

# Agentic AI

## Langraph

Background

Learning AI

Lecture - 01

## Generative AI vs Agentic AI

refers to a class of AI models that can create new content - such as text, images, audio, code, video etc. that resembles human - generated data.

Traditional AI → It is about finding patterns in data and giving predictors.

Generative AI is about learning the distribution of data so that it can generate a new sample from it.

The best part of GenAI is that it mimics human. (their functions.)

### Applications Areas

- Creative and Business writing
- Software Development
- Customer Support \*
- Education
- Designing

So, The solution is ?

→ RAG Chatbot. (Retrieval Augmented Generation) By giving data from us, It stores that and make a parametric knowledge.

So it solving a -

- Specific Advice

21-08-2025

## AGFTIC-AI.

→ foundations

→ Langraph fundamentals

→ Advance Langraph

→ AI-Agent.

→ Agentic RAG

→ Productization

So, there is solution of rest problem!

## → Tool Augmented Chatbot (RAG + Tool)

It solves problems. →

- Specific Advice
- Can take action

Natireactive  
No memory  
Can't Adapt

So Now, most advanced chatbot for these

## → Agentic AI Chatbot (AI agents)

- It is proactive that can plan and execute.
- Obviously it has memory so it can plan.
- Specific Advice
- It can take actions.
- It's adaptable

Gen AI is about creating content

Gen AI is about solving a goal

Agentic AI is about doing things autonomously

Gen AI is reactive, Agentic AI is proactive

(autonomous)

Gen AI is building block of Agentic AI

(superset)

(subset) → Capability Behavior

\* trigger \* promoted \*

Lecture - 02

What is Agentic AI.

→ It is a type of AI that can take up a task or goal from a user and then work toward completing it on its own, with minimal human guidance \*

\* depends on previous lec 02

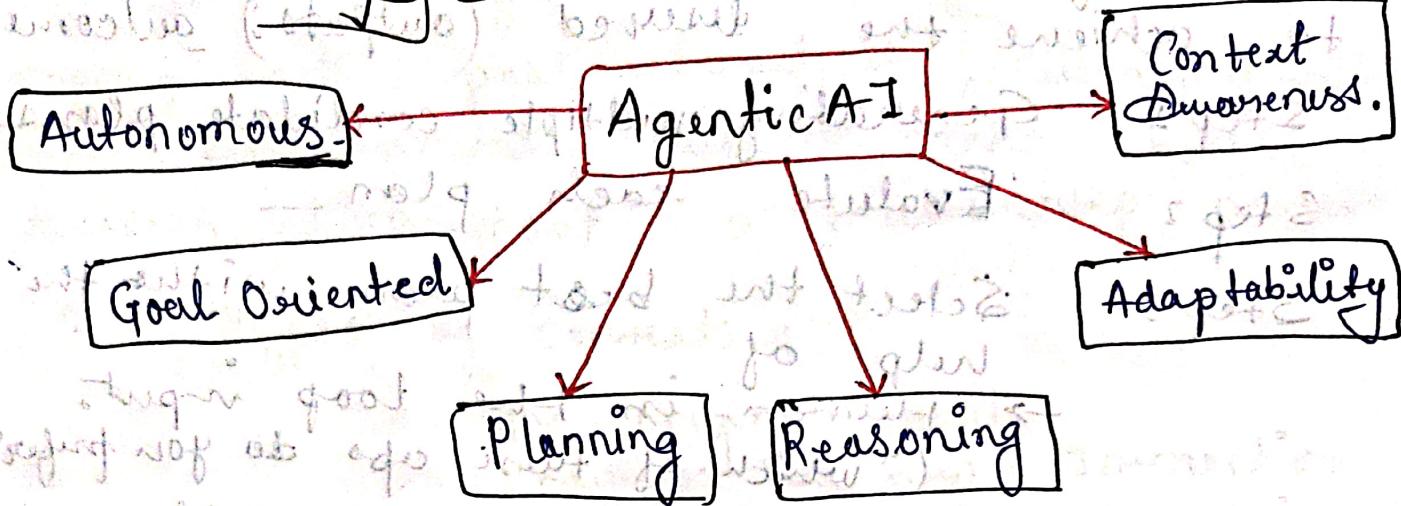
\* which is different



Scanned with OKEN Scanner

→ It takes plan, take actions, adapts to changes and seeks help only when necessary.

### Key Characteristics



Autonomy → Autonomy refers to the AI system's ability to make decisions and take actions on its own to achieve a given goal, without needing step by step human instructions.

But it can be controlled by Human in the Loop (HIL) through Override Controls or Guardrails / Policies. It can be dangerous if the autonomy is lost.

Goal Oriented → It means that the AI system operates with a persistent objective in mind and continuously directs itself to action to achieve that objective, rather than just responding to ad-hoc prompts.

1. Goals act as a compass for autonomy
2. Goals come with constraints
3. Goals are stored in long-term memory.
4. Goals can be modified.

Planning → Planning is the agent's ability to break down high-level goal into a structured sequence of actions or subgoals and decide the best path to achieve the desired (outputs) outcome.

Step 1: Generating multiple candidate plans

Step 2: Evaluate each plan

Step 3A: Select the best plan with the help of

→ Human in the loop input.  
→ Which of these ops do you prefer?

→ A pre-programmed policy.

Reasoning → It is the cognitive process through which an agentic AI system interprets information, draws conclusions and make decisions — both while planning ahead and while executing actions in real time.

Reasoning during planning / during executing

- 1. Goal decomposition
- 2. Tool selection
- 3. Resource estimation

- 1. Decision-making
- 2. HTL Handling
- 3. Error Handling

Adaptability → It's agents' ability to modify its plan, strategies or actions in response to unexpected conditions — all while staying aligned with the goal.

Failure of calendar API

External feedback (no. of application)

Changing goals (Hiring a broker)



## Context Awareness →

It's the agent's ability to understand, retain and utilize relevant info. from the ongoing task, past interactions, user preferences and environmental cues to make better decisions throughout a multi-step process.

Environment Agent work in  
Eg - Agent → play chess  
so, Env → Chessboard

## 1 Types of Context →

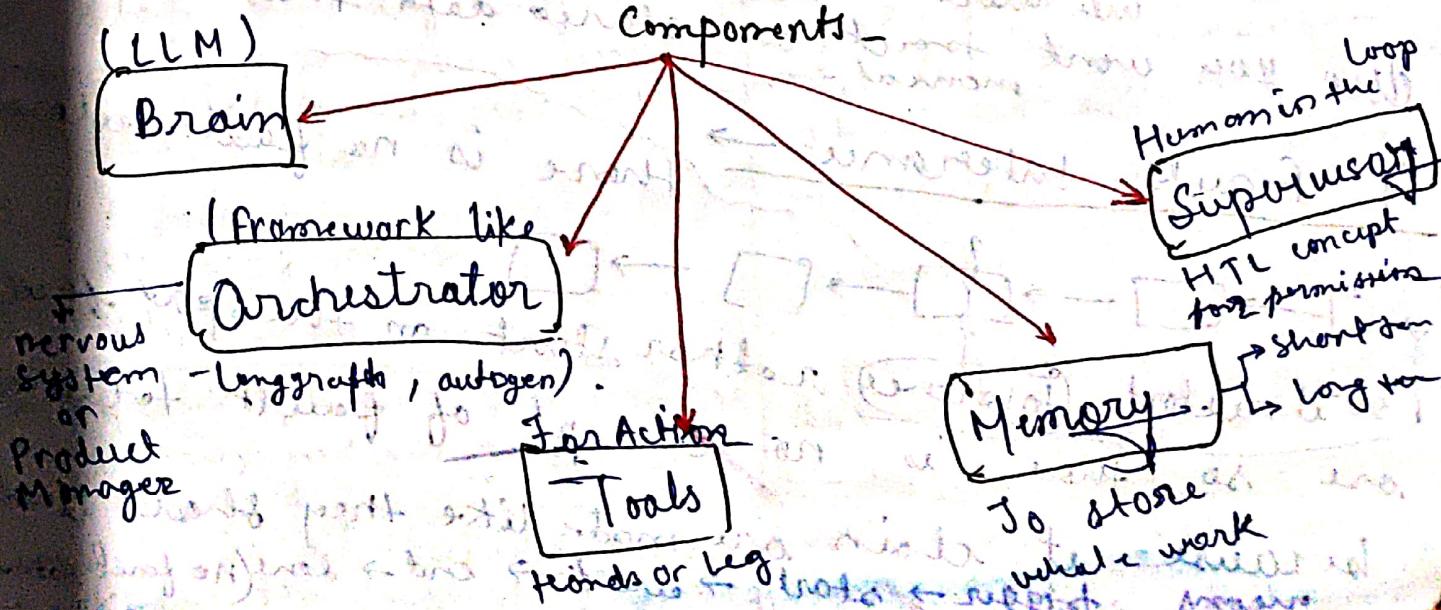
- The original goal
- Progress till now + Interaction history
- Environment state
- Tool response
- User specific preferences
- Policy or guardrails

2 Context implement through memory

3 Short term memory:

4 Long term memory:

## Components of Agentic AI [top 5]



ecture → 03

## LongChain & LongGraph

flow → long A → B  
early  
download → very  
long time

alternates of longGraph → Autogen, Creual etc.

why LongGraph → We create graph of nodes and connect by its edges  
Workflow → static, predefined paths

Agent → Dynamic, that creates their own path to execute process

longChains Challenges →

1. Control Flow Complexity

[No logic to write in longChains  
need to use python glue code]

that effect maintainability

It only efficient for linear workflow  
not conditional or parallel

2. Handling State

It has only conversational  
memory. there is no storing  
or tracking of data so it is

stateless. Think when you need to change its step trees  
we were have to do manually -  
and you went track every state data that is also  
manual -

conditional branch

loops

Jumps

glue code

State must set of datapoint  
and group or set of state is  
workflow

3. Fault tolerance

There is no fault tolerance

so it restart (tolerance) rather start on the spot fault  
there is no concept of fault tolerance  
because it chain are made like they short lived  
Linear → start → work → end → done/no fault toler

why longgraph →

1. No (zero) glue code

2. Stateful → graph

State object

→ Pydantic [mutable]  
TypedDict [is it]

3. check-pointer → to some/track  
(State)

there is possible event-driven execution in the workflow.

4. Fault Tolerance

is there →

retry logic & recovery  
(when server  
(for short time down or  
problem & system issue  
faults) & errors)

It is possible because  
of concept of  
check-pointers

also timer [resume func]

to resumes the work where  
is stop by defining a  
previous state - then automatically  
identify state and also their  
problem

Fault tolerance is very high  
in longgraph.

5. Human in the loop

longgraph allows you to pause  
execution indefinitely for  
min, hours, days even until human

input is received [because it saves step forward later resume  
to point the graph state after each step without time constraints.]

Challenge 3 →

Event driven Execution

Workflow will pause in  
conds and wait for external  
trigger to resume workflow

It is not build for this  
longchain is build for - at -

Challenge 4 →

Fault Tolerance

Challenge 5 →

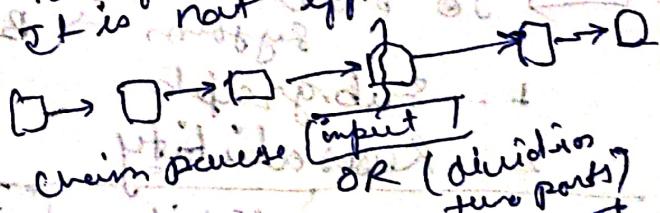
Human in the loop

pauses in workflow for  
human permission and  
their approval after it

start work

But in longchain ~~is~~  
not interaction in workflow  
you can take input

But the thing is if applicable  
for short running workflow  
not for long running  
it is not efficient and feasible



chain pieces OR (dividio  
two parts)

But not maintainable and glue  
code is fine up now

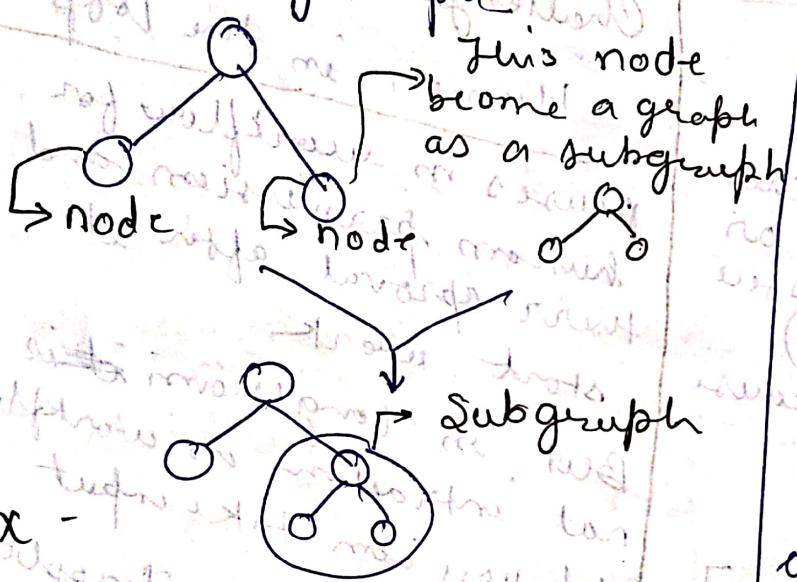
forward later resume  
the workflow without  
time constraints.

## Nested workflow

longGraph's Feature

→ It use it in like  
a subgraph →

A subgraph is graph feature  
use ~~as~~ ~~as~~ a node in another  
graph - this is concept  
of encapsulation applied  
to longGraph



Ex - Multigint System

like self driving cars

To build Multagent system → By

+ subgraph

2. Reusability

aut. Observability

when you deploy workflow

you can monitor very efficiently

with longsmith they are  
very tightly coupled to this feature

longChain's Challenge

It have no capability  
to do this

no concept of  
nested workflow

Challenge 7  
Observability →

refers to how easily  
you can monitor,  
debug, and understand  
what your workflow  
is doing at runtime

But → It is  
longChain but via  
longSmith. But et can  
only monitor longChain  
code not glue code.  
the bad things is in  
longChain a glue code is  
there a lot to build  
good workflow  
so, it is partial observability

It store whole timeline of workflow via longsmith in longGraph



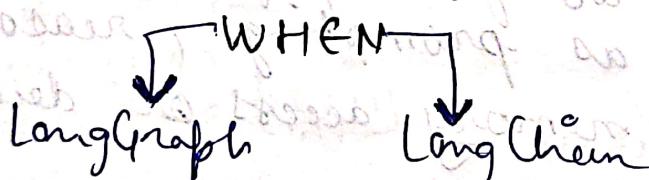
## LANGGRAPH

→ It is an orchestration

framework that enables you to build stateful multi-step, end event driven workflaess using LLM's. It's ideal for designing both single agent and multi agent agentic AI Applications.

Trunk of LangGraph is flowchart engine for LLM's

You define the steps (nodes) in how they're connected (edges) and the logic that governs the transition. LangGraph take care of state management, conditional branching, looping, pausing/resuming and fault memory — features essential for building robust & production-grade AI system.



- Conditional Path ← ◦ Simple
- Loops
- Human in the loops
- Multi Agent coordination
- Asynchronous or event-driven execution

→ LongChain Advance → is long Graph.

→ Not replacing the long chain.

Lecture → 04

## Long Graph

21-08-

### Core Concept

It enables advanced features like parallelism, loops, branching, memory and reusability making it ideal for agentic and production grade AI applications.

Long graph is an orchestration framework for building intelligent, stateful, and multi-step LLM workflows.

LLM Workflows → Series of tasks to achieve a goal

They are a step process using which we can build complex LLM applications.

Each step in workflow performs a distinct task - such as prompting, reasoning, tool calling, memory access or decision making.

Common Workflows →

1. Prompt Chaining

2. Routing

3. Parallelization

4. Orchestrator Workers

5. Evaluator Optimizer

Graph → Nodes + Edges →



Scanned with OKEN Scanner

State → State is the shared memory that flows through your workflow. It holds all the data being passed between nodes as your graph runs.

- Every node have access of states or shared input
- It is mutable and special type dictionary called TypedDict (alternate Pythonic) use via library.

Reducers → It define how updates from nodes are applied to the shared state.

- Each key in the state can have its own reducer, which determine new data replaces merges, or adds to the existing value.

### LangGraph Execution Model [Google Project]

1. Graph Definition
2. Compilation
3. Invocation
4. Super Step begin
5. Message Passing
6. Halting Condition

Execution stop when:  
→ No nodes are active  
→ No message were in transit

Lecture  $\rightarrow$  05 Sequential Workflows  $\rightarrow$  Practical

~~BM1~~ basic workflow  
Now adding some functionalities  
and features  $\star$

In this pass a  
partial state

$\rightarrow$  LLM workflow  $\rightarrow$  ~~partial State~~

Prompt Chaining  $\star$

Lecture  $\rightarrow$  06 Parallel Workflow  $\rightarrow$  Practical

$\rightarrow$  Simple Parallel Workflow  $\rightarrow$  In this passing  
batchness - workflow a partial-state

Parallelization Workflow  $\rightarrow$

UPSC Essay Evaluation

Workflow

Prompt Chaining  $\rightarrow$  22 ~ 0.8

Lecture  $\rightarrow$  07 Conditional Workflow  $\rightarrow$  Using manual way not command.

$\rightarrow$  Quadratic Equation Workflow

$$ax^2 + bx + c = 0$$

$\hookrightarrow$  roots  $\rightarrow$  2

$\rightarrow$  Discriminant  $= \frac{b^2 - 4ac}{d}$

$d > 0$  (2 real roots)  $d = 0$  (repeated root)  $d < 0$  (don't have any root)

$$x_1, x_2 = \frac{-b + \sqrt{d}}{2a}, \quad x_2 = \frac{-b - \sqrt{d}}{2a}$$

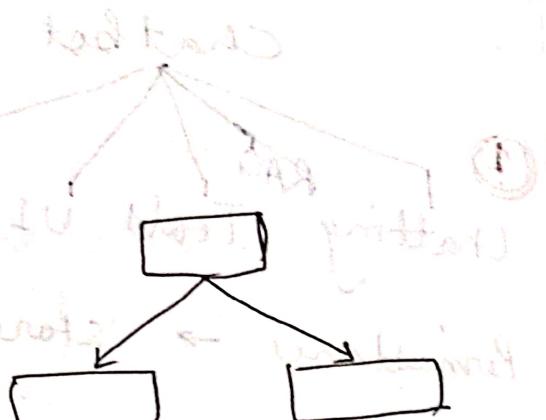
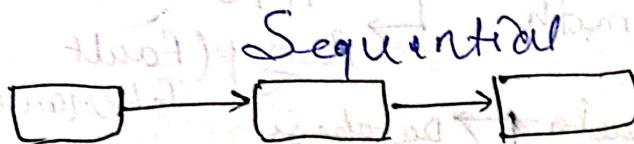
$$d < 0$$



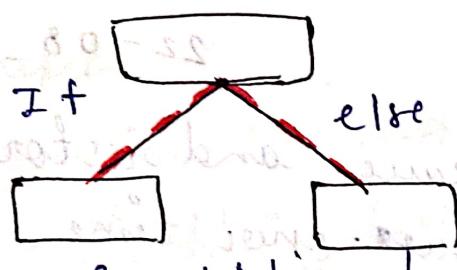
→ LLM Based Conditional workflow →  
↳ Customer Support Assistant Workflow.

## Lecture → 08

## Iterative Workflow

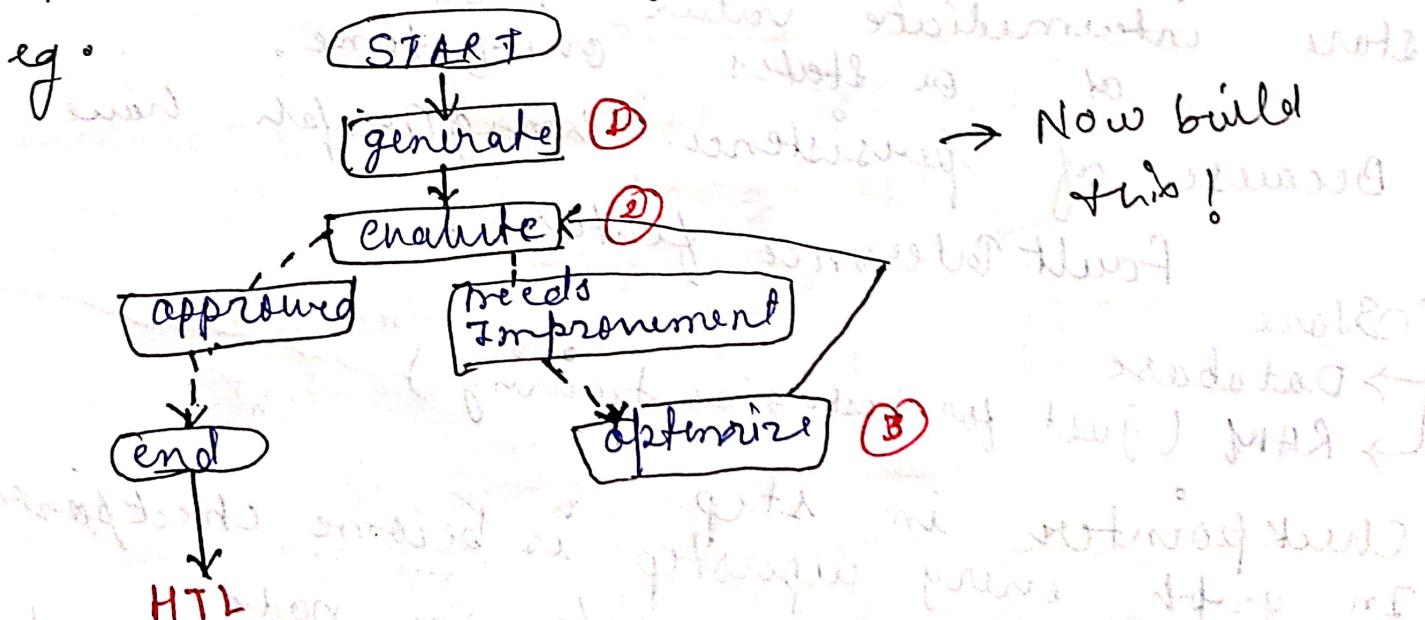


## Workflows



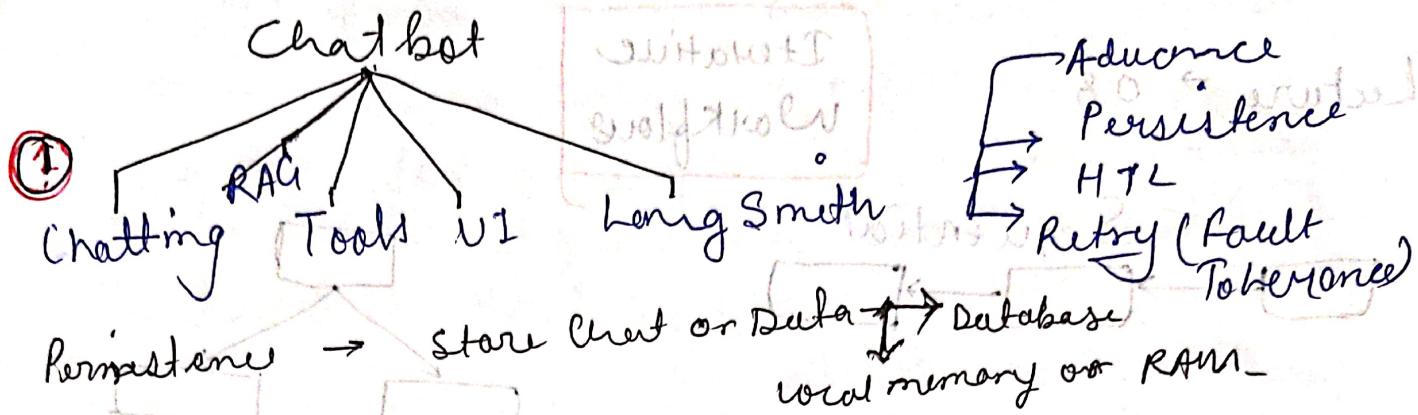
## Parallel

Iterative or looping workflow → Automated



Chatbot  
is longChain

1.



Persistence

In long graphs it refers to save and restore the state of a workflow over time.

It naturally only stores final value it also stores intermediate value means all values over time.

Because of persistence changing Graph frame

Fault Tolerance feature

- Store
- Database
- RAM (just for testing or training)

Checkpointer in step In graph every superstep is become checkpointer Superstep → when in graph's one node going multiple nodes or this process or step is called one superstep

Each time on one superstep is check by checkpointing

Threads Persistence → While you store data in database after multiple state or multiple execution, so to don't occur the issue and conflicts we provide an id to each execution and that is known as Thread ID.

By defining the executions by assigning thread ID user can resume their history to check and see it's results.

Benefits →

- Short term memory.
- fault Tolerance
- HTL (Human in the loop)
- Time Travel. [for replay the execution of code]

Lecture - 11

Chatbot with UI  
in LangGraph

LangGraph + Streamlit  
Backend Frontend

Done

Lecture 12

Streaming  
In LangGraph

Done

In LLMs, streaming means the model starts sending tokens (words) as soon as they're generated, instead of waiting for the entire response to be ready before returning it.

Why →

- Faster response time - low drop off rates
- Mimics human-like conversation (Builds trust, feels alive and keeps user engaged)

3. Important for news model etc
4. Better UX for long autosearch as code articles etc
5. You can cancel midway & saving taken.
6. You can interleave UI updates -  
eg, show a "thinking..." , show tool result.

**Generator** → In Python, a generator is a special type of iterator that allows you to generate values on the fly, one at a time, using the `yield` keyword instead of `return`.

### Lecture - 13

[chatbot history]

Resume Chat

feature

23-08

Done

(4)

like ChatGPT or any other models

### Lecture - 14

SOLITE  
In hangGraph

Done

(5)

### Lecture - 15

Lang Smith

23-08

**Observability** → is the ability to understand a system's internal state by examining its external output. like logs, metrics and traces. It allows you to diagnose issues, understand performance and improve reliability by analyzing data generated by the system.

Essentially, it's about being able to see what's happening on this specific system - even if you don't have access to the problem.

LangSmith → It is a unified platform where teams can debug, test and monitor AI app performance.

What does LangSmith Trace?

1. Input/Output

2. All the intermediate trace

3. Latency

4. Token usage

5. Error cost

6. Tags

7. Metadata

8. Feedback

9. Execution of each component in project

Core Concept



Project

Trace

Run

LLM Application

WHILE In Production -

Reg have two big failure modes:

1. Retriever errors - wrong/irrelevant docs retrieved
2. Generator errors - model hallucinates or misinterprets context

In production, it's unclear where it occurs

So, LangSmith helps in → User Query automatically gets records → Retrieved → LLM prompt (with inserted docs) → LLM response

for trace every component of the LLM  
you have define library of import  
traceable execution



## o LongGraph + LongSmith

- Every graph execution can be logged in LongSmith as a trace.
  - Each node becomes a run inside a run inside the trace (eg - retriever, LLM, tool call, subgraph)
  - You can visualise the path taken:
    - o START → Retriever → Reranker → LLMans → END
  - If a workflow branches (conditional) / parallel/ subgraph), LongSmith was captured ~~multiple~~ which path was executed.
- Other features of LongSmith → under LLMOP's

## Monitoring & Alerting

Monitoring in this took across many traces at once to track the overall health of your LLM system. It aggregates key metrics (latency (P50, P95, P99), token usage, cost error rates, and success rates). You can set up alerts to notify you unless these metrics drift outside acceptable ranges (eg. spike in latency, higher error rates, major unexpected cost growth).

In production, issues often appear first as patterns across multiple runs, native to one is single trace.

Observability

+  
single trace

Monitoring



having multiple trace across and studied all

Evaluation: It helps you systematically measure the quality of your LLM's outputs.

You can test against gold-standard datasets or apply custom evaluation such as faithfulness, relevance, or completeness. It supports multiple approaches:

LLM as a judge  
semantic similarity checks

custom Python Evaluators.

Try out LLM on some datasets this called evaluation.

Prompt Experimentation:

It allows you to systematically test and compare different prompt revisions.

You can run A/B tests across prompts on the same datasets and track their performance against evaluation metrics and record the outcomes.

Results are stored over time, giving you a clear history of which prompt variations worked best under what conditions.

Also provides prompt versioning.

Dataset Creation & Annotation:

- Provides tools to build datasets for evaluation and fine-tuning
- Supports manual annotations (e.g. labeling whether an LLM's response is correct)
- Stores datasets versioned for reuse across project.

Why → High-quality datasets are critical

for evaluation and feedback loop.

## User feedback Integrations

- Let's you capture thumbs up/down, rating or structured feedback from users.
- supports bulk analysis of what like/dislike
  - feedback is alongside traces → tied to the exact prompt, model and state

## Collaboration

- Team members can view, share and comment on traces datasets and evaluation
- provides a web UI where non engg (PMs, QA, annotators) can inspect and annotate runs.
- enables shared experiment dashboard.

Lecture → 16

### Observability In LongGraph

Done

6

Adding this feature in Chaffcut and connect it to LongSmith.

Lecture → 17

### Tools in LongGraph

Done

7

Tool node → parties a prebuilt node type that acts as a bridge between your graph and external tools (functions, APIs, utilities etc)

tools\_condition → It is prebuilt conditional edge function that helps your decide.

Should the flow go to the Tool Node next, or back to the LGM?