

AI Agents

Week 1 → foundation

- Setting up environment →
 - IDE → Cursor → Cursor is an AI platform that's powered by LLMs that allows us to be so productive and it's built on VSCode.
 - UV → used in a virtual environment. It is really fast and simple. It is a package manager
- steps →
 - ① installing git
 - ② clone the Repo
 - ③ install cursor
 - ④ setting up uv.
- Using dotenv module and its load_dotenv class for reading the api-key stored in .env file.
- load_dotenv is an utility function that looks for a file with a particular name in project root directory and it happens to be always .env or name.env.
- .env consists of the environment variables.
- And load_dotenv function load the variables for using in the project.
- load_dotenv utility looks at the file, interprets it, and sets the environment variables as a result.
- a parameter override = True tells that if you already have set the environment variable in the past then override that variable value.
- import os → Os module provides us a function getenv() that read the api-key from the .env file.
- import openai → since in this project we are using gptq-api-key instead of Openai-api-key some functions & classes and modules are different.

→ Plan an agent, M N
coming up with number of queries.

- ① OpenAI →
 - OpenAI is a lightweight library for connecting to OpenAI's endpoints in the cloud.
 - OpenAI format to represent messages →
 - messages = [{"role": "user", "content": "What is 2+2?"}]
 - It is represented as list of dictionaries where each dictionary has a role and a content
- ② Calling OpenAI →


```
response = openai.chat.completions.create(
    model = "gpt-3.5-turbo", // model to be used,
    messages= messages // the message in openAI
) // the message in openAI
print(response.choices[0].message.content)
```

③ Day 2 → Agents & Agentic Architecture

→ What is an Agent?

Ans → AI Agents are programs where LLM outputs control the workflow. → huggingface
that simply means that one output from LLM is able to decide what tasks are carried out in a sequence.

→ People tends to often use the term Agentic AI, if any or all of the hallmarks are met →

- ① multiple LLM calls
- ② LLMs with ability to use tools
- ③ An environment where LLMs interact with each other, a kind of coordination/orchestration environment.
- ④ A planner to coordinate activities.
- ⑤ Autonomy → this term really defines the essence of agentic AI.

→ Plan for an agent, coming up with naming scenarios

① Agentic Systems

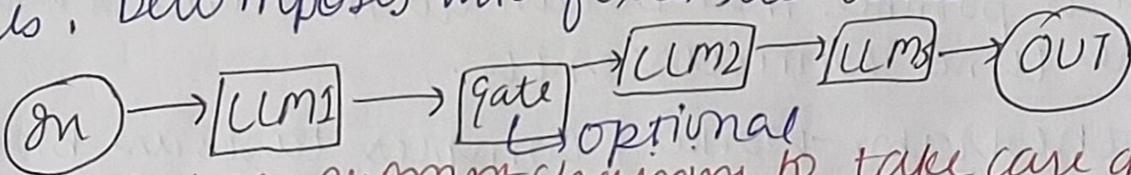
Anthropic distinguishes two types:

- ① **Workflow** → are systems where LMs and tools are orchestrated through predefined code paths.
- ② **Agents** → are the systems where LMs dynamically direct their own processes and tool usage, maintaining control over how they accomplish tasks.

① **Workflows** → anthropic identifies 5 different design patterns that you'll find when building agentic system that has workflows in them.

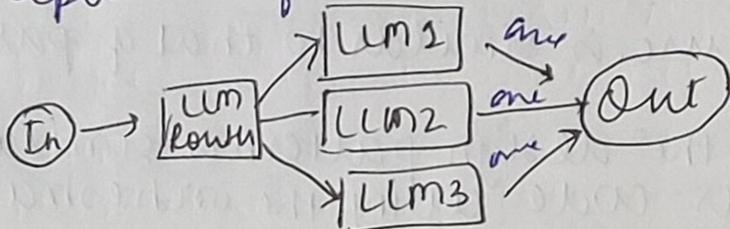
5 types:

① **Prompt chaining** → creating a chain of LM calls. Decomposes into fixed sub-tasks.



and we use prompt chaining to take care of each LM call very precisely, get the most effective LM response.

② **Routing** → Directs an input into specialised sub-task, ensuring operation of concern.



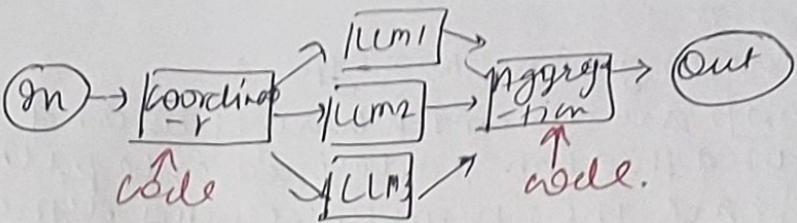
→ LM1, LM2, LM3 are specialised to do some different tasks.

→ LM router job is to classify the tasks, which of the specialists will be best equipped to classify the task.

→ In the model, there is some artificial agent used.

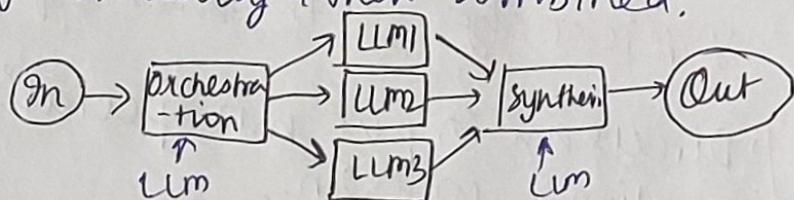
③ **Parallelization** → Breaking down tasks and running multiple subtasks concurrently.

R
R
B
P



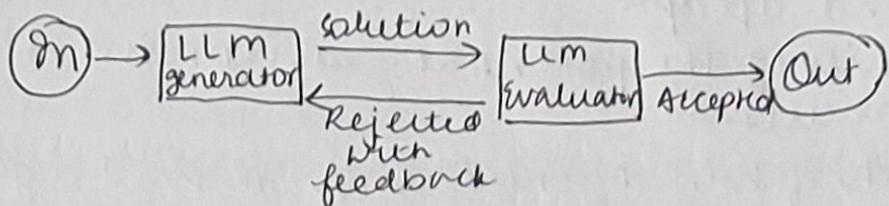
- Coordinator is a node that takes a tasks and break it down into multiple ~~stages~~ pieces that should all run in parallel.
- so, in routing workflow one of the LLM model was executing the ~~task~~ task at one time but in parallelization we are giving to three ~~task~~ elements to carry out three different activities concurrently.
- It is possible that we can do the same task for three times concurrently.

① Orchestrator-worker Complex tasks are broken down dynamically & then combined.



- this structure is similar to that of parallelization structure.
- Although the design pattern is similar but it is not longer code doing the orchestration it is an LLM
- Now, instead of code, we are using an LLM model to break down the task into smaller steps and then again an model to ~~to~~ combine the results.
- we can debate on the topic that orchestrator-worker design is an agents rather than an workflow since, the model can choose the way to divide the task and it could choose how many LLM's get assigned any activity.
- ~~LLM~~ break down the task → LLM carry out the task → LLM combines → Out

→ ⑤ Evaluator-optimizer → um output is decided by another.

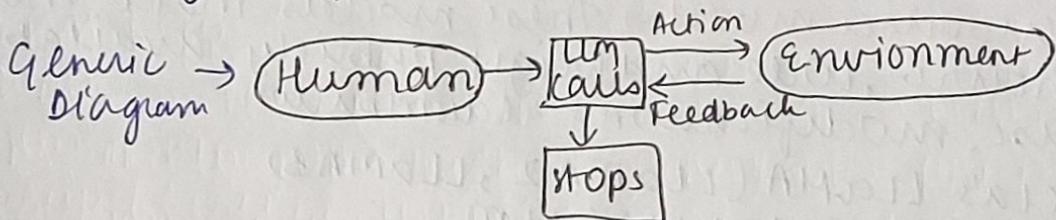


→ in this design two um's are used one for doing some specific task and the other one for evaluating the output of the first um's.

- LLM evaluate only evaluate the work of the first um, and based on that it choose either it is accepted or it is rejected.
- If ~~the~~ it accepts it, then it goes to the output.
- If it is rejected then rejection with the reason goes back to the um generator.
- It uses a powerful feedback loop.
- It is an effective pattern.

② Agents → By contrast with the workflow patterns agents are →

- ① open-minded open-ended
- ② it has feedback-loop
- ③ NO fixed-path → It does not use serious of steps



- human initiates a request or input to um
- um processes the human's input & processes it to generate an "Action" that is executed or applied with the function.

→ Please come up with more

⑤ Stable Diffusion (by Stability AI)

use → Text-to-image generation

⑥ Whispr (by Cohere)

→ environment provides feedback back to the LLM
based on the action executed, allows the LLM to learn and adapt.

→ After achieving the potential goal, LLM can stop the process execution

① Drawbacks of using Agents framework and agents design patterns →

① Unpredictable path

② Unpredictable output

③ Unpredictable cost → because we do not know how long it will take

④ Monitoring → visibility & understanding of what's going on behind the scene.

⑤ Guardrails → it ensures your agents behave safely consistently, and within your intended boundaries.
- it makes sure that model is doing what it is supposed to be doing.

Day 3 → Orchestration LLMs

① Calling multiple LLMs →

→ We ~~will~~ be using either Paid APIs and open source models in the cloud and locally throughout this course.

Notes → Open source model → These models make their architecture, code and weight publicly available for free use, modification & redistribution.

① Meta's LLaMA (LLaMA2, LLaMA3)

use → General purpose LLM

② Mistral & mptmal

use → General purpose LLM, in chatbots & coding assistants.

③ Falcon

use → Language generation and reasoning

④ Bloom

use → Multilingual language generation

④ can do some modeling (), Pw

⑤ Stable Diffusion (by Stability AI)

use → Text-to-image generation

⑥ Whisper (by OpenAI)

use → speech-to-text transcription

⑦ Closed source Models →

These models are proprietary - source code, training data, and weights are not shared publicly. They are typically accessed via paid APIs.

⑧ ① GPT-4 / GPT-4o (By OpenAI)

- use → chatbots, text generation, reasoning, vision, coding

② Claude (by Anthropic)

use → ethical AI assistant + Enterprise Application

③ Gemini (formerly Bard, Google DeepMind)

use → Web assistant, coding + research

④ Cohere Command R+

use → Retrieval-Augmented Generation (RAY) - NLP

⑤ Amazon Titan (by AWS)

use → Enterprise NLP tasks.

→ We can use 'Gemini-2.0-falcon' with Gemini API key

→ DeepSeek is a upstart startup that has powerful model.

→ Qwen - open-source LLM's including Llama3.3

→ Qwen - open source LMs on the system locally very similar to openAI

Lab 2 →

① In this Lab, they are experimenting with open-source closed source of LM for seeing the difference in the way of answering the question.

② And the second objective is to do orchestration between models.

→ openai + got - 4.0-mini

↳ we use client.messages.create

we use openai.chat.completions.create(
model = " ",
messages = messages);

3

answer = responses.choice[0].message.content

→ anthropic Claude - 3.7 - sonnet-later

Claude = Anthropic, we use claude.messages.create(
model = " ",
message = " ",
max_tokens = " ")

answer = response.content[0].text

→ gemini (gemini-2.0-flash).

gemini = OpenAI(api-key = __, base-urI: ' ')
→ same as method you did in openai.

OpenAI is used with google.gemini, inside. OpenAI
function you will mention the base-urI for
building an endpoint.

and by giving the base-urI we are switching to
google's endpoint instead of openAI.

→ Deepseek is executed as that of gemini google.
like base-urI.

→ groq is again similar to gemini & deepseek.

→ "!" → shortcut to execute terminal commands directly
from python cell.

→ ("") triple quotes in python is used for
writing multi-line strings and docstrings

and integrations

winding up with number 6 statements.