

13 NOVEMBER 2024

Session 6: Chains and Agents

GenAI BootCamp

Northeast

Presentation contents

01	Round Robin - How'd your homework go?	5 min.
02	Introduction to Chains and Agents	30 min.
03	Chaining in Action	15 min
04	Challenge Project Intro + Teaming	10 min.

In this session

We'll cover

- Round robin
- Chaining and Agents Basics - LangChain
 - Schemas
 - Models
 - Prompts
 - Chains
 - Agents
- Chaining in Action Demo
- Challenge Project

You will

- Learn how & when to use chains specifically LangChain
- Understand the basics of Chaining
- Understand the basics of Autonomous Agents
- Understand Use Cases for both
- Witness chaining in action

Come prepared

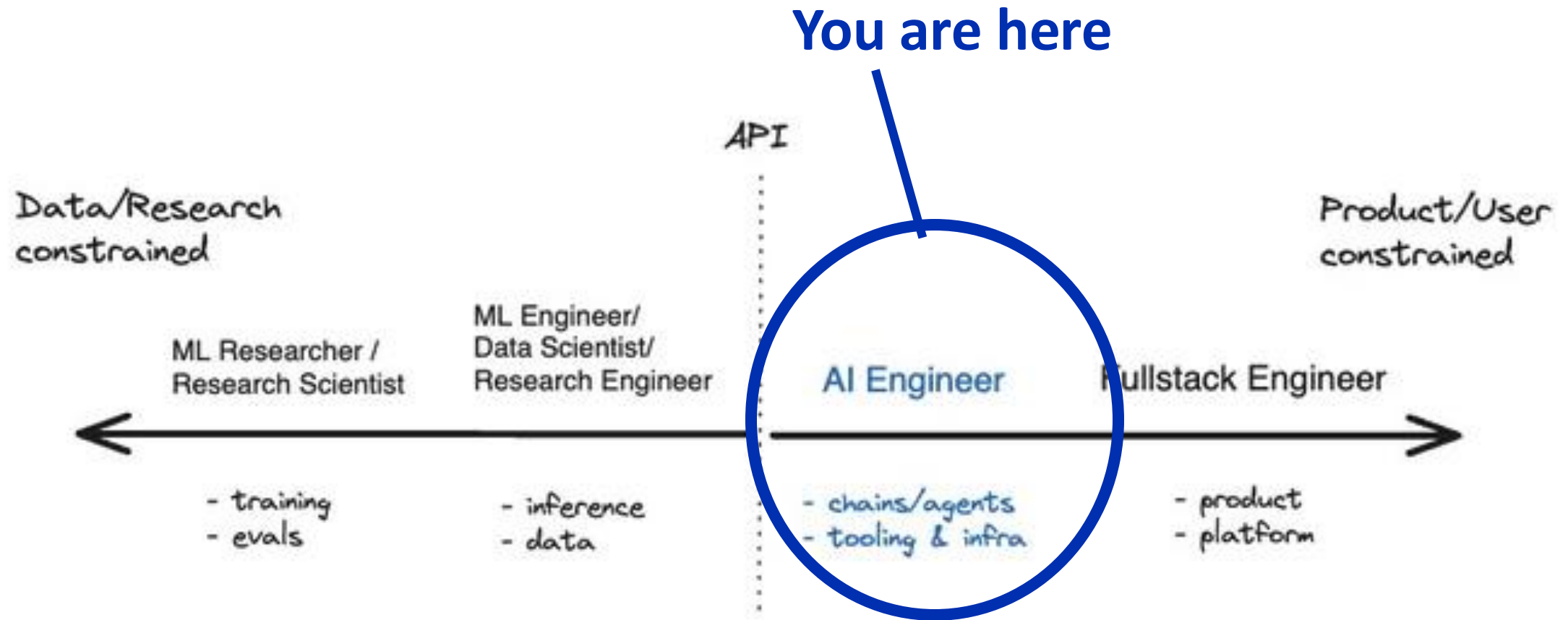
- Bring your homework from last session and any questions
- Begin thinking about your Challenge Project

01 Round Robin – How'd your offline work go?

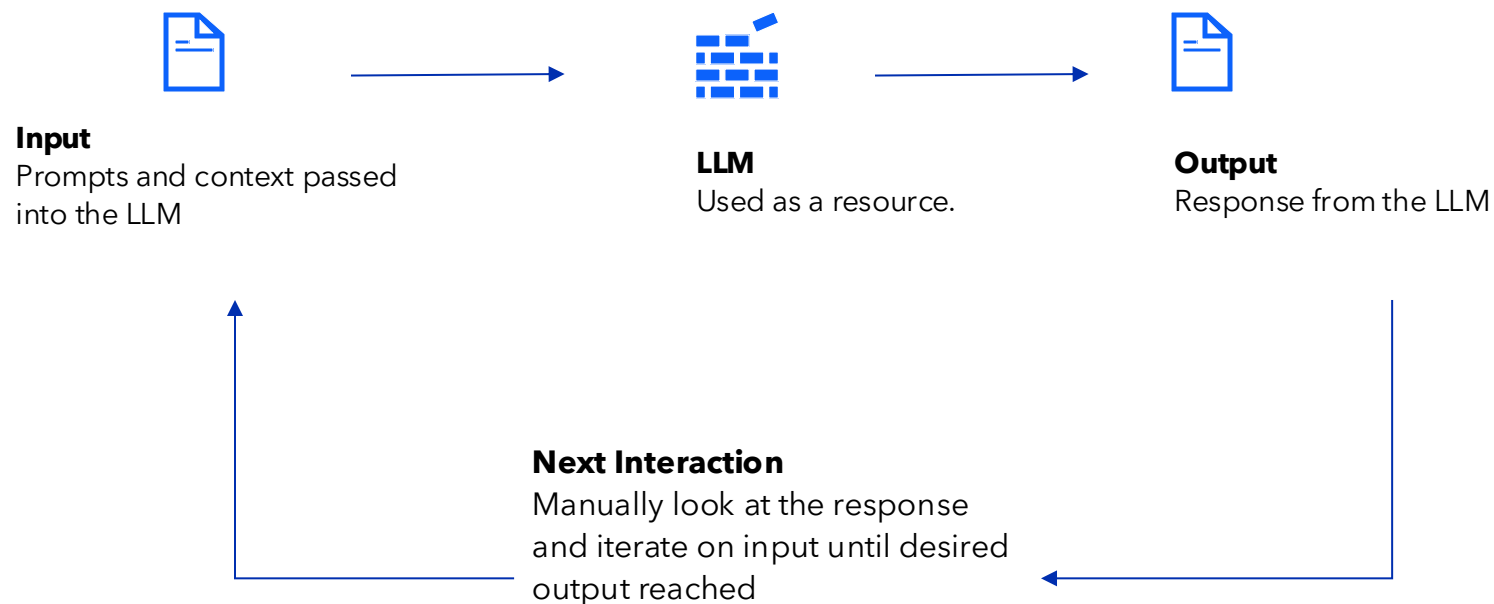
02 Chains and Agents



Fair warning



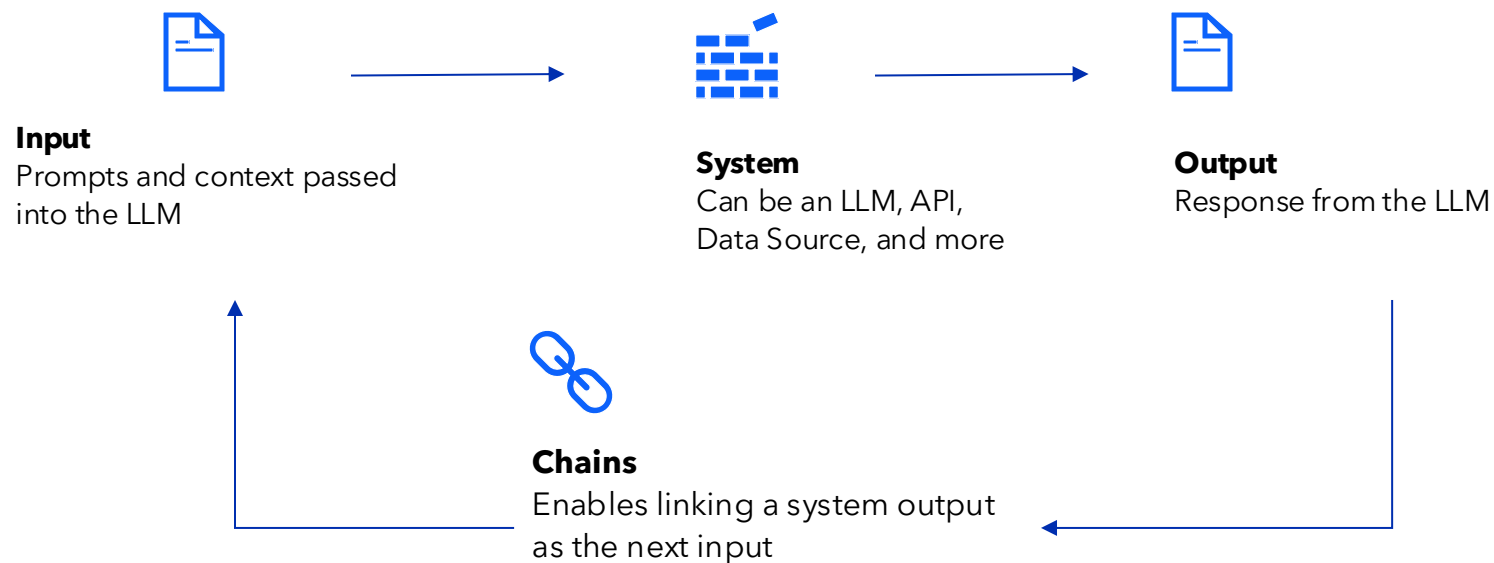
LLM Interaction Loop



Highlights

- LLM interaction loop is usually a single interaction
- LLM doesn't have "memory", requires context to be passed into with interactions
- May run into token limitations (e.g., 4096)

LLM Interaction Loop with Chains



Highlights

- Enables an abstraction layer to plug in different LLMs and systems
- Input Chain = Input -> System -> Output
- Different types of Chains available that can be linked together

Chain Examples

Using LangChain

LLM Chain

Calls different LLMs

Transformation Chain

Run a function on the input or output

Sequential Chain

Join chains together

Conversation Chain

Assumes a conversation is taking place

PAL Chain

Rewrites natural language into python code

SQL Database Chain

Rewrites natural language into a SQL Query

Bash Chain

Call bash commands

Request Chain

Requests an HTML page

API Chain

Call APIs

Chaining and Agents - LangChain

Chaining tools such as LangChain provide modular abstractions for the components necessary to work with models. LangChain also has collections of implementations for all these abstractions.



Schemas

Blueprint guiding the interpretation and interaction with data, mainly utilizing a "text in, text out" principle. This foundation lays out the ground rules before initiating the interaction.



Models

Primarily relies on three types of models: Large Language Models (LLMs), Chat Models, and Text Embedding Models, each playing a unique role in enhancing LangChain's versatility and strength.



Prompts

Functioning like the steering wheel guiding the model's direction based on the input questions or statements



Indexes / Memory

Efficiently organize and retrieve information to support easy and effective user interactions. Indexes in LangChain are crucial for processing and locating specific data amidst a massive amount of information. Memory persists context.



Chains

The system's master orchestrators, bringing together diverse elements to create responses from the language models that are meaningful and applicable, creating a system that is both unified and workable.



Agents

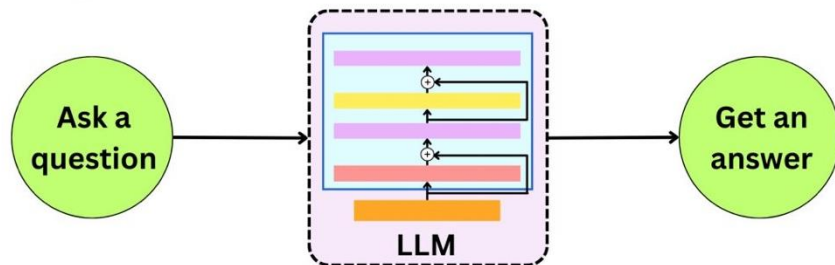
Action Agents perform like sprinters in a race, doing swift, precise actions, making them ideal for minor tasks. In contrast, Plan-and-Execute Agents are strategic and endurance-focused and excel at tackling difficult or long-term activities.

There are two main value props that chaining frameworks such as LangChain provides:

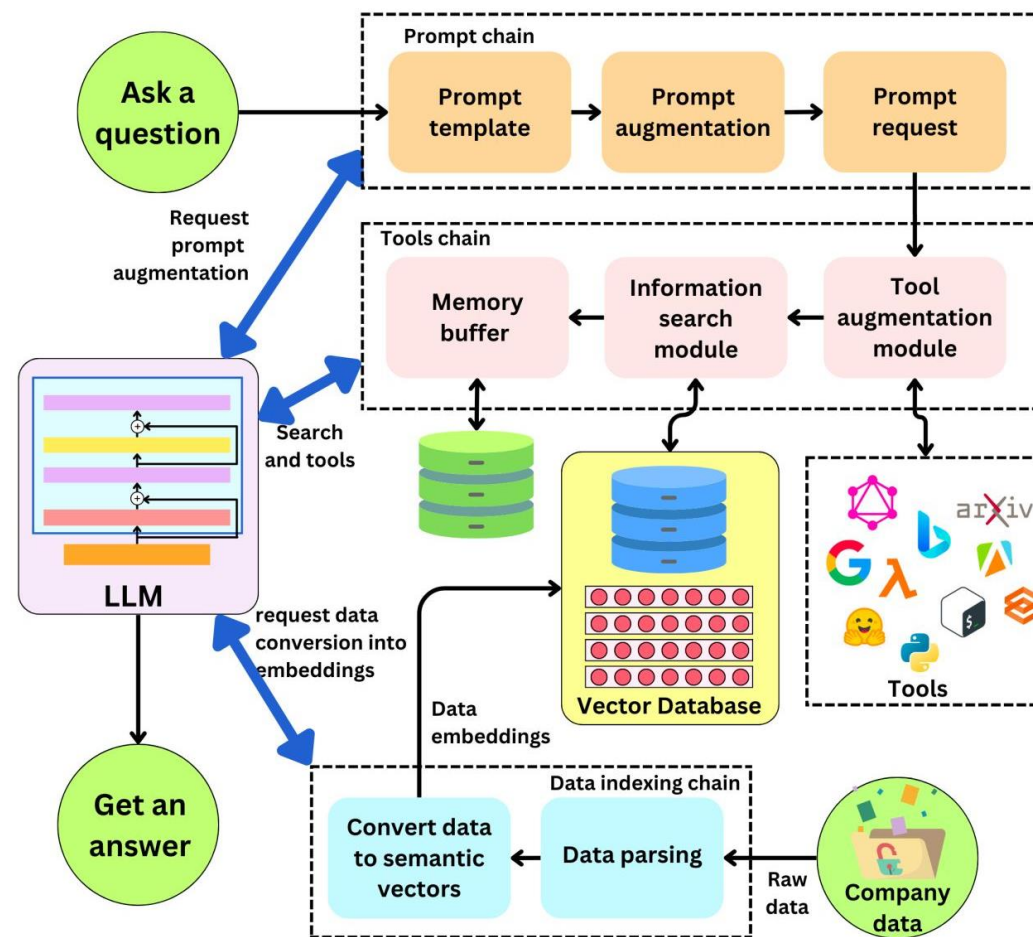
- *LangChain provides modular abstractions for the components necessary to work with language models.*
- *Use-Case Specific Chains: Chains can be thought of as assembling these components in particular ways in order to best accomplish a particular use case.*

Building Real Applications with LLMs

TheAiEdge.io Typical flow to interact with LLMs



Flow to build applications with LLMs



Agents

An agent uses an LLM to **choose** a sequence of actions to take, building a chain ***dynamically***.

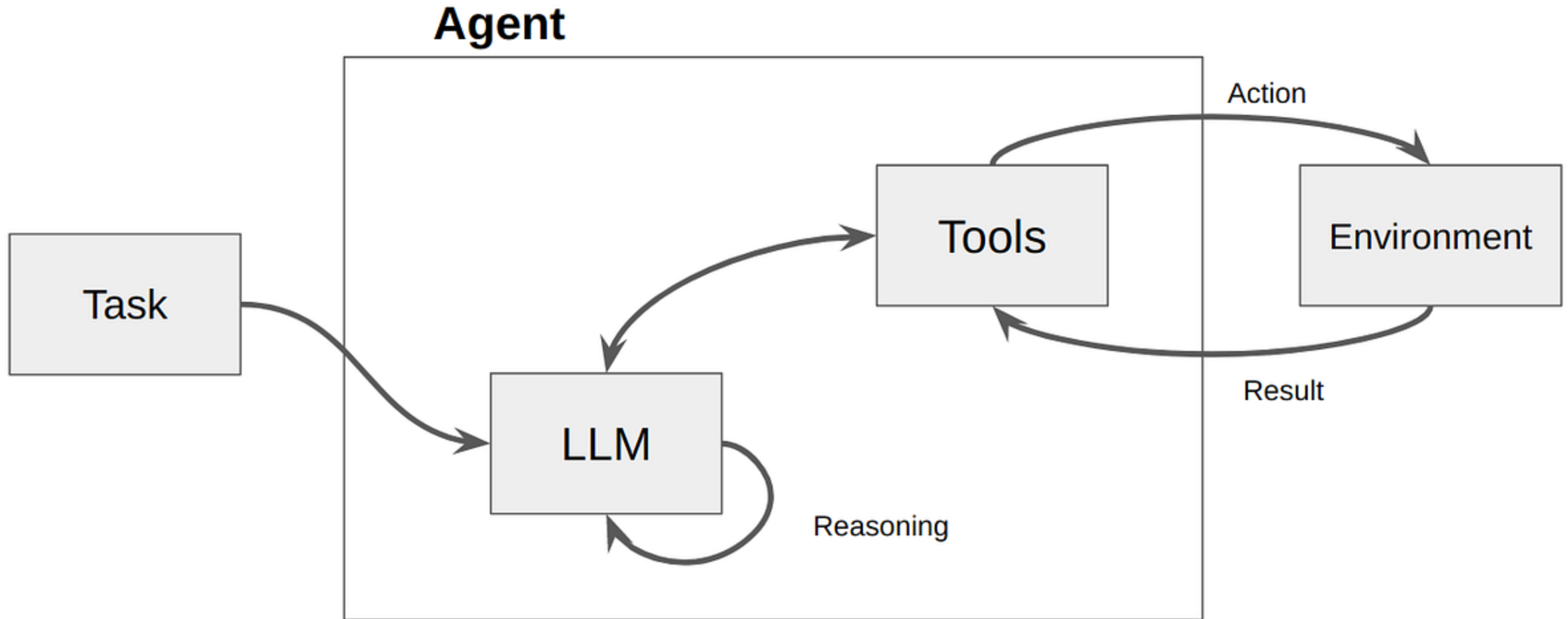
Some applications require an unknown chain of calls to LLMs/other tools. In these types of chains, there is a “agent” which has access to a suite of tools. Depending on the user input, the agent can then decide which, if any, of these tools to call. Chains are hardcoded, an agent reasons on what chain to execute.

Agent Components

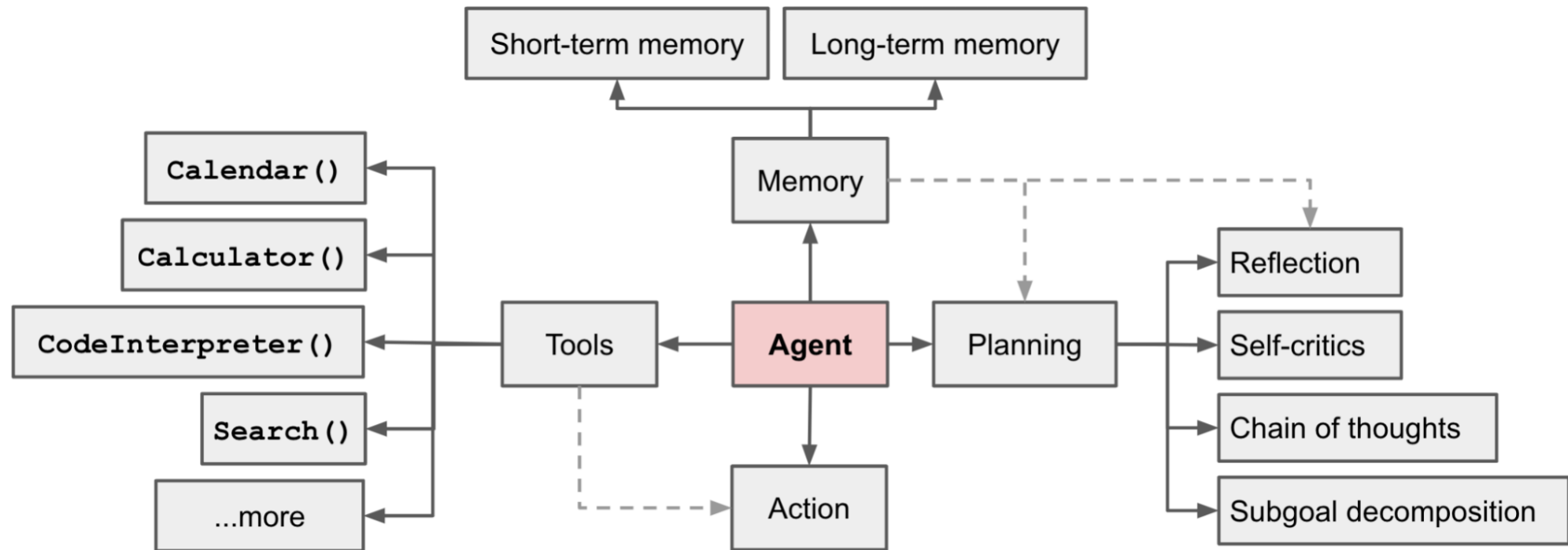
Components of an “agent”:

- **Agent** – Wraps a model, takes input, and produces an “action” to take.
- **Tools** – Functions that an agent executes, along with a description of when to use.
- **Memory** – The prior interactions with the agent, such as a “chat history”.
- **Agent executor** – combines the agent and tools to perform problem solving.

Agent Components



Agent Components



Use Cases for Chaining and Agents with LangChain

Personal Assistants

1. Prompt Template: Defines the personality of your personal assistant
2. Memory: Equip with short-term conversation retention
3. Tools: Differentiate your assistant by selecting specific capabilities
4. Agent: Design an efficient agent that understands and performs actions.
5. Agent Executor: Establish an environment where the agent can effectively utilize its tools.

Question Answering Over Docs

1. Transform data to a compatible format: Index
2. Ingestion process into a Vectorstore:
3. Load documents using a Document Loader
4. Split documents utilizing a Text Generation
5. Find pertinent documents in the index based on the query
6. Retrieve and return the generated result to the user

Querying Tabular Data

- Tabular data storage: csvs, excel sheets, SQL tables
- Utilize document loaders (e.g., CSVLoader) for loading text data in tabular formats
- Create an index for efficient data querying and interaction
- Direct interaction with numeric tabular data using a language model
- Chain for simple and small datasets
- Agents for complex multi-query processes with the Language Model

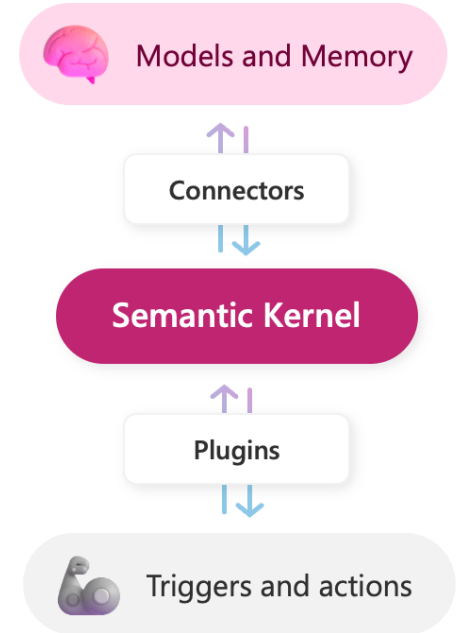
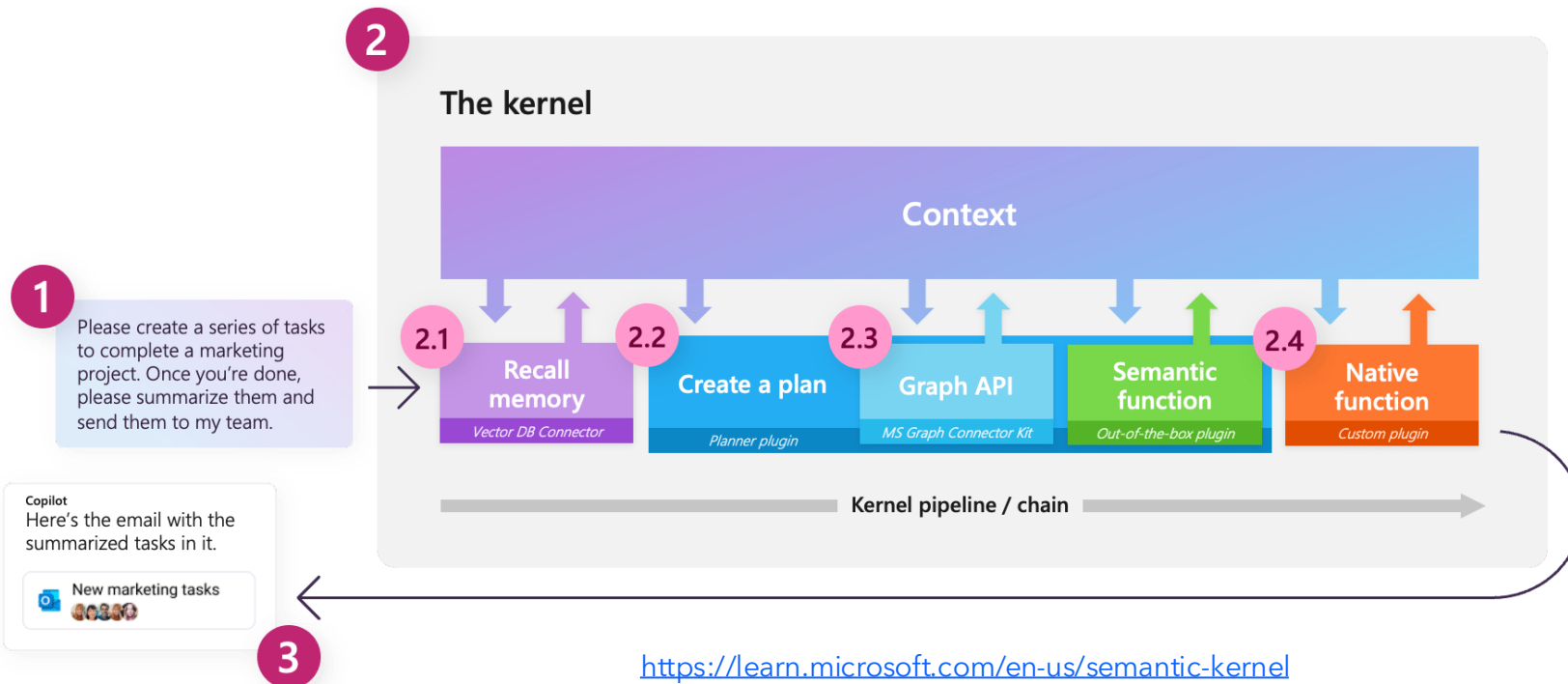
Interacting with APIs

- Integrate LLM with external APIs
- Retrieve context for LLM usage
- Interact with APIs through natural language
- Interface LLMs with external APIs
- Functions (e.g., OpenAI functions)
- LLM-generated interface utilizing API documentation



Semantic Kernel (SK)

- [Oversimplified] Microsoft's open-source competitor to LangChain
- Integrate GenAI models, memory, agents, and tools (terminology differs)
- Supported Languages: C#, Python, Java
- "With SK, you can leverage the same AI orchestration patterns that power Microsoft 365 Copilot and Bing in your own apps, while still leveraging your existing development skills and investments."



03 Chaining in Action

Chaining in Action

Conversation Bot

<https://huggingface.co/spaces/kimadams/ai-kit> - Prompt Builder Tab

Audio-to-Text:

- Separate audio from video, transcribe audio using Whisper-1 speech recognition

Embedding Query:

- Leverage embeddings for best answer primer (context)

Assistant Service:

- OpenAI ChatCompletion with prompt and context to push conversation along

AI Voice Service:

- Pass OpenAI response to 11Labs for voice response

The screenshot shows the 'Prompt Builder' tab of the 'ai-kit' space on Hugging Face Spaces. The interface is divided into several sections:

- Purpose:** Explains how generative AI can be used to create content and explore personas. It mentions translation options for language and voice options for AI voices.
- Directions:** Provides instructions on how to use the interface, including selecting a persona, language, and voice, and using the microphone to record a question.
- You can ask questions like:** A text input field with a placeholder question: "What benefits are provided? Is there a 401k match? How many vacation days can I expect each year?"
- Persona:** A section titled "What role would you like the system to play?" with radio button options: HR Expert (selected), Customer Service, Financial Expert, IT Expert, Insurance Agent, and Investigator.
- Translation:** A section titled "Which language would you like responses in?" with radio button options: English (selected), Spanish, and French.
- Check to enable sentiment analysis:** A checkbox for "Sentiment" (unchecked).
- Check to enable emotion detection:** A checkbox for "Emotion" (unchecked).
- Voice:** A section titled "Which voice would you like to use?" with radio button options: Sally (selected), Earl, Luke, Vin, Ebony, Matilda, and Serena. Below these are buttons for "Sasha" and "Adam".
- Use your microphone to record a question:** A section with a microphone icon and a button labeled "Ask a question". Below this is a checkbox for "Record from microphone" (checked).
- Buttons:** "Clear" and "Submit" buttons.
- OpenAI Communication Log:** A table with columns 1, 2, and 3, and a "New row" button.

Chaining in Action

Recording Analysis

<https://huggingface.co/spaces/kimadams/ai-kit> - Recording Analysis Tab

Audio-to-Text:

- Separate audio from videos, transcribe audio using Whisper-1 speech recognition

Summarization:

- Condense the transcribed text into clear and concise summaries

Topic Generation:

- Leverage topic modeling algorithms to extract meaningful topics from the summaries

The screenshot shows the 'Recording Analysis' tab in the Hugging Face Spaces AI Kit. The interface is divided into several sections:

- Purpose:** A paragraph explaining that the tool uses Generative AI to summarize and identify key concepts from videos. It details a three-step process: separating audio, transcribing it with OpenAI Whisper-1, summarizing with OpenAI Chat Completion, and extracting key concepts with OpenAI Chat Completion.
- Directions:** Instructions on how to use the tool, including uploading a video file and navigating the interface to find summaries and topics.
- Upload a Video(.mp4):** A section with a video player showing 'John's Voice' by Microsoft Teams. The video is dated 2023-07-28 17:45 UTC and is owned by Kim Adams.
- Summary:** A text box containing a promotional message for Slalom Build, highlighting their passion for technology and their commitment to creating transformative change.
- Topics:** A horizontal list of extracted topics, including 'Slalom Build', 'technology', 'transformative change', 'makers', 'planners', 'creatives', 'coders', 'better world', 'fearless approach', 'challenges', 'co-creating', 'cutting-edge products', 'experiences', 'future of impact', and 'join'.
- OpenAI Communication Log:** A table showing the interaction between the system and the assistant. The log includes the system's instructions, the transcribed audio, the assistant's summary, and the system's request for keywords.

The bottom of the interface features a 'New row' button and a 'New column' button.

Disambiguation: Chains, Tools, Agents

Chains

- Provide the ability to connect multiple LLM calls and additional steps in a deterministic workflow
- May or may not use tools or agents

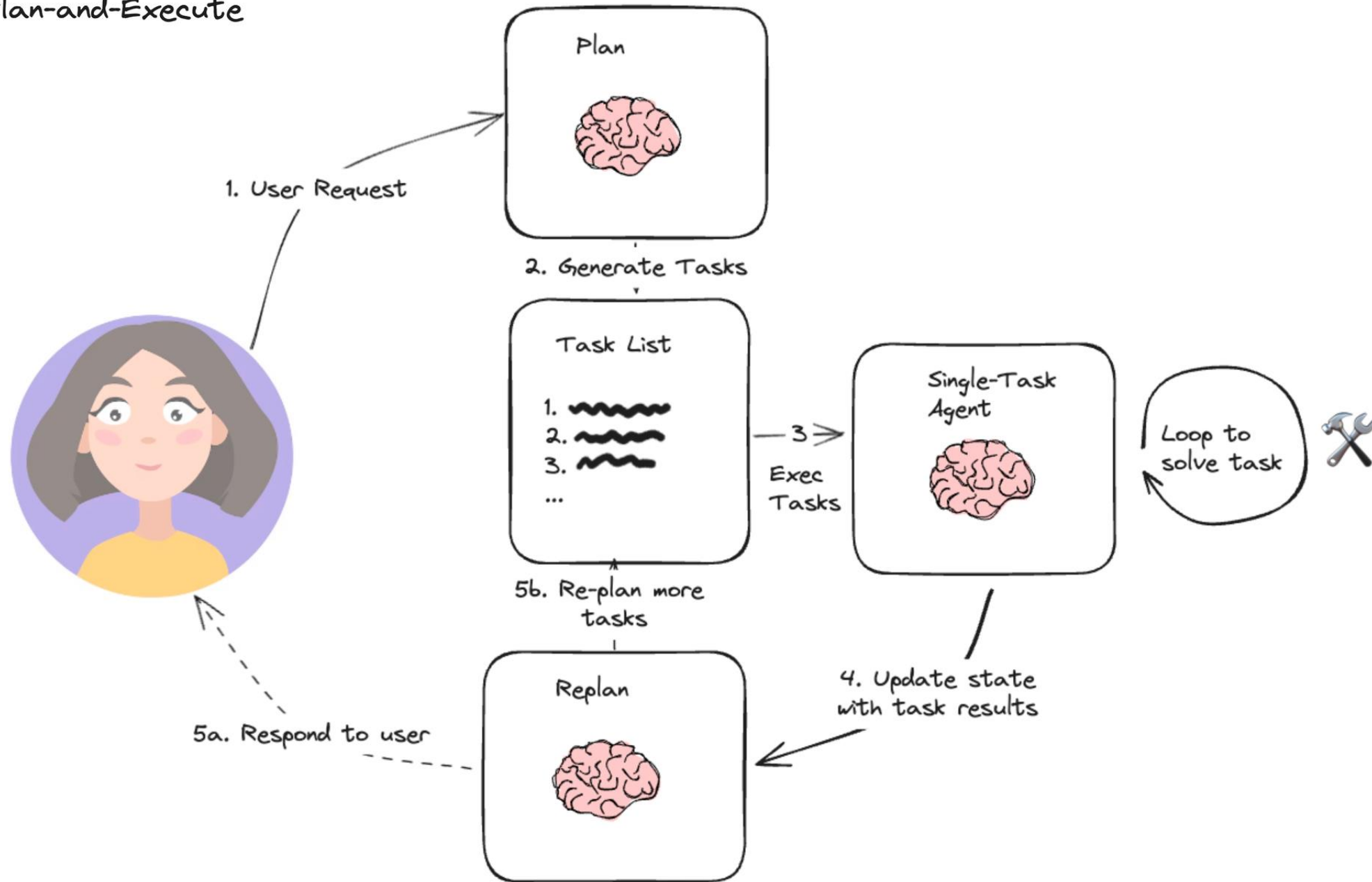
Tools

- Enhance LLM capabilities by providing the LLM with additional actions they can take to pull context
- Aka Function Calling

Agents

- Smarter tools that use LLMs to enhance results
- provide a natural language interface for a tool
- Autonomous Agents are smarter Agents that can self-evaluate and run independently from user input or create their own adhoc / non-deterministic workflows

04 Chains and Agents Example



Chain Example

```
def ask_sql_chain(question: str, verbose: bool = True):
    """
    This function runs a SQL query using a chain of operations against a local sample database.

    The database is a sample SQLite database called Chinook. The database contains information about
    a record store. The database has 11 tables: albums, employees, invoices, playlists, artists, genres, media_types,
    tracks, customers, invoice_items, playlist_track

    Parameters:
    question (str): The SQL query to be executed.
    verbose (bool): If True, the function will print additional debug information. Default is True.

    Returns:
    str: The result of the SQL query.
    """

    # Create an instance of SQLiteDatabase with the database URI and the tables to include.
    db = SQLiteDatabase.from_uri(
        database_uri=database_uri,
        include_tables=["Album", "Artist", "Track"],
        sample_rows_in_table_info=5,
    )

    llm = AzureChatOpenAI(temperature=0, verbose=True, deployment_name=deployment_name)

    # Create an instance of SQLiteDatabaseChain with the AzureChatOpenAI instance, the SQLiteDatabase instance,
    # and the verbose mode set.
    db_chain = SQLiteDatabaseChain.from_llm(
        llm,
        db,
        verbose=verbose,
    )

    # Run the SQL query and get the result.
    response = db_chain.run(question)

    # Return the result of the SQL query.
    return response

if __name__ == "__main__":
    q = "What is the name of the album that has the most tracks?"
    result = ask_sql_chain(q)
    print(result)
```

```
> Entering new SQLiteDatabaseChain chain...
What is the name of the album that has the most tracks?
SQLQuery:SELECT "Album"."Title", COUNT("Track"."TrackId") AS "NumTracks"
FROM "Album"
JOIN "Track" ON "Album"."AlbumId" = "Track"."AlbumId"
GROUP BY "Album"."AlbumId"
ORDER BY "NumTracks" DESC
LIMIT 1;
SQLResult: [('Greatest Hits', 57)]
Answer:Greatest Hits
> Finished chain.
Greatest Hits
```

https://bitbucket.org/slalom-consulting/chain-agent-example/src/main/examples/sql_chain.py

Agent Example

```
def ask_sql_agent(question: str, verbose: bool = False):
    """
    This function creates a sql agent. It can generate and execute one or many SQL queries to
    acquire the data needed to answer the given question against a local sample database.

    The database is a sample SQLite database called Chinook. The database contains information about
    a record store. The database has 11 tables: albums, employees, invoices, playlists, artists, genres, media_types,
    tracks, customers, invoice_items, playlist_track

    Parameters:
    question (str): Question asked in about the database
    verbose (bool): If True, the function will print additional debug information. Default is False.

    Returns:
    str: The result of the SQL query.
    """

    # Create an instance of SQLiteDatabase with the database URI and the tables to include.
    db = SQLiteDatabase.from_uri(
        database_uri=database_uri,
        include_tables=["Album", "Artist", "Track", "Customer"],
        sample_rows_in_table_info=5,
    )

    # Create an instance of AzureChatOpenAI with temperature set to 0, verbose mode enabled,
    # and the deployment name set.
    llm = AzureChatOpenAI(
        temperature=0, verbose=verbose, deployment_name=deployment_name
    )

    # Create a SQL agent with the AzureChatOpenAI instance, the SQLiteDatabase instance, and the verbose mode set.
    agent_executor = create_sql_agent(
        llm=llm,
        toolkit=SQLiteDatabaseToolkit(db=db, llm=llm),
        verbose=verbose,
        agent_type=AgentType.ZERO_SHOT_REACT_DESCRIPTION,
    )

    # Run the SQL query and get the result.
    agent_response = agent_executor.run(question)

    # Return the result of the SQL query.
    return agent_response

if __name__ == "__main__":
    q = "What are the names of the 5 artists with the most albums?"
    result = ask_sql_agent(q)
    print(result)
```

```
> Entering new AgentExecutor chain ...
Action: sql_db_list_tables
Action Input: ""
Observation: Album, Artist, Customer, Track
Thought:I can query the Album and Artist tables to find the names of the artists with the most albums. I should check the schema of these tables to see which columns I need to select.
Action: sql_db_schema
Action Input: Album, Artist
Observation:
CREATE TABLE "Album" (
    "AlbumId" INTEGER NOT NULL,
    "Title" NVARCHAR(160) NOT NULL,
    "ArtistId" INTEGER NOT NULL,
    PRIMARY KEY ("AlbumId"),
    FOREIGN KEY("ArtistId") REFERENCES "Artist" ("ArtistId")
)

/*
5 rows from Album table:
AlbumId Title      ArtistId
1  For Those About To Rock We Salute You    1
2  Balls to the Wall      2
3  Restless and Wild      2
4  Let There Be Rock      1
5  Big Ones               3
*/

CREATE TABLE "Artist" (
    "ArtistId" INTEGER NOT NULL,
    "Name" NVARCHAR(120),
    PRIMARY KEY ("ArtistId")
)

/*
5 rows from Artist table:
ArtistId Name
1  AC/DC
2  Accept
3  Aerosmith
4  Alanis Morissette
5  Alice In Chains
*/
Thought:I can use the ArtistId column from the Album table to join with the Artist table and get the names of the artists. I should order the results by the count of albums in descending order and limit the results to 5.
Action: sql_db_query
Action Input: "SELECT Artist.Name FROM Album JOIN Artist ON Album.ArtistId = Artist.ArtistId GROUP BY Artist.Name ORDER BY COUNT(Album.AlbumId) DESC LIMIT 5"
Observation: [('Iron Maiden',), ('Led Zeppelin',), ('Deep Purple',), ('U2',), ('Metallica',)]
Thought:I now know the final answer.
Final Answer: The names of the 5 artists with the most albums are Iron Maiden, Led Zeppelin, Deep Purple, U2, and Metallica.

> Finished chain.
The names of the 5 artists with the most albums are Iron Maiden, Led Zeppelin, Deep Purple, U2, and Metallica.
```

https://bitbucket.org/slalom-consulting/chain-agent-example/src/main/examples/sql_agent.py

05 Challenge Project



Challenge Project

Demo at Week 7 session

- If you can't make it, record your demo!
- Get creative
- Demonstrate one or more skills that you've learned
- Doesn't need to be polished or have a UI
 - Python notebooks are fine
 - Command line is fine
- Individuals or small teams
- Sign up with your team (or solo) → [Demo Sign-up Sheet](#)
 - (so that we can plan enough time for demos)

Sample Ideas:

- Soccer rules chat
- Generate a kitten picture based on the weather



Thank You



Resources

YouTube (Sam Witteveen):

- LangChain - Basics: https://www.youtube.com/watch?v=J_0qvRt4LNk
- LangChain - Tools & Chains: https://www.youtube.com/watch?v=hl2BY7yl_Ac
- LangChain - Conversations with Memory: https://www.youtube.com/watch?v=X550Zbz_ROE
- LangChain - Agents: <https://www.youtube.com/watch?app=desktop&v=ziu87EXZVUE>

LangChain Memory

Assists in providing “memory” to LLM interactions

ConversationBufferMemory

Buffers the entire conversation which can be retrieved programmatically. Best for limited number of interactions given token restrictions.

ConversationBufferWindowMemory

Similar to ConversationBufferMemory but limits number of previous interactions. May run into token issues if configured for high number of interactions.

ConversationSummaryMemory

Instead of taking the last X interactions, continues to summarize the interactions so far. Takes more tokens initially, but less as the conversation extends.

ConversationKGMemory

Builds a Knowledge Graph as the conversation is happening.

Entity Memory

Extracts information keyed by entities

Custom Memory

Combines different types of memories

Why Memory?

- *Passing the previous context for a long-running conversation may run into token issues*
- *People expect Bots to remember long conversations*
- *Co-reference resolution during natural conversation. Example: “My friend Andrew...” later referencing Andrew by “him”*