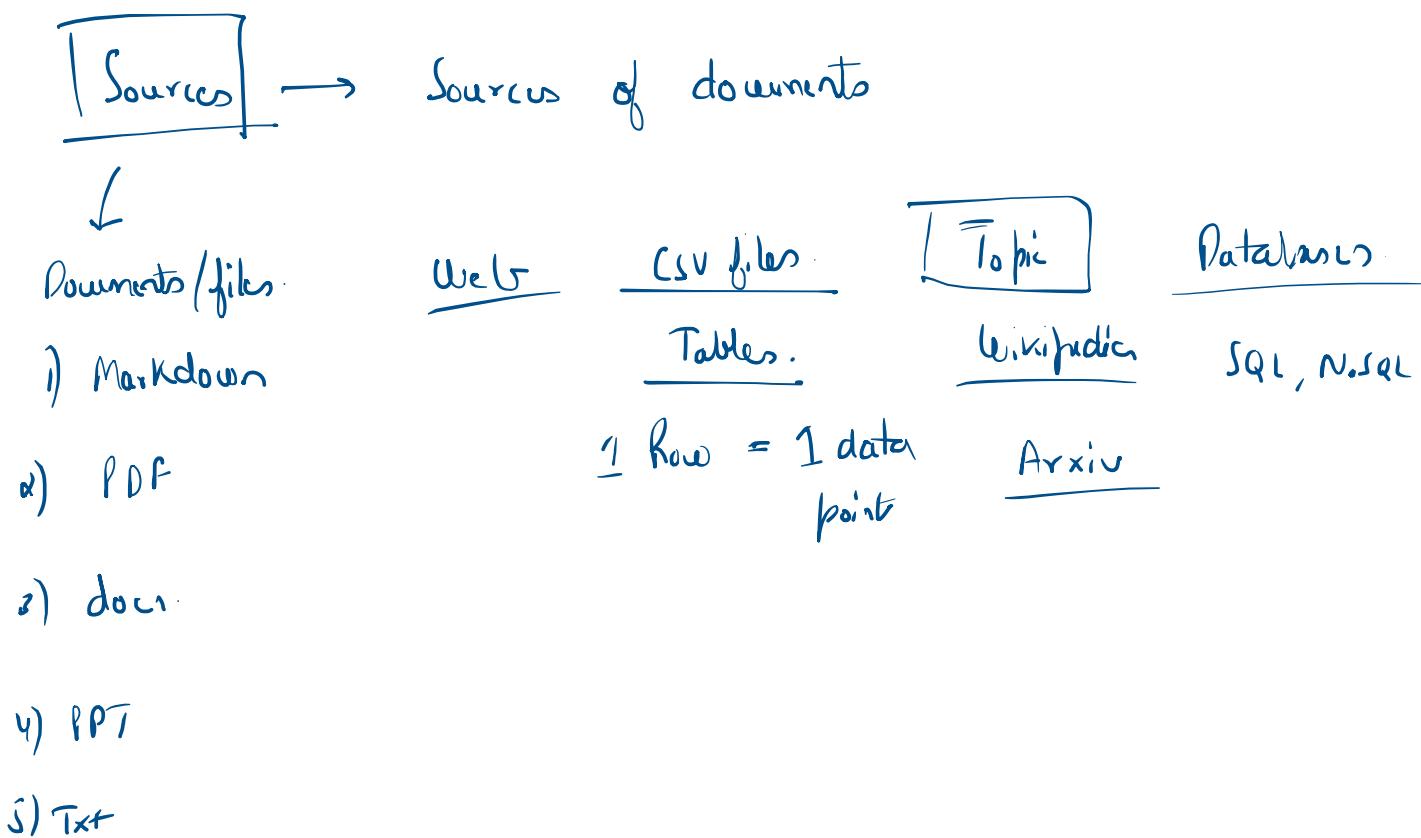
4) Storage.

2) chunking

5) Retrieval.

3) Embedding

1. Readers/Loaders and Sources



Readers / Loaders → load. the document and
parse it.

Hard drive

pdf

PDF Reader

show the document
in paper format

Format

Text, images,
graphs, tables

2. Chunking

Chunking → To break into chunks.

500 pages .

Page wise .

Chapter wise

Agno →

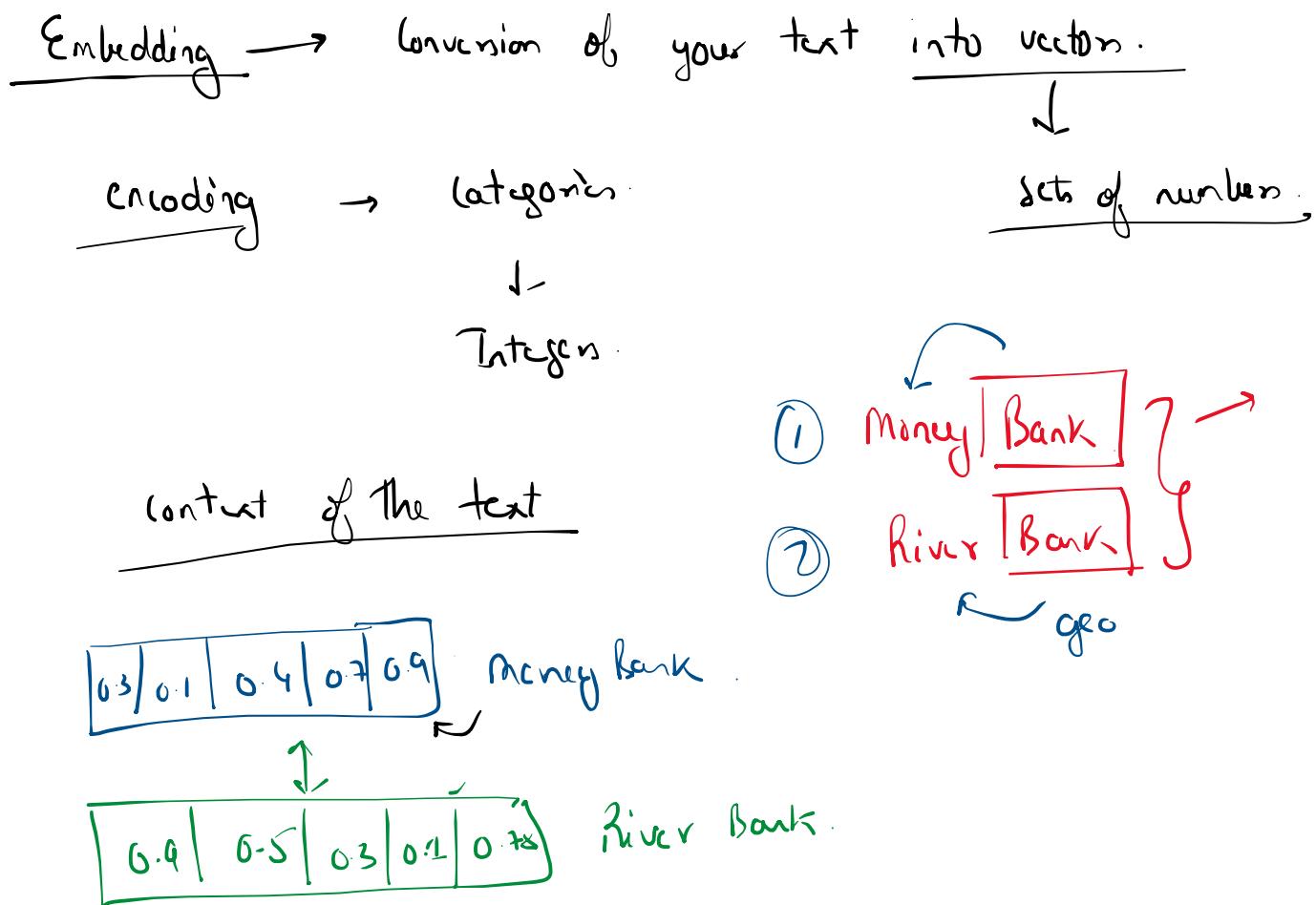
1) fixed chunking → 2000 words .

2) semantic chunking → similar text .

3) Document chunking → Document format

4) recursive text chunking →

3. Embedding



- ① Document load .
② broke the document into
multiple sub documents

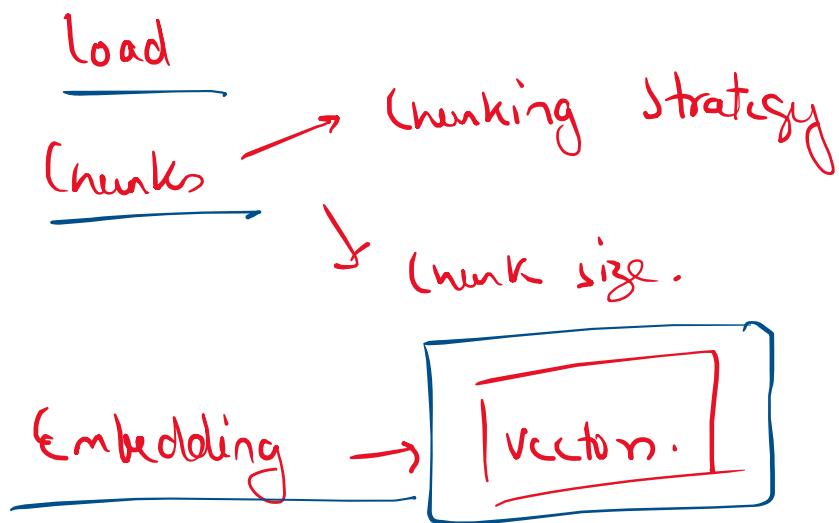
③ Numbers convert → LLM → embedding
→ coherence
→ chunk size ↑ → context → numbers.
large context → ↓ numbers.

lunch size ↓ → number.

↓ in detail →

↓ number of vectors.

4. Storage

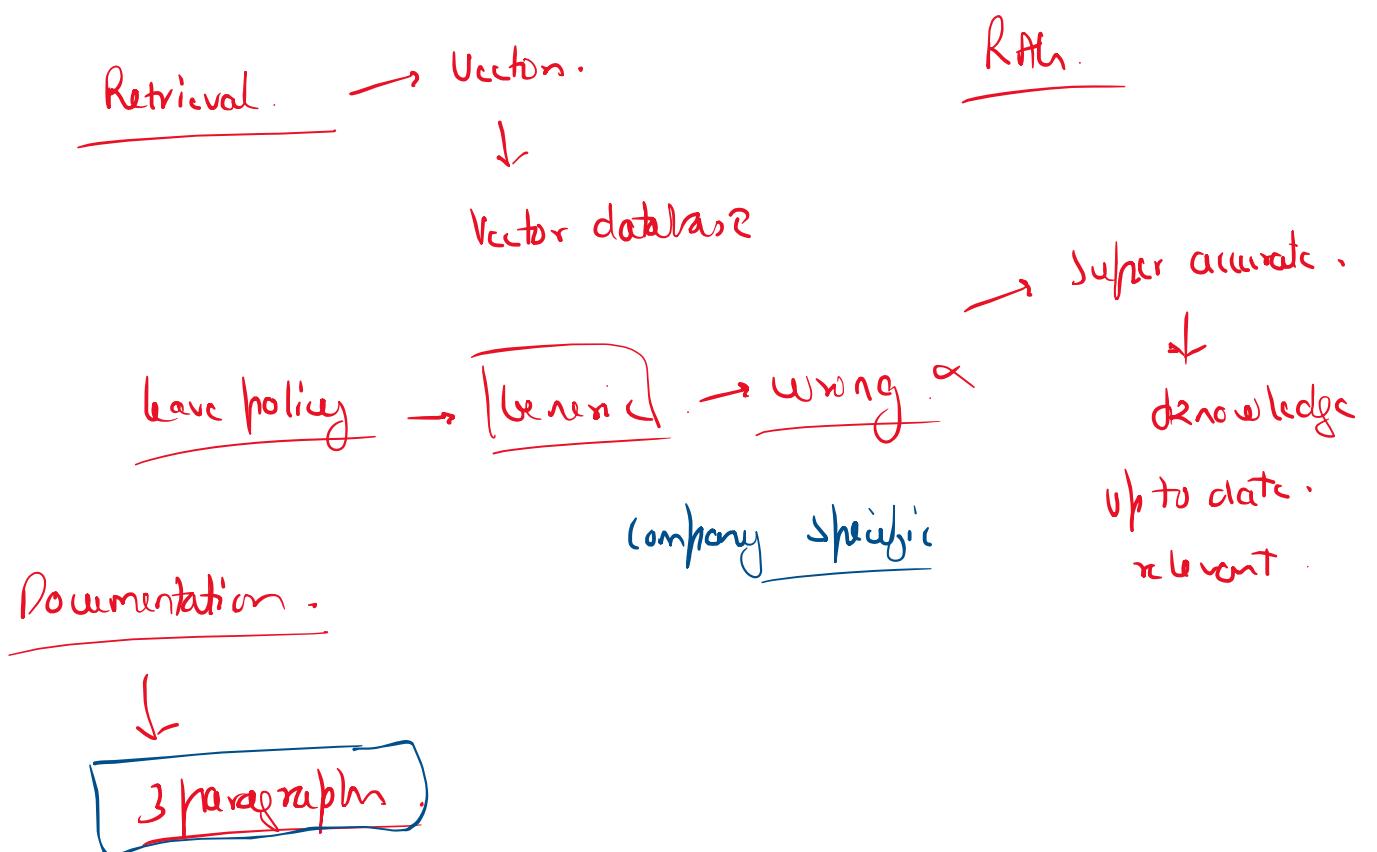


Vector databases → embedding vectors stores

vector · metadata → Text

indexing

5. Retrieval



Retrieval pipeline:-

- Input
- Output
- Processing

Input → What is the leave policy of my

company!



embedding component → numbers.

input vector

0.1 | 0.3 | 0.2 | 0.9 | 0.6

5 numbers

similarity search

→ cosine similarity

Semantic context

Top 3 matches → return. → 3 paragraphs

similar to my input

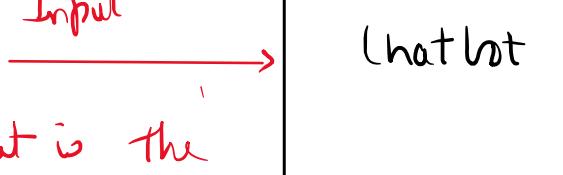
leave policy

RAG → retrieval.

Retrieval →

↓ 3 paragraphs → leave policy.

Input



Output

Specific for
my company.

query:

what is the
leave policy of
my company?

+

additional
context

+

3 paragraphs

(RAG)

generate the
response

Augmentation

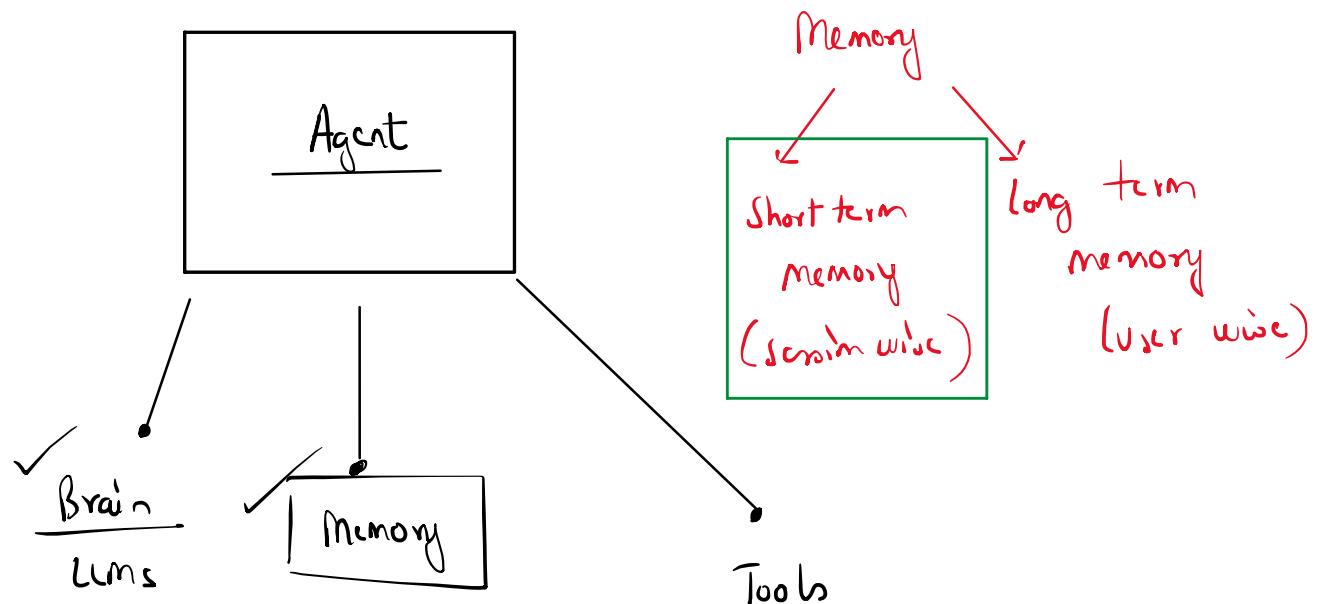
What is Agentic RAG

Knowledge Pipeline in Agno

Code

Recap

✓ Knowledge → Superpower for your agents



- ✓
1) Generate 1) Context Perform tasks / Do actions
2) Think/Reason 2) Retain context functions.
(cognitive abilities)

Which tool to call?
Which inputs to give to
a tool?

Memory → short term memory

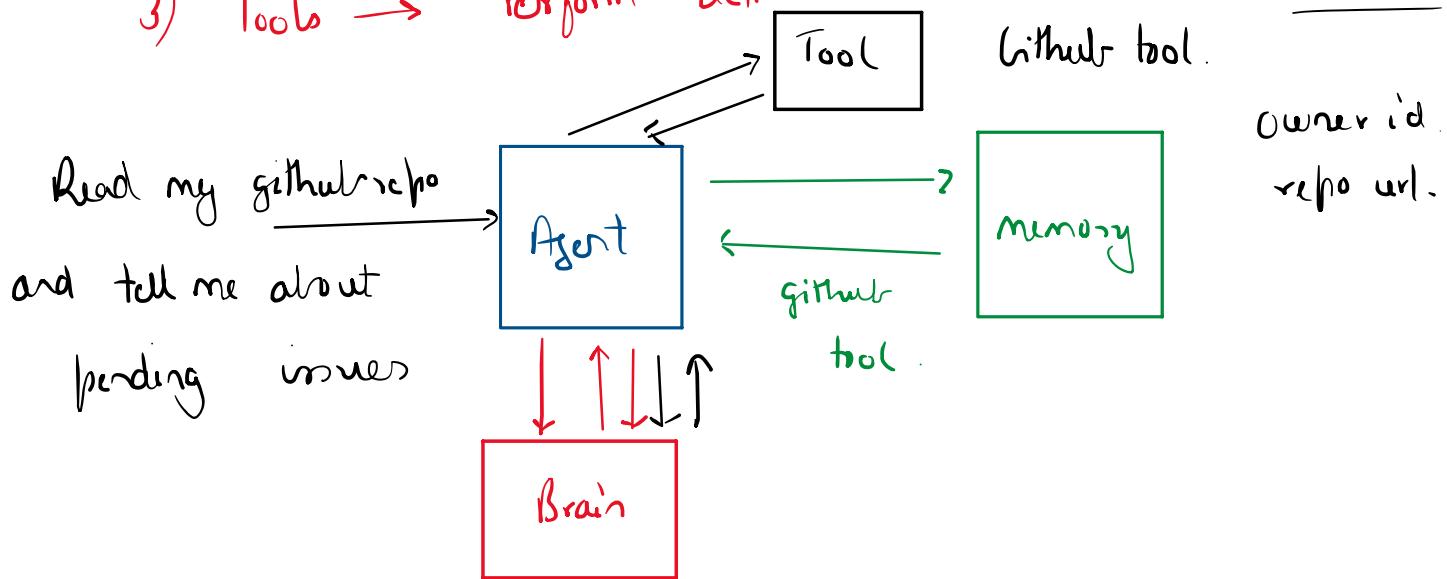
- 1) Conversational history / Chat history → Scratches
2) Agent state / Session state → Working memory that can
be accessed across your runs.

{ shopping list : [] } → CRUD

Agent → 3 components
... run ... eval ...

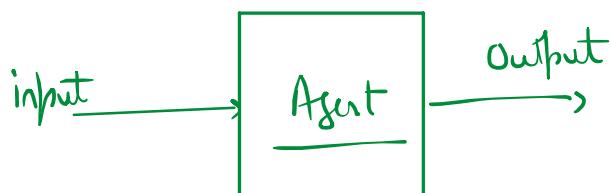
Agent

- 1) Brain → Cognitive / thinking ability
- 2) Memory → retains context / meaning across the entire conversation
- 3) Tools → Perform actions.



These are some of the pending issues for the repo.

- 1) Issue name Issue Id → Description



4th component :-

knowledge → Information / Data. → Relevant info

Agent → LLM Knowledge out of date.

Dec 2024

Who won Asia Cup 2025 ?

→ I don't have info
about this event



M hallucinate → Was won by
Bangladesh

Proper information

Relevant Info and up to date information -



Knowledge →

R A t h

Retrieval augmented generation .

Agentic R A t h

Introduction

Flow :-

RAn

Diffr components in RAn

Agentic RAn

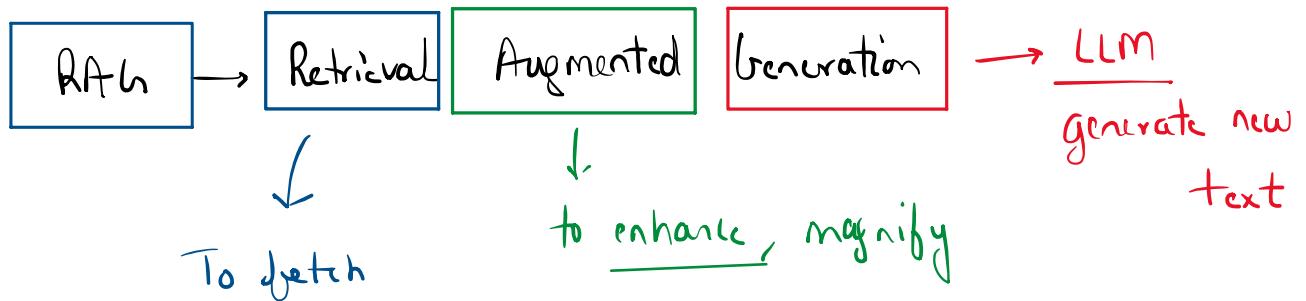
Diffr b/w Agentic RAn vs Traditional RAn

knowledge → Agentic RAn



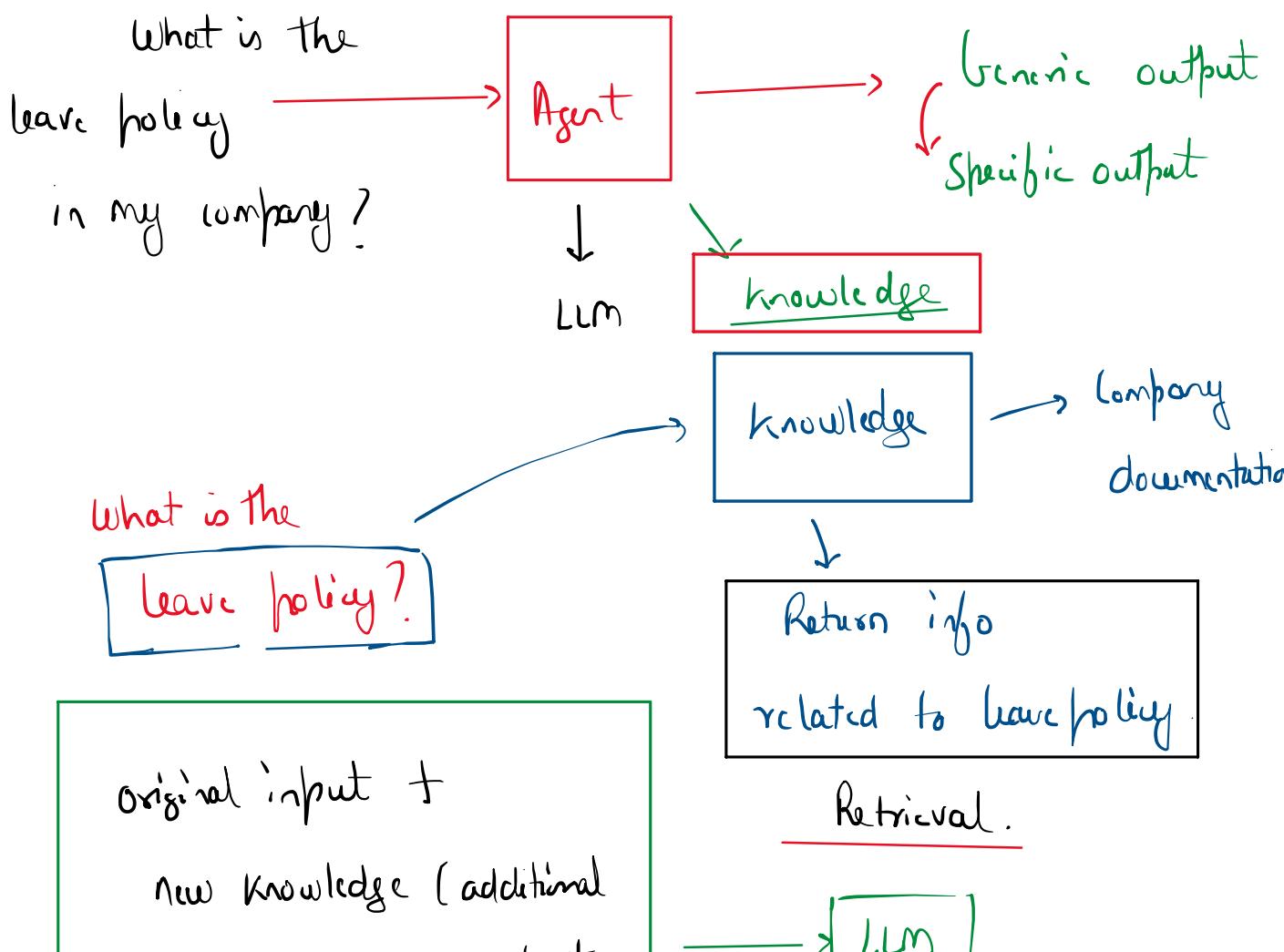
Code .

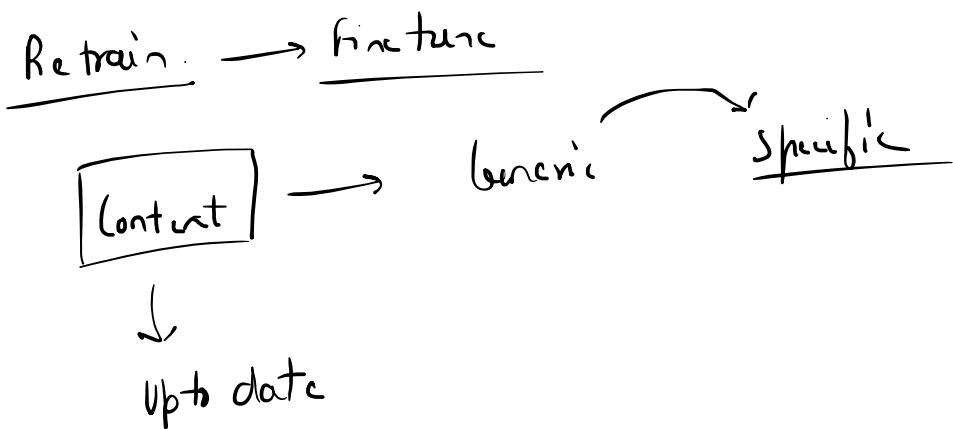
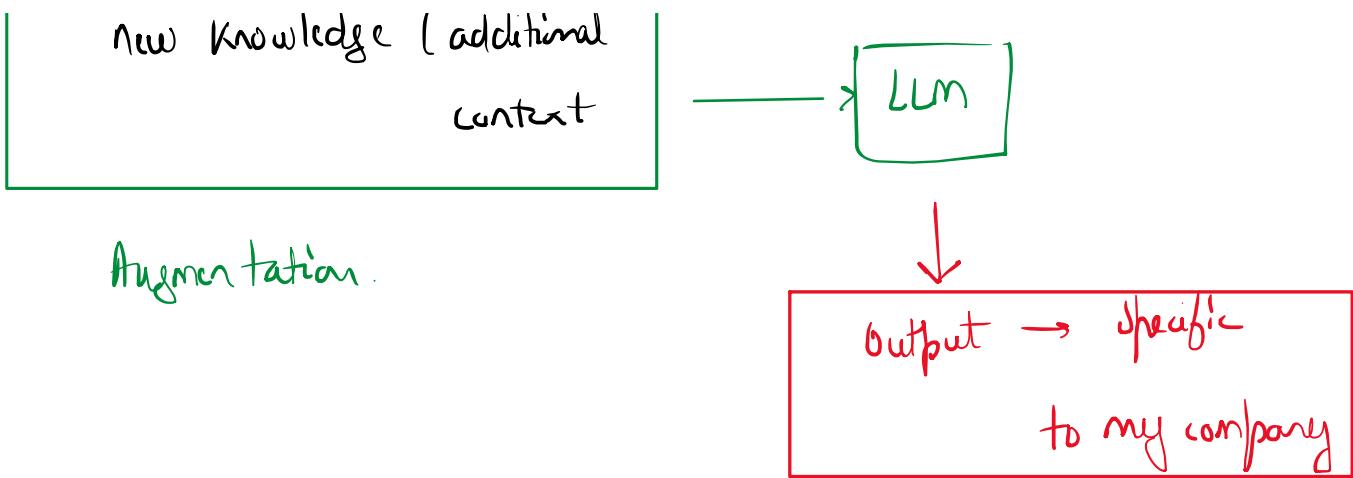
What is RAG?



fetch relevant information

more specific from more generic. RAG & parameters





5 main components of RAG:-

- 1) loaders / Readers
- 2) chunking strategy
- 3) embeddings
- 4) knowledge base / storage
- 5) Retrieval

1. Loaders/Readers

Components → Which help in loading the document.

Knowledge / content sources:-

- 1) file → Markdown, doc, txt, pdf, ppt
- 2) Web → Web page
- 3) Database → SQL, NoSQL
- 4) topic → Wikipedia, Arxiv

Loaders

Readers

Agno

.pdf

→ S3

Hard drive.

- 1) Load the document into memory.
- 2) Parse the document → text, images, tables, links
- 3) Convert it into document type. PDF.
↓ ↓ ↓
type content metadata
 { no. of pages = __,
 creation = __,

saved on = — /

modified : — {

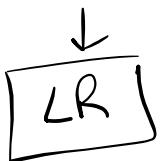
2. Chunking

We have the loaded document.

→ PDF → read . extract

(chapter 2)

Document → Very large → Into chunks.
logical.



Chapter wise → Topic wise.

↓
20 pages

Paragraph wise , Sentence wise .

20 times

Some of the chunking strategies:-

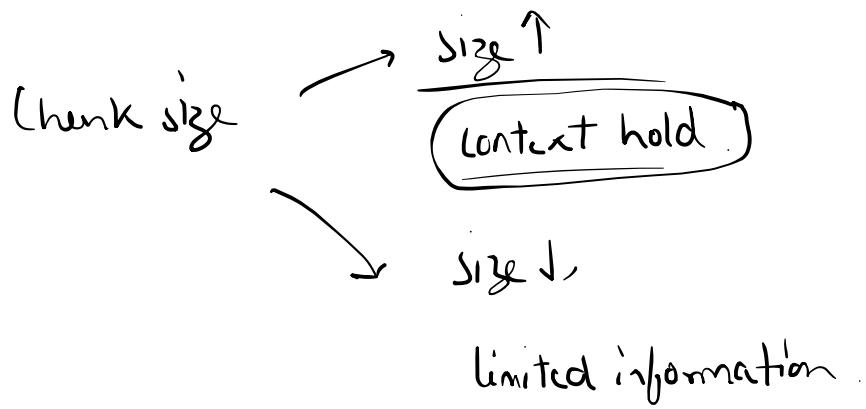
- 1) fixed type → 100 words.
- 2) semantic chunking → semantic meaning similar.
max chunk size .
- 3) Document chunking → well formatted .
- 4) Recursive chunking →

2 selections -

- 1) chunking strategy → Experimentation.

1) chunking strategy → Experimentation

2) chunk size



3. Embedding

Load, chunks

Into numbers → Text α

Categorical

ONE



encode → 0 1 2 3

embedding

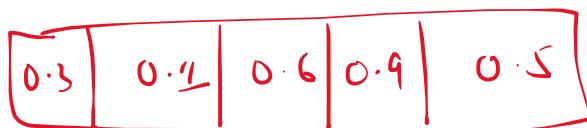
→ Semantic meaning of the text

dynamic → context

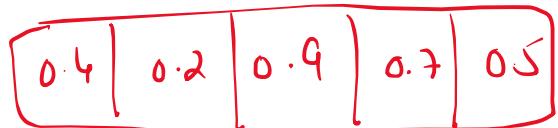
→ intuition.

→ location.

Money Bank



River Bank



Semantic meaning

dynamic contextual

embeddings

Load 500 page

→ Page wise chunk

↓ 500 documents

embeddings 500

chunk size → 1000 words → context

chunk size → 1000 words → —
↓
S numbers

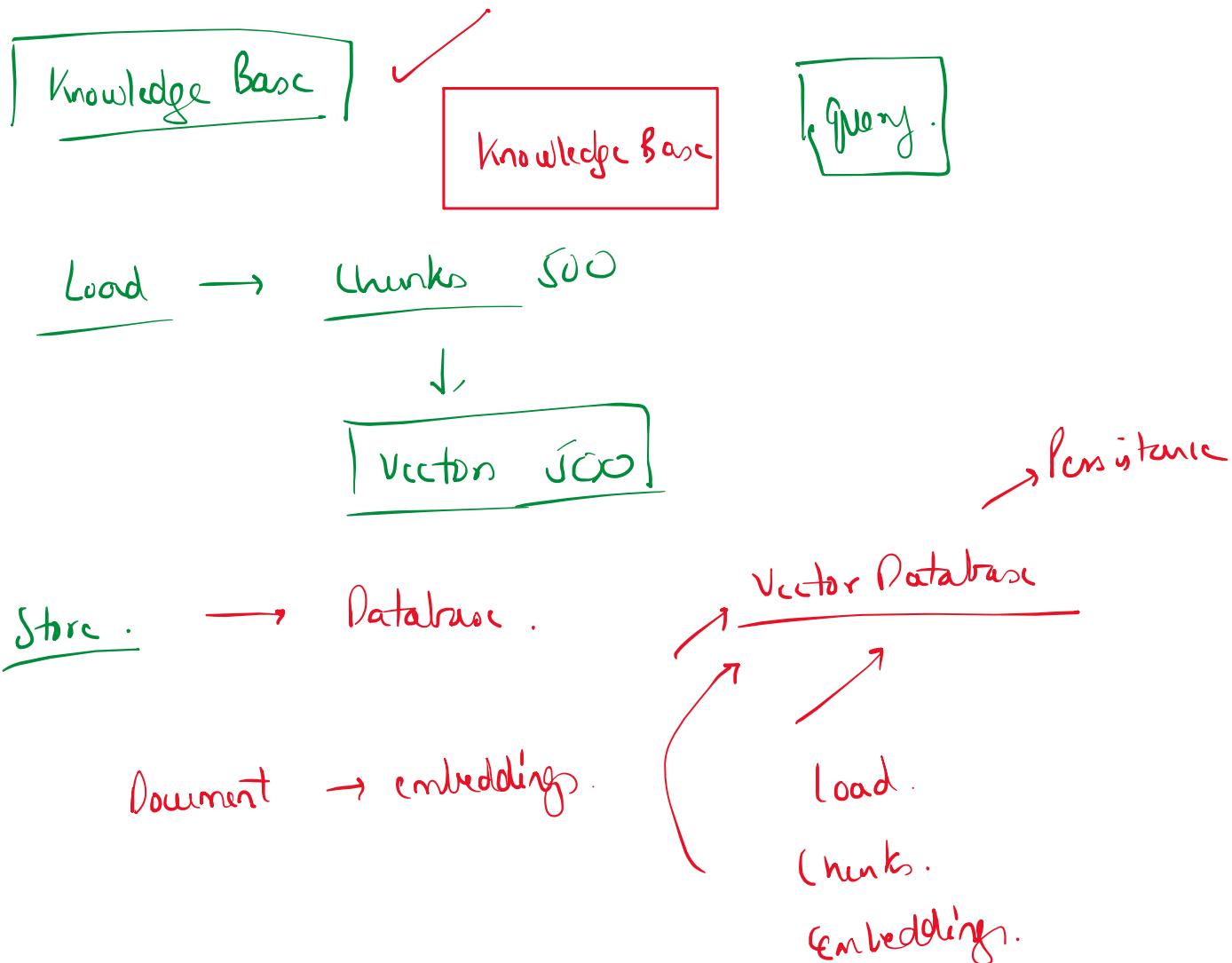
200 words → context

↓
information

loss α.

50 words → context
↓
S numbers

4. Storage



5. Retrieval

Documentational \rightarrow Info

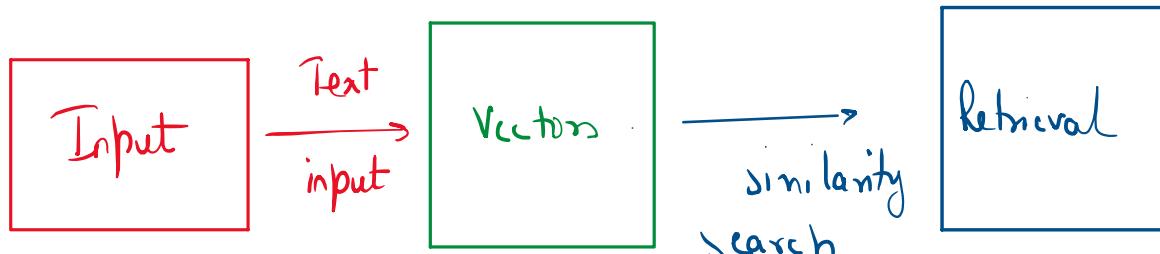
Very specific \rightarrow based on the user input

leave policy \rightarrow leave
Salary

similar to the
input query.

knowledge base . Similarity search.
Input \rightarrow Vector.

✓ embedding
model.
content



Similarity search \rightarrow cosine similarity. -1 to 1
vector distance

What is the leave policy of \rightarrow Embedding model.

my company

Context
and intent

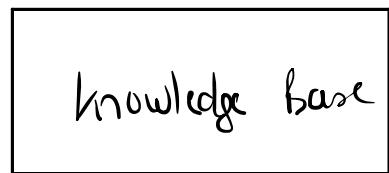
↓

Input vector

0.3	0.2	0.46	0.78	0.93
-----	-----	------	------	------

context
of my input

0.3	0.2	0.46	0.78	0.93
-----	-----	------	------	------



→ 500 vectors → compare
similarity search

500 → similarity → sort descending order
using similarity score.

Top 3, top 5 documents retrieve

top 5 vectors.

metadata

↓
text → Page

Page 50, 63, 128, 250,
367

Retrieval phase.

R&H : Retrieval Augmented generation
enhance

Prompt → Input

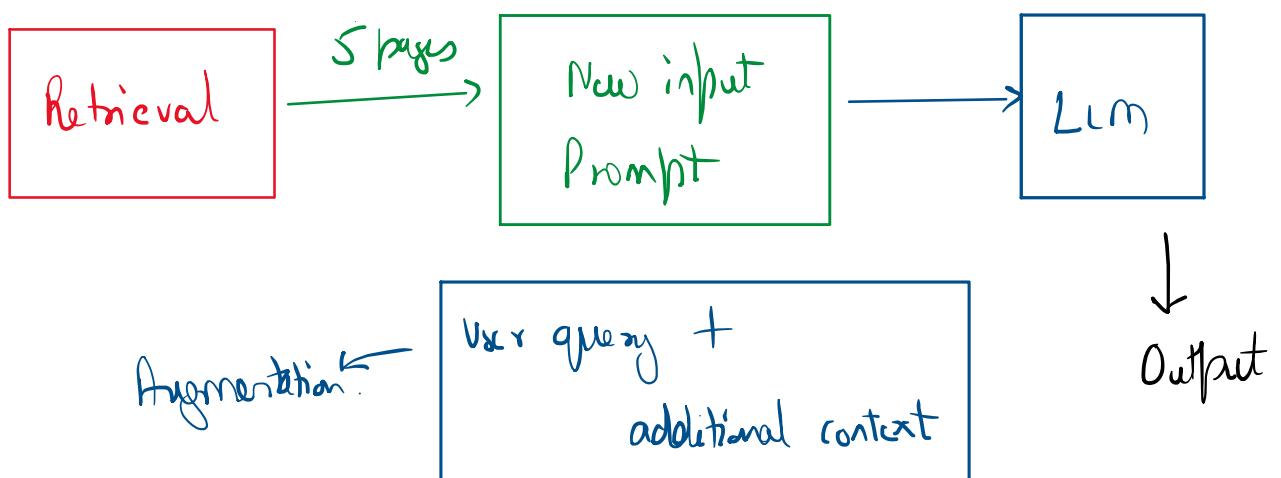
query : ,

Based on the input query and context,

Context :

answer the question. Make sure

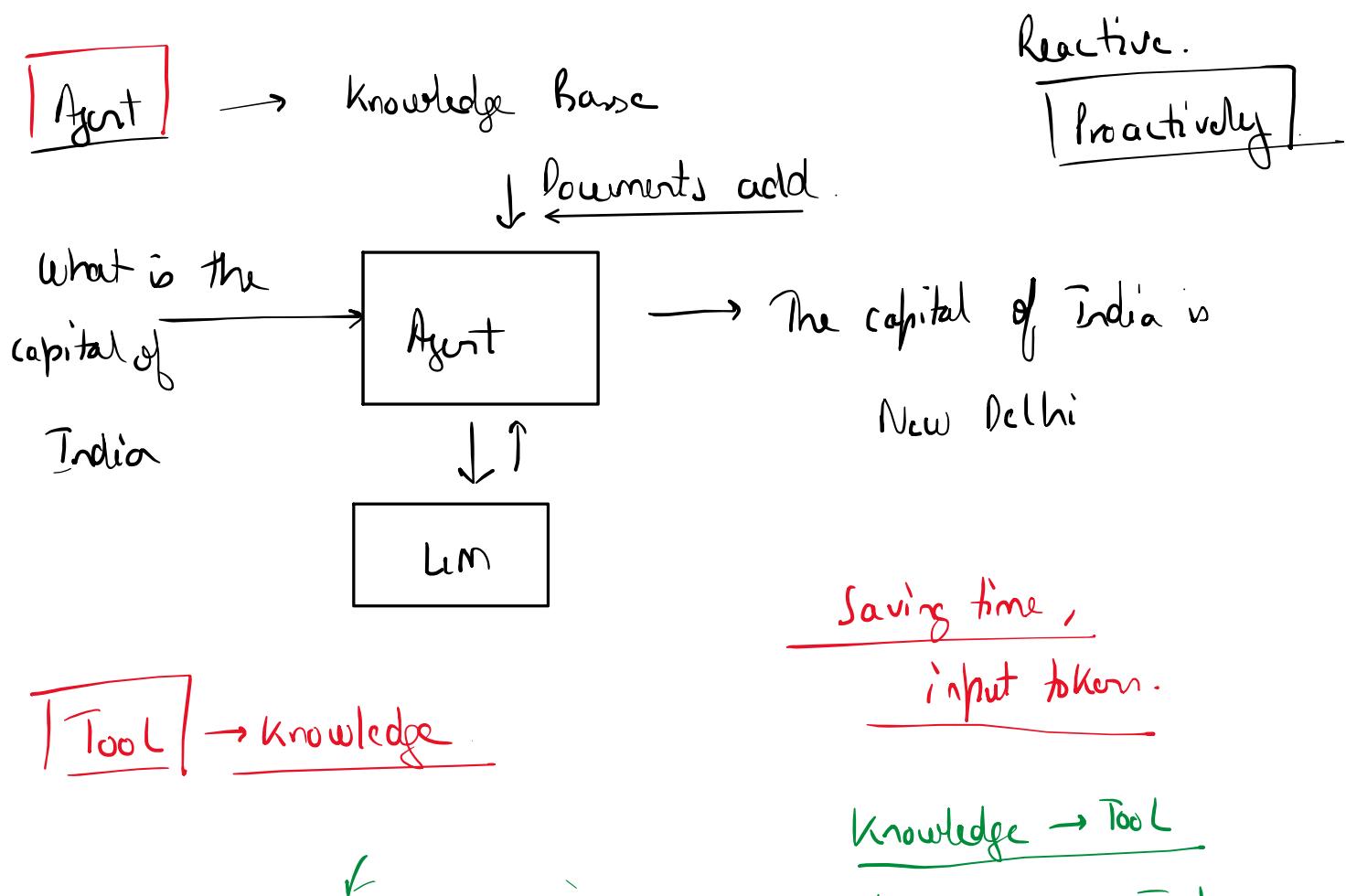
to answer from the context.



Agentic RAG → Retrieval
↓
Smart Retrieval.

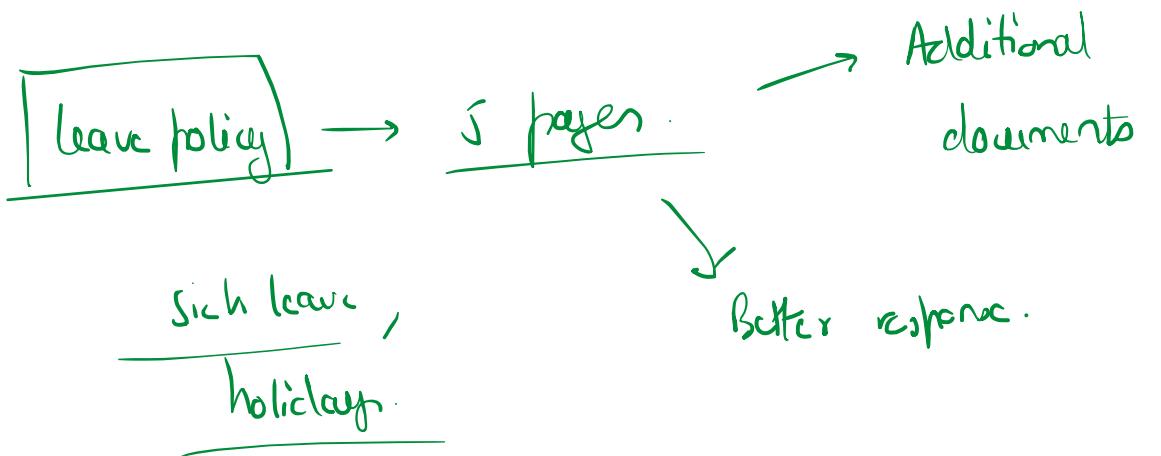
Agent has the capability.-

- 1) Agent autonomously decides when to search knowledge base.
- 2) Determines what to search.
- 3) Controls how to use the retrieved information.



knowledge base & Web

Knowledge → Tool
Web search → Tool.



Characteristics of Agentic RAI:-

- 1) Active decision making
- 2) Dynamic
- 3) Full control
- 4) Query Reformulation
- 5) Selective usage
- 6) Iterative search

Advantages :-

- 1) Efficiency →
- 2) Better Search Quality

→ knowledge base'

Compare b/w X and Y → knowledge Box 2

3) Iterative Refinement

4) Context aware decisions

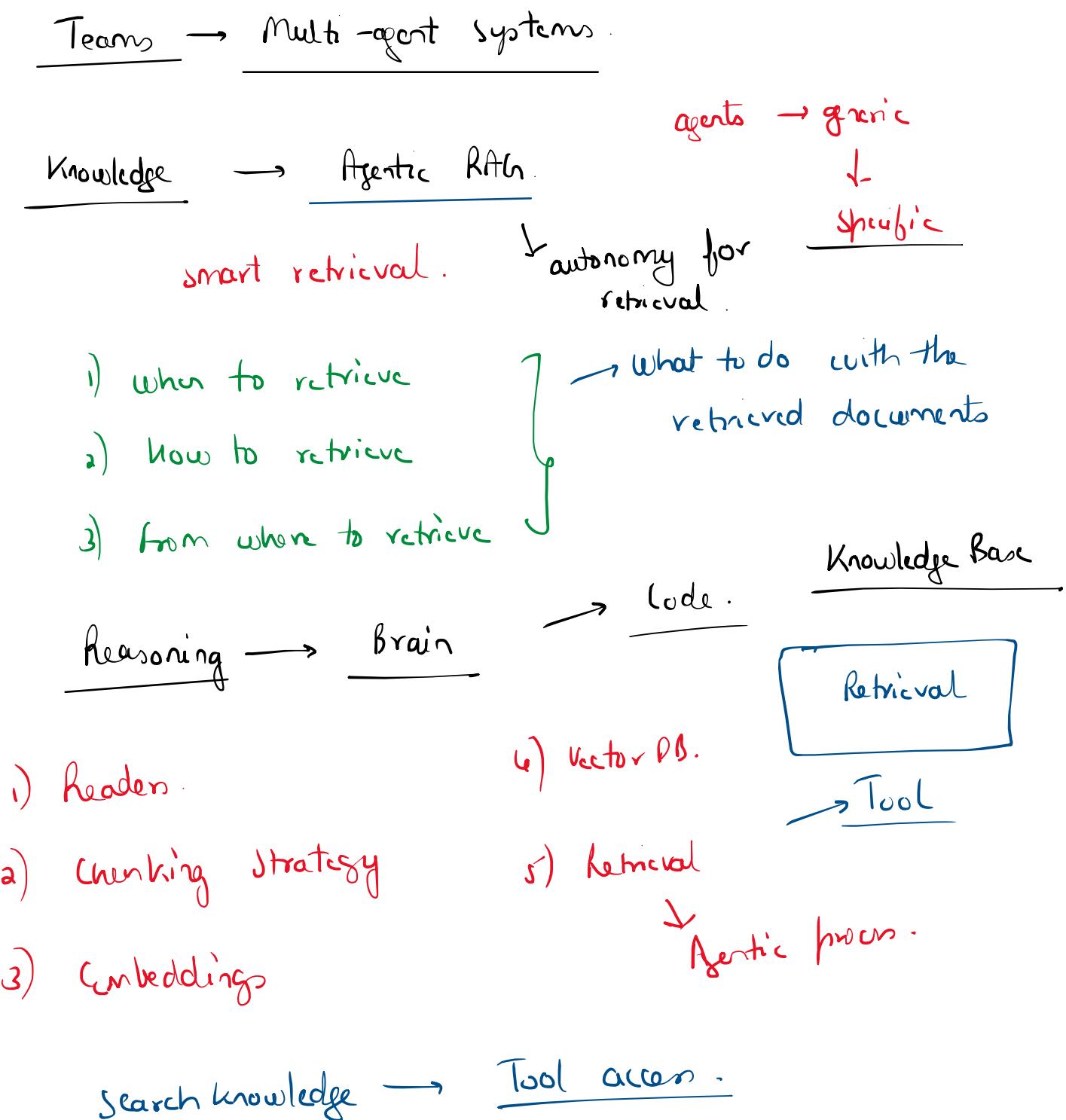
5) Selective information use

Diff between Traditional vs Agentic RAG

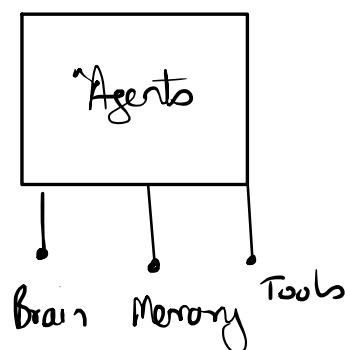
	Traditional RAG	Agentic RAG
Decision Making	Reactive	Proactive
Searching	Always	Agent decides
Query formulation	Direct	Reformulate for better retrieval.
Number of searches	Usually 1	More than 1
Result Processing	Use all of the retrieved info	Agent selects the relevant info
Efficiency	less efficient	more efficient
Adaptability	No fixed workflow	Dynamic workflow
Intelligence	Basic retrieval	Smart retrieval

Code

Recap



Introduction



Complex task
→ Brain → Planning
↓
Into multiple sub tasks



Reasoning → from the list of tools

↓
Selection of appropriate tool

Reasoning → forms the input for the tool

example - tool (name : str, age : int)

str

↓
Response

Suitable inputs to the tool

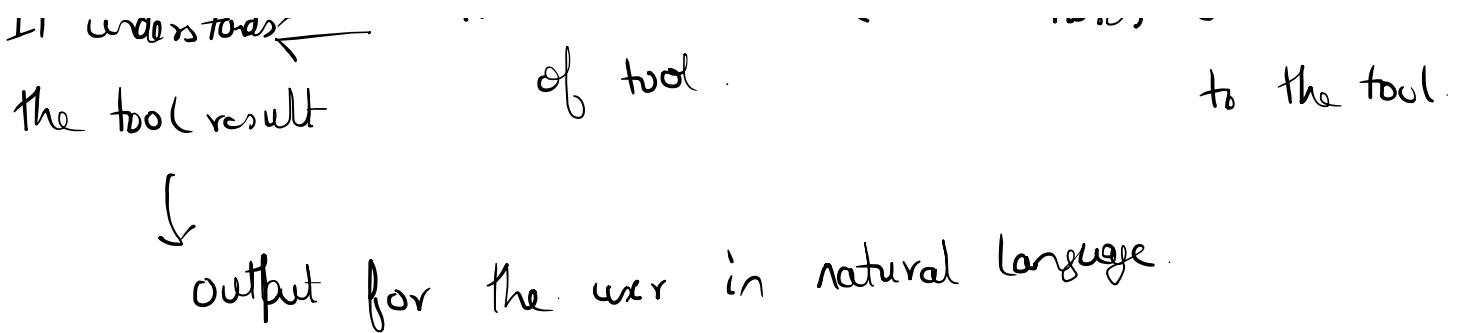
query → Natural language query

Convert to proper dtypes

It understands
the language

Tracks the execution
of tool

Parses on the inputs
to the tool



Planning phase → 25 subtasks.

2 tools → 50 tool calls.

A lot to handle for a single agent.

Separation of concern.

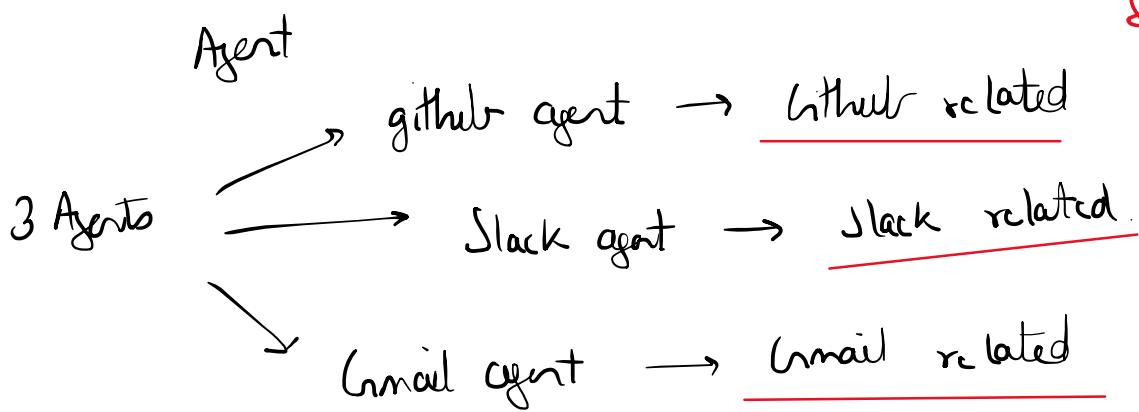
5 Agents

selected tools

{ complex task ✓
 ↓
 25 subtasks.
 ↓
 50 tool calls.

✓ Github Repo, ✓ Slack manager, ✓ Email send

Jack of all trades.



3 Agents which are experts in what they

do

Manager → Orchestrator model → Query.

tools → Action

Captain of a
Team

members attached.

↓
Planning
↓
Subtasks.

Team leader → Knows about the capabilities of all of its members.

3 members.

a) Github agent

b) Email agent

c) Slack Agent

Delegates tasks to members. → delegate task to members tool → Power.

1) Member → github.

2) task → Read the Readme file of this repo

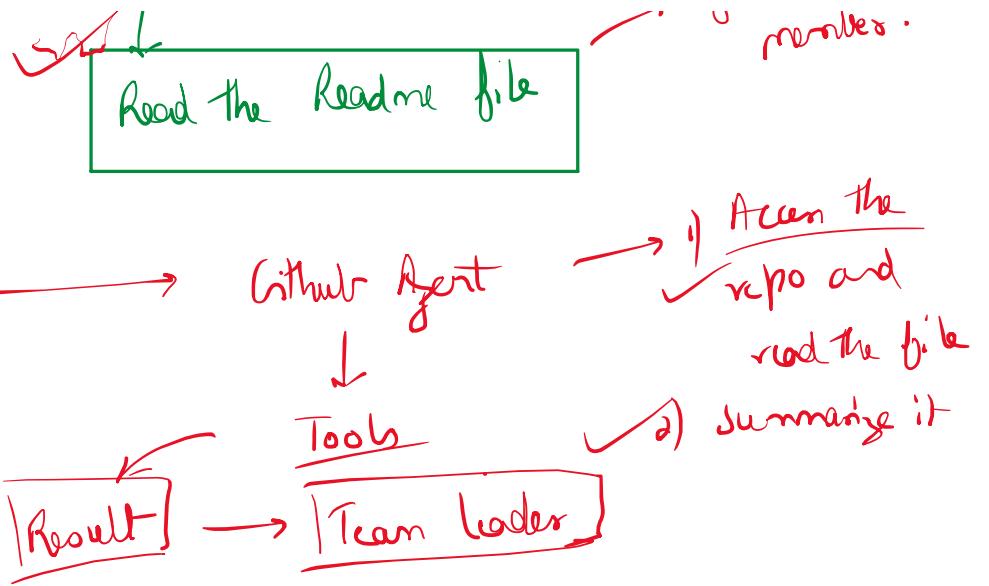
Github related tools

Team leader

Complex goal → 10 subtasks

... ↗ Find the Readme file ↗ github members.

delegate task to



How Execution Works?

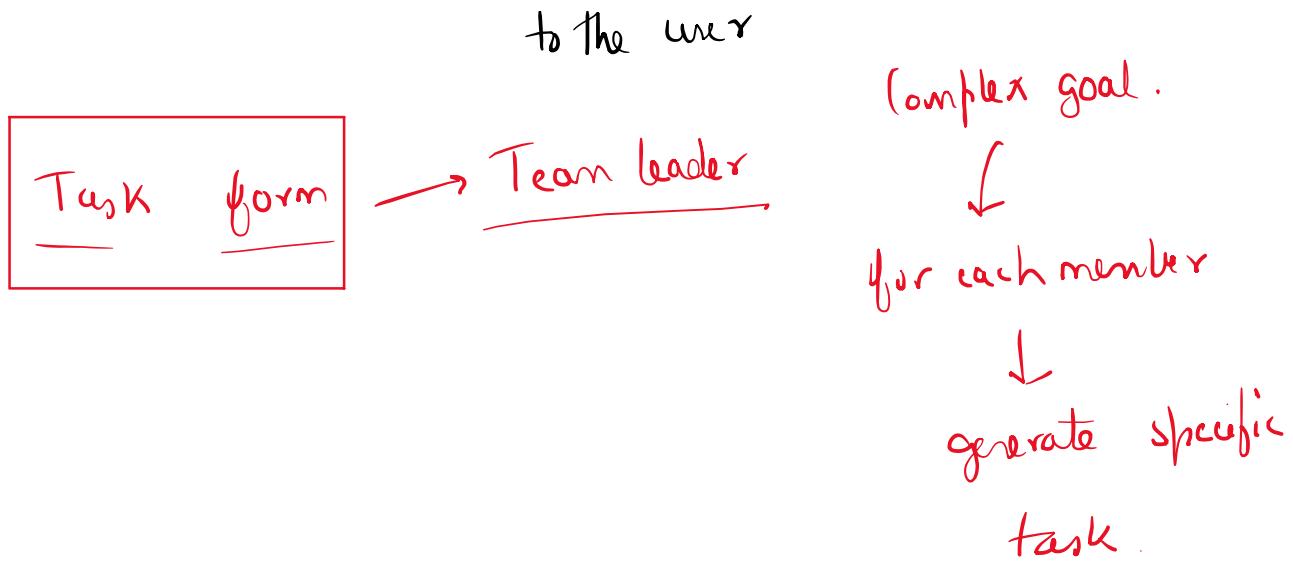
Team → Team-leader (Orchestrator)
↓
delegate tasks to its
members.

Flow:-

- 1) Team input → user query
 - 2) Team leader → reads the input → decides on how to break it into subtasks.
 - 3) Team leader delegates specific tasks to appropriate members
 - 4) Team members complete their assigned tasks and return their results. Agents
 - 5) Decide → whether to delegate further to members, or else it synthesizes all outputs into a fine cohesive response
- user response

↓

to the user
- Complex goal.



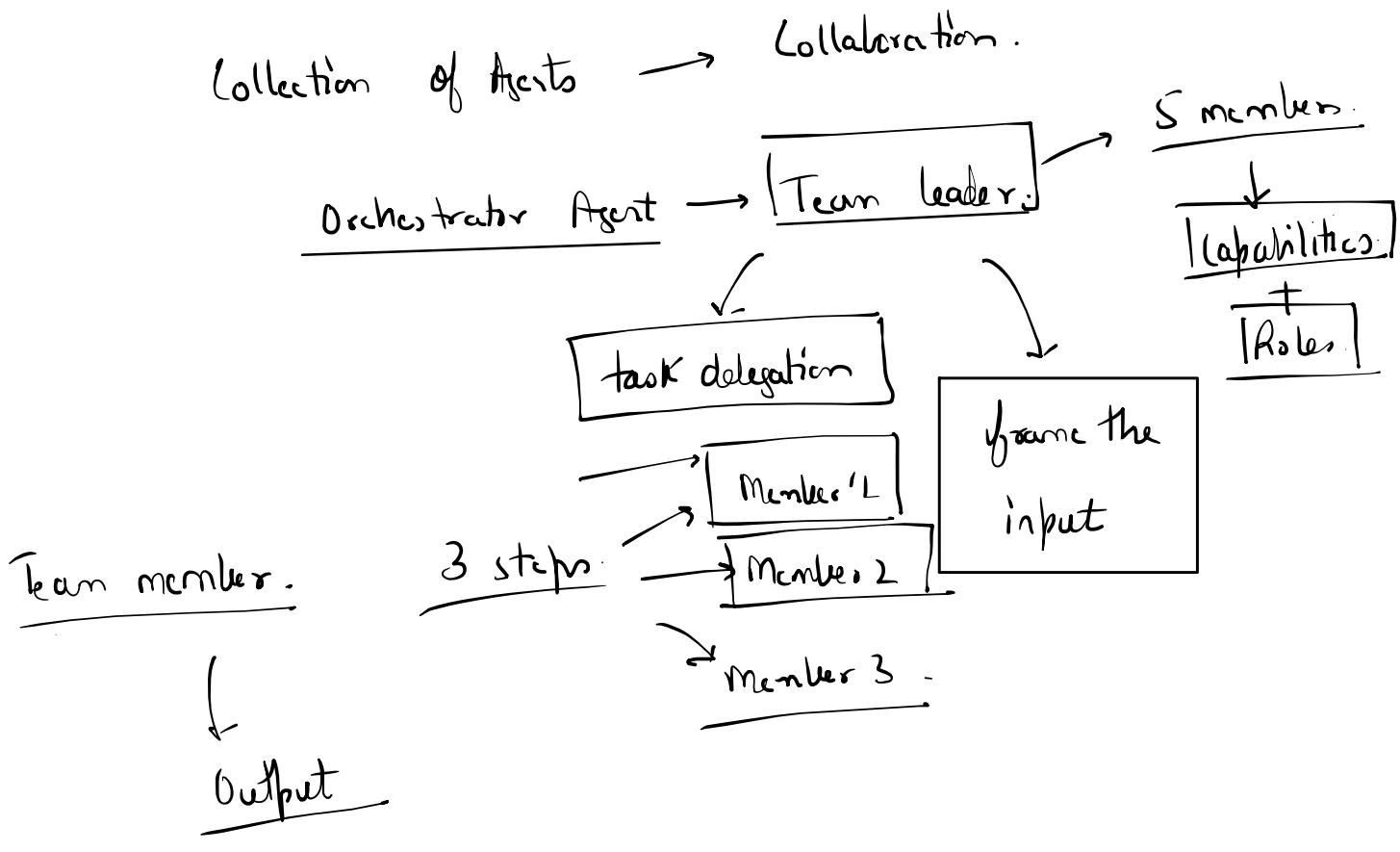
Team leader has 4 imp roles:-

- 1) Team leader understands the context of input query and creates subtasks.
- 2) Team leader knows about its members in the team, their capabilities and also which tools each member has.
- 3) Team leader has to delegate task to the member. The task itself is generated by the team leader and is generally diff. from the user query. This allows the member to perform the task efficiently.
- 4) Team leader compiles all the outputs from its members.

and send the response back to the user.

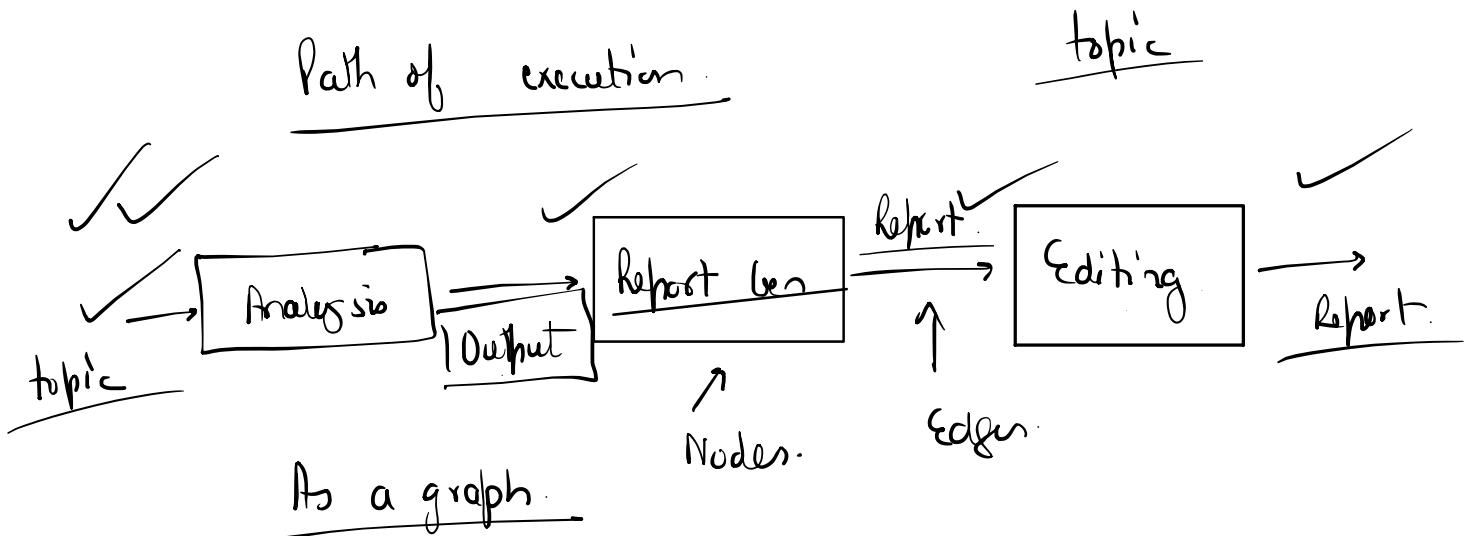
Recap

Teams → Multi-Agent system



Why Workflows

Deterministic Multi agent system.



Node → Steps.

Edges. → Connections

Agents → Steps. Report generation step.
→ Agent

Boxes. → Nodes

Workflow → Steps → order

1) Sequential.

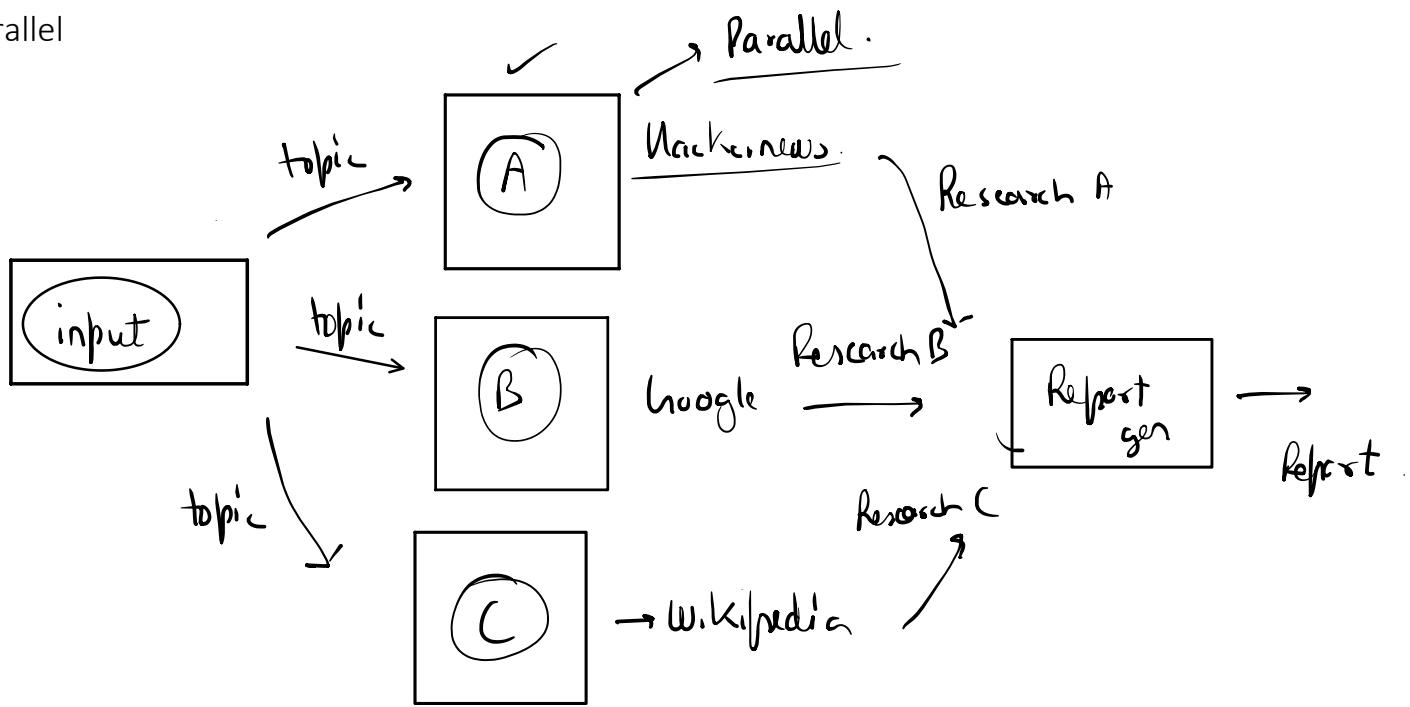
2) Parallel.

3) Conditional

Predetermined Path of execution

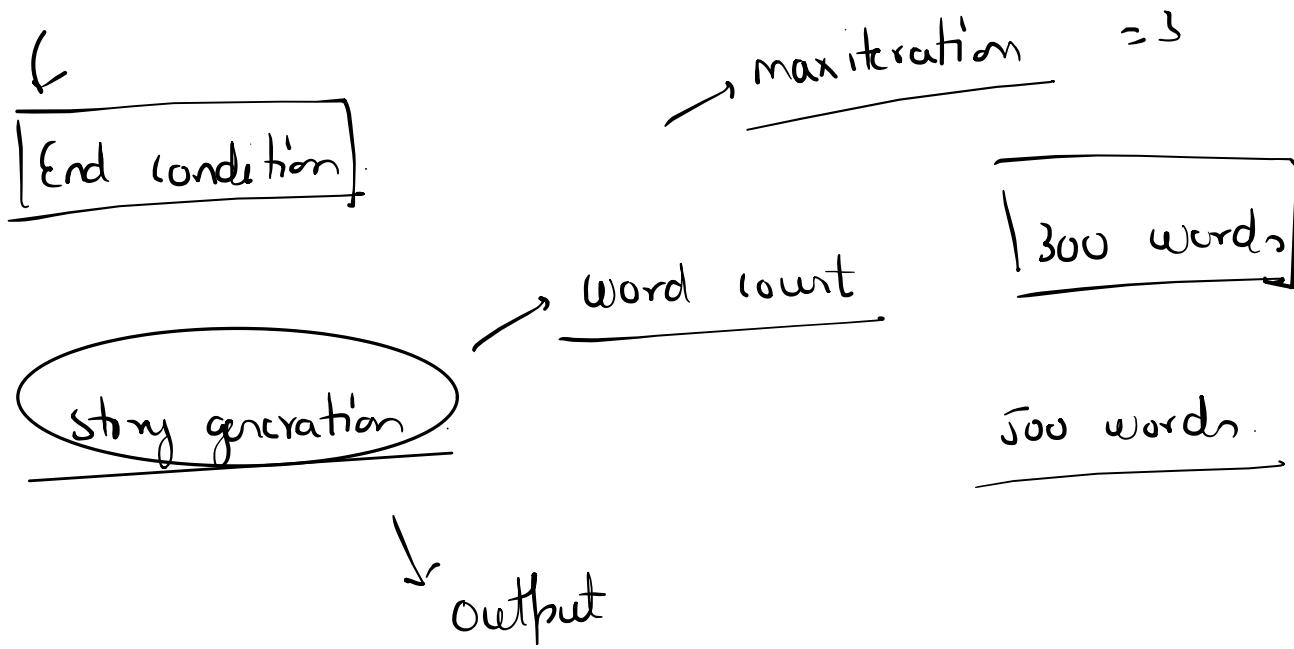
4) looping

Parallel



Looping

loop → Multiple Iterations.

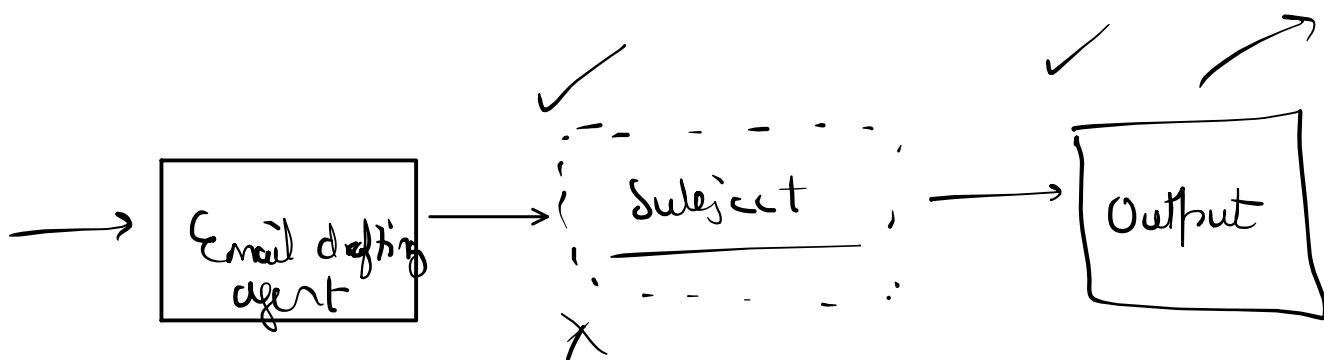
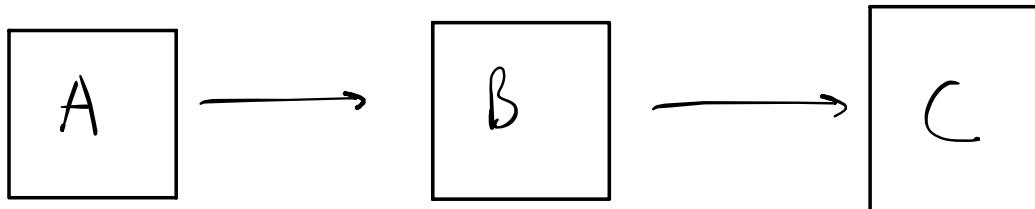


less than 300 words.

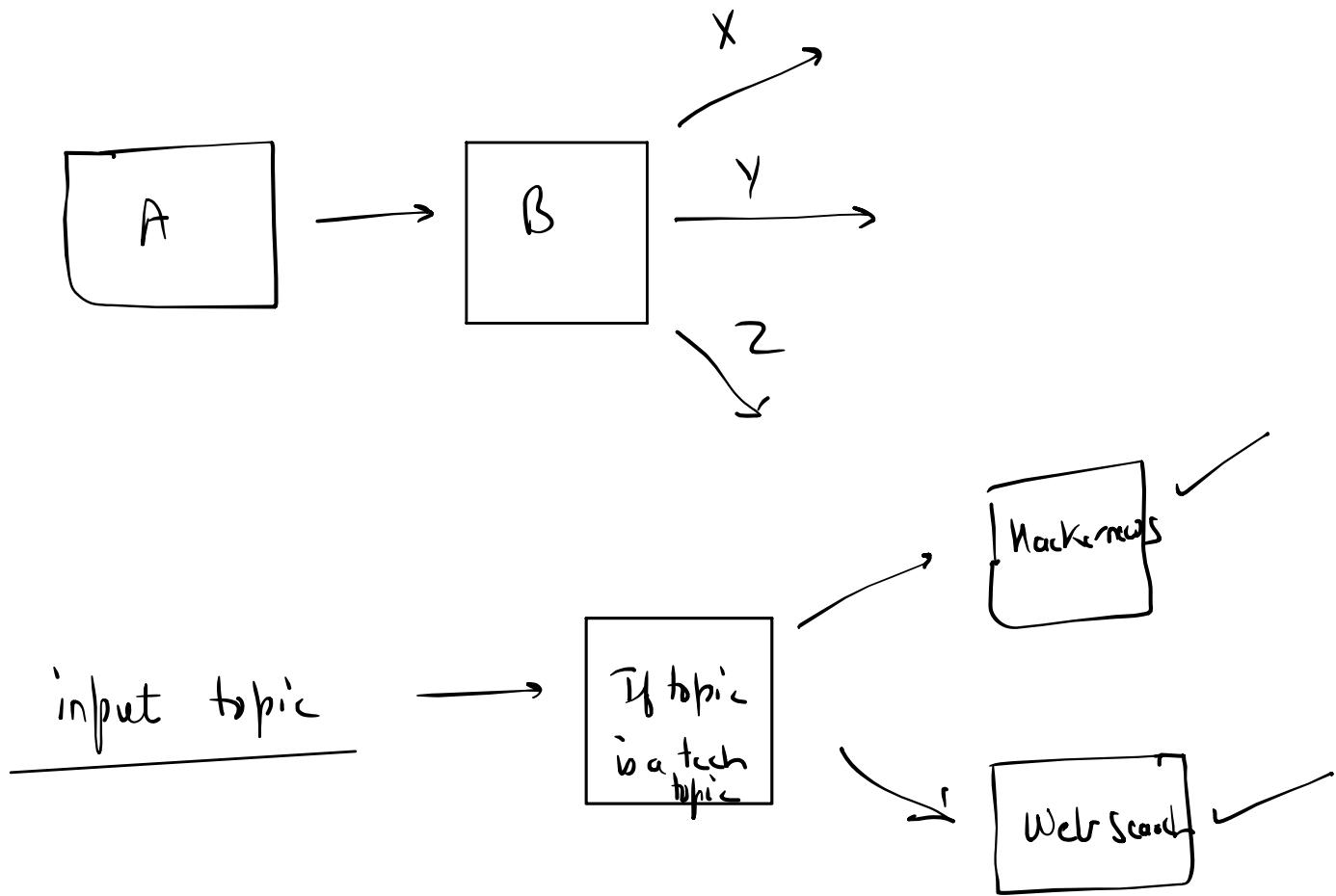
Conditional

less than 100 words.

↓ Post on X



Router/Branching



Workflows in Agno. → Predetermined series of steps.

Steps → Agent → task

connect → Sequential, parallel,
Conditional (dynamic workflows),
looping

Graphs → Nodes → Step

Edges → Connections,
Combinations

Medium Article generation Project

from where to read about the topic → Resources

Topic → Resources → Validate



Agno.

Validate ← Read



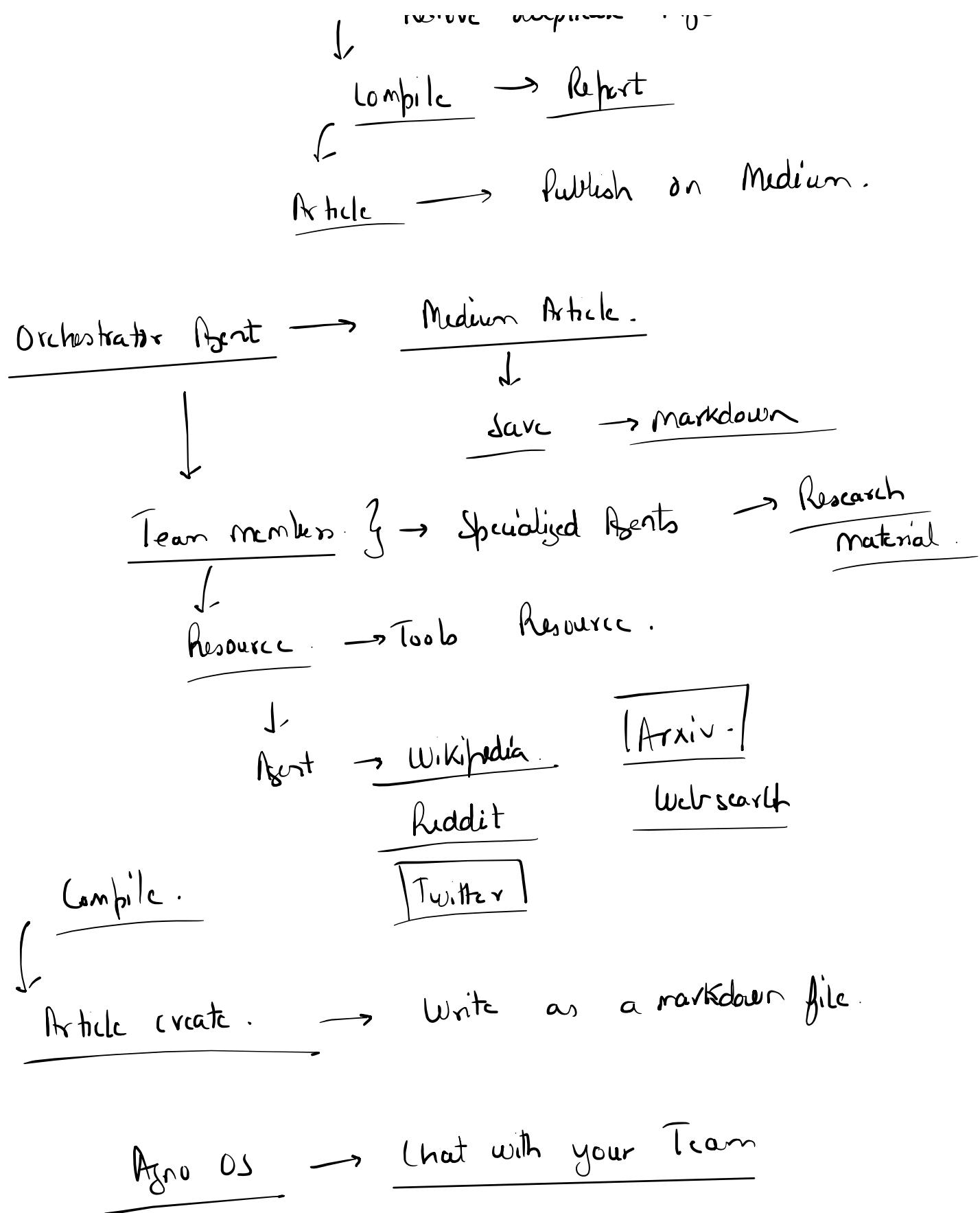
Compile

Summarize

Tearrs

, Remove duplicate info

Laplability



Data Science Team

Data Ingestion

Data Understanding

EDA

Data cleaning

Feature Engg

Model training

Model Evaluation → Metric

NP tuning stage

Orchestrator → Manager

Data ingestion

Tools

EDA

Plots

FE

Model training

Model Evaluation

Pandas

matplotlib

Chat interface → Dataset → dir

EDA → Summarized

→ artifacts

Web capability

Plots → savc

FE → Execute

train.csv
test.csv

Model training → Code

→ ← ↑ → ↓ model

