

# IERG5050 AI Foundation Models, Systems & Applications Spring 2025

## Agentic AI Applications

Prof. Wing C. Lau

[wclau@ie.cuhk.edu.hk](mailto:wclau@ie.cuhk.edu.hk)

<http://www.ie.cuhk.edu.hk/wclau>

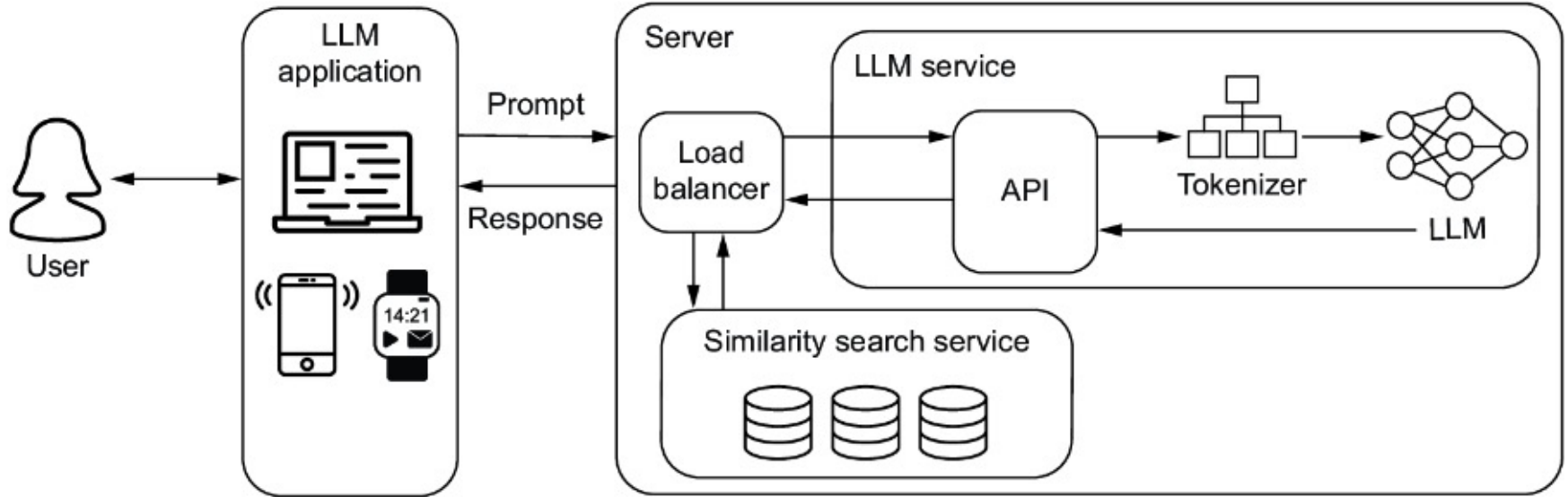
# Acknowledgements

Many of the slides in this lecture are adapted from the sources below. Copyrights belong to the original authors.

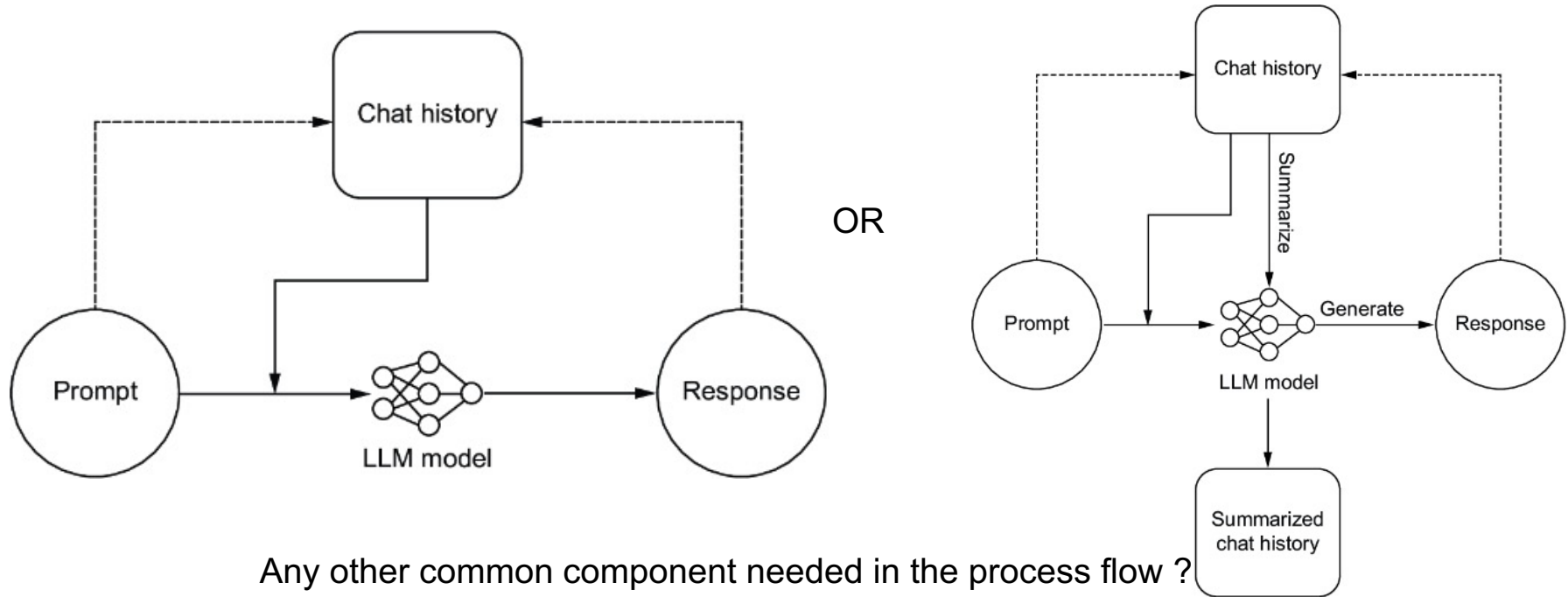
- LangChain for LLM Application Development – a short course by Harrison Chase, Andrew Ng, <https://www.deeplearning.ai/short-courses/langchain-for-llm-application-development/>
- LangChain Chat with your Data – a short course by Harrison Chase, <https://www.deeplearning.ai/short-courses/langchain-chat-with-your-data/>
- AI Agents in LangGraph– a short course by Harrison Chase, Rotem Weiss, <https://www.deeplearning.ai/short-courses/ai-agents-in-langgraph/>
- Building LLM Powered Applications By Valentina Alto, Packt Publishing, May 2024
- AI Agents in Action, by Michael Lanham, Feb 2025, Manning Publications
- LLMs in Production by Christopher Brousseau, Matthew Sharp, Jan 2025, Manning Publications
- Shunyu Yao, “LLM agents: brief history and overview,” Guest Lecture for UC Berkeley MOOC on Large Language Model Agents, Fall 2024, [https://rdi.berkeley.edu/llm-agents/assets/llm\\_agent\\_history.pdf](https://rdi.berkeley.edu/llm-agents/assets/llm_agent_history.pdf), <https://www.youtube.com/watch?v=RM6ZArd2nVc>
- Chi Wang, “Agentic AI Frameworks & AutoGen,” Guest Lecture for UC Berkeley MOOC on Large Language Model Agents, Fall 2024, <https://rdi.berkeley.edu/llm-agents-mooc/slides/autogen.pdf>, <https://www.youtube.com/watch?v=OQdImCMSOo4>
- Jerry Liu, “Building a Multimodal Knowledge Assistant (LlamaIndex),” Guest Lecture for UC Berkeley MOOC on Large Language Model Agents, Fall 2024, <https://rdi.berkeley.edu/llm-agents-mooc/slides/MKA.pdf>, <https://www.youtube.com/watch?v=OQdImCMSOo4>
- Ruslan Salakhutdinov, “Multimodal Autonomous AI Agents,” Guest Lecture for UC Berkeley MOOC on Advanced Large Language Model Agents, Spring 2025, <https://rdi.berkeley.edu/adv-llm-agents/slides/ruslan-multimodal.pdf>, <https://www.youtube.com/live/RPINOYM12RU>
- Nicolas Chapados, “AI Agents for Enterprise Workflows,” Guest Lecture for UC Berkeley MOOC on Large Language Model Agents, Fall 2024, <https://rdi.berkeley.edu/llm-agents/assets/agentworkflows.pdf>, <https://www.youtube.com/live/-yf-e-9FvOc>
- Ben Mann, “Measuring Agent capabilities and Anthropic’s RSP,” Guest Lecture for UC Berkeley MOOC on Large Language Model Agents, Fall 2024, <https://rdi.berkeley.edu/llm-agents/assets/antrsp.pdf>, <https://www.youtube.com/live/6v2AnWol7oo>
- Caiming Xiong, “Multimodal Agents – from Perception to Action,” Guest Lecture for UC Berkeley MOOC on Advanced Large Language Model Agents, Spring 2025, [https://rdi.berkeley.edu/llm-agents-mooc/slides/Multimodal\\_Agent\\_caiming.pdf](https://rdi.berkeley.edu/llm-agents-mooc/slides/Multimodal_Agent_caiming.pdf), [https://www.youtube.com/live/n\\_\\_Tim8K2LY](https://www.youtube.com/live/n__Tim8K2LY)
- Omar Khattab, “Compound AI Systems & the DSPy Framework,” Guest Lecture for UC Berkeley MOOC on Large Language Model Agents, Fall 2024, <https://rdi.berkeley.edu/llm-agents-mooc/slides/dspy Lec.pdf>, <https://www.youtube.com/live/JEMYuzrKLUw>
- Graham Neubig, “Agents for Software Development,” Guest Lecture for UC Berkeley MOOC on Large Language Model Agents, Fall 2024, <https://rdi.berkeley.edu/llm-agents-mooc/slides/neubig24softwareagents.pdf>, <https://www.youtube.com/live/f9L9Fkg-8Kd>



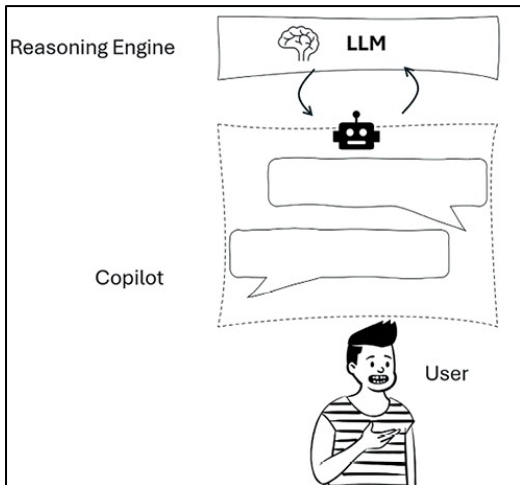
# An LLM-powered Application



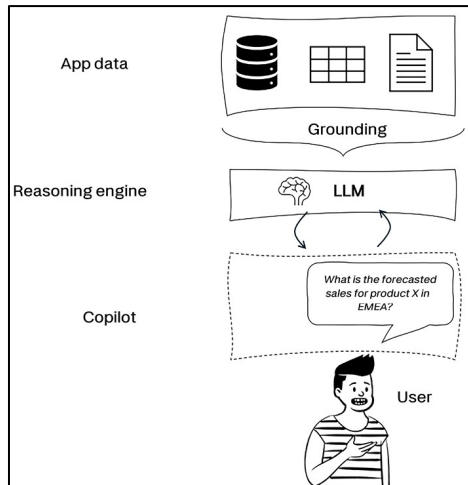
# Process-Flow of a Sample LLM-powered Application: a Minimalist ChatBot



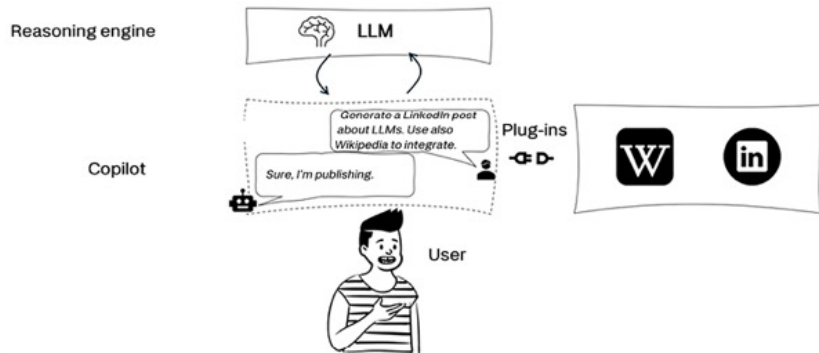
# Copilot (aka AI Assistant): An LLM-powered Application



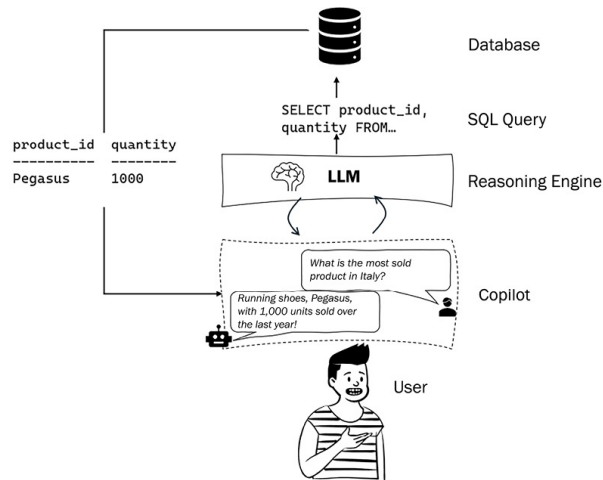
An LLM-powered Copilot



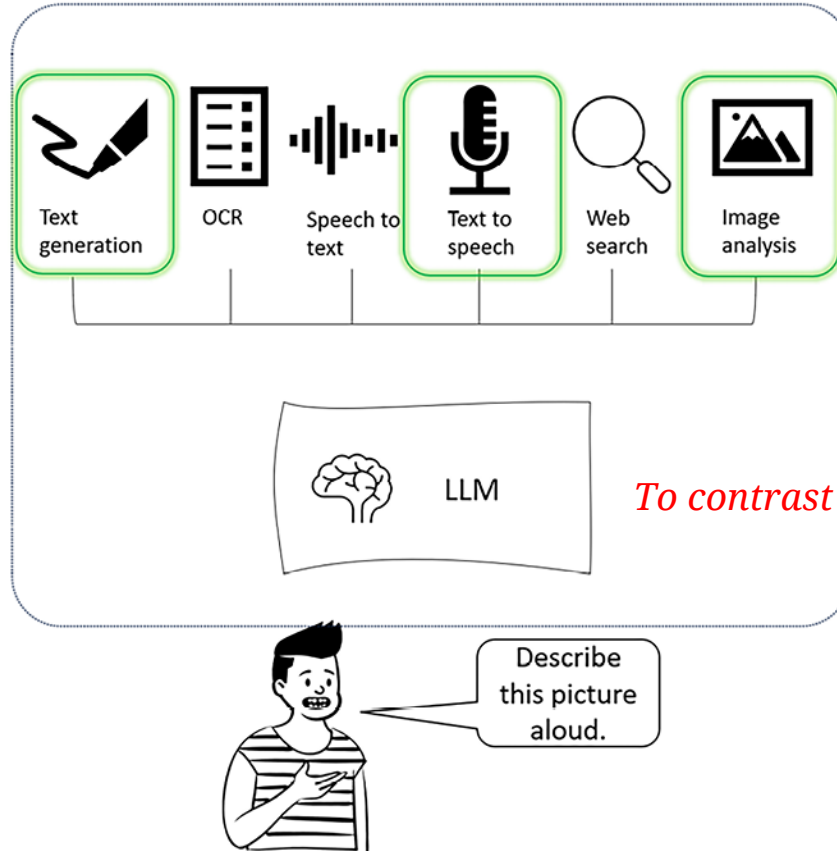
A Copilot w/ grounding



A conversational UI to reduce the gap between the user and the database or ext. tools

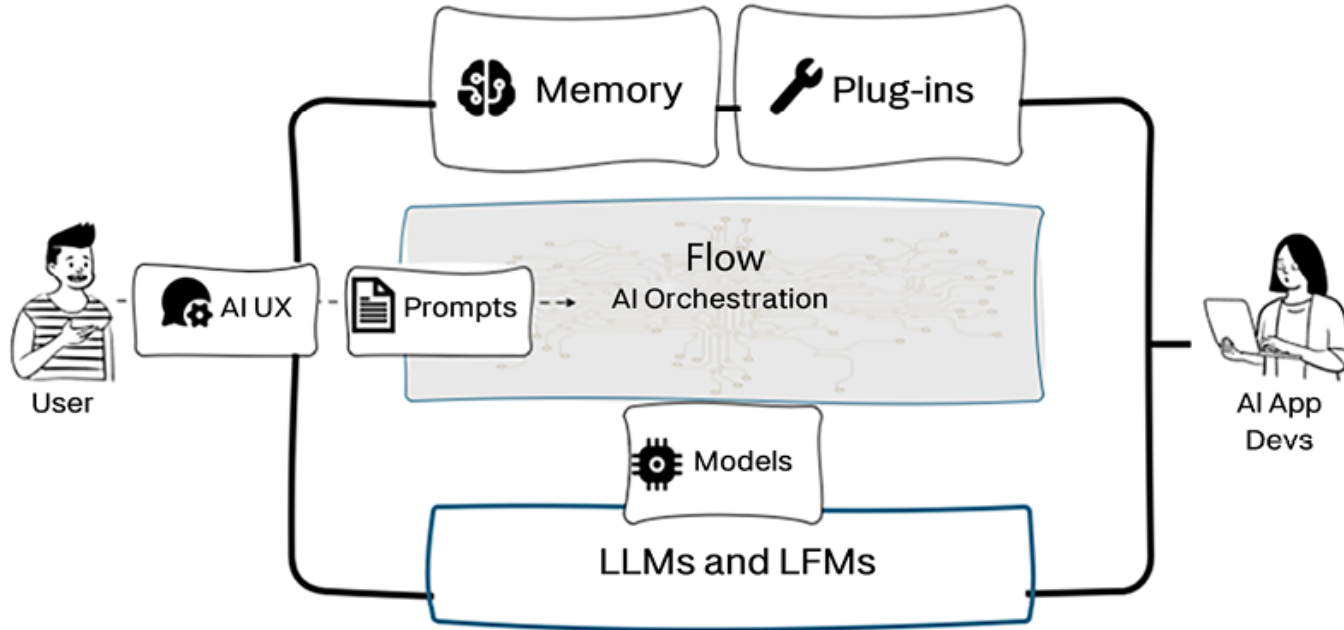


# A Multimodal Application via an LLM w/ Single-modal Tools



*To contrast with a Multimodal-NATIVE LLM (later)*

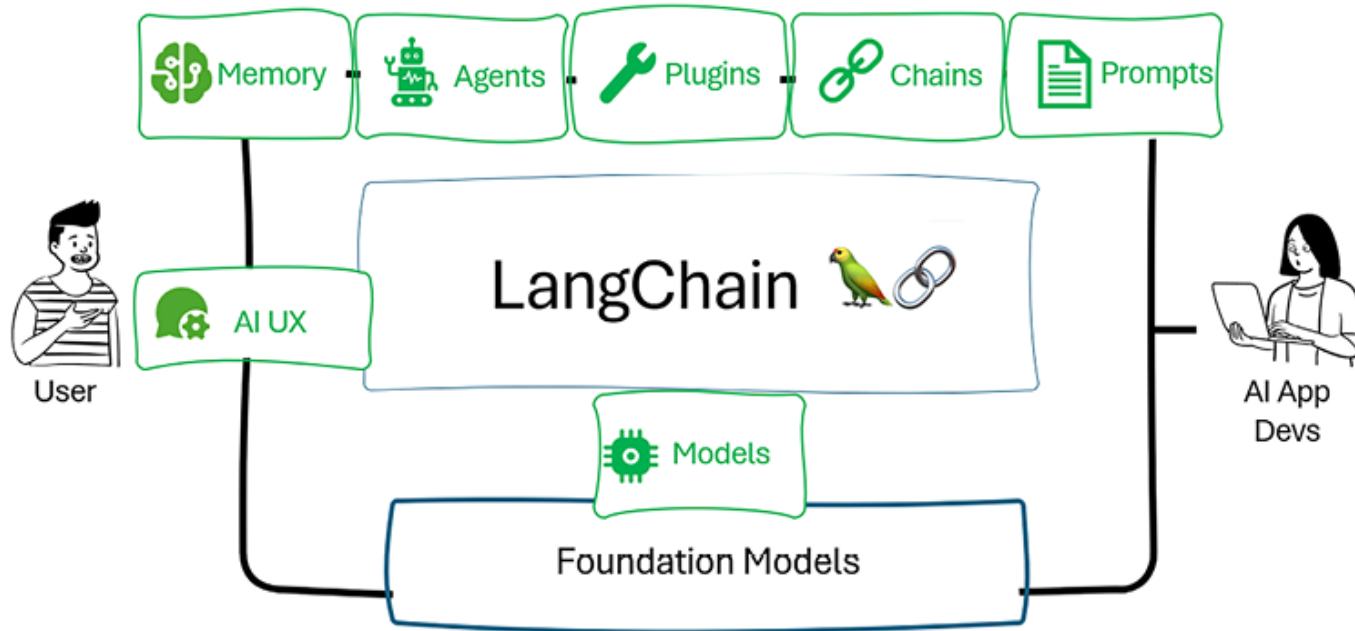
# High-level Architecture of an LLM-powered Application



# LLM-application/ Agent Development Frameworks

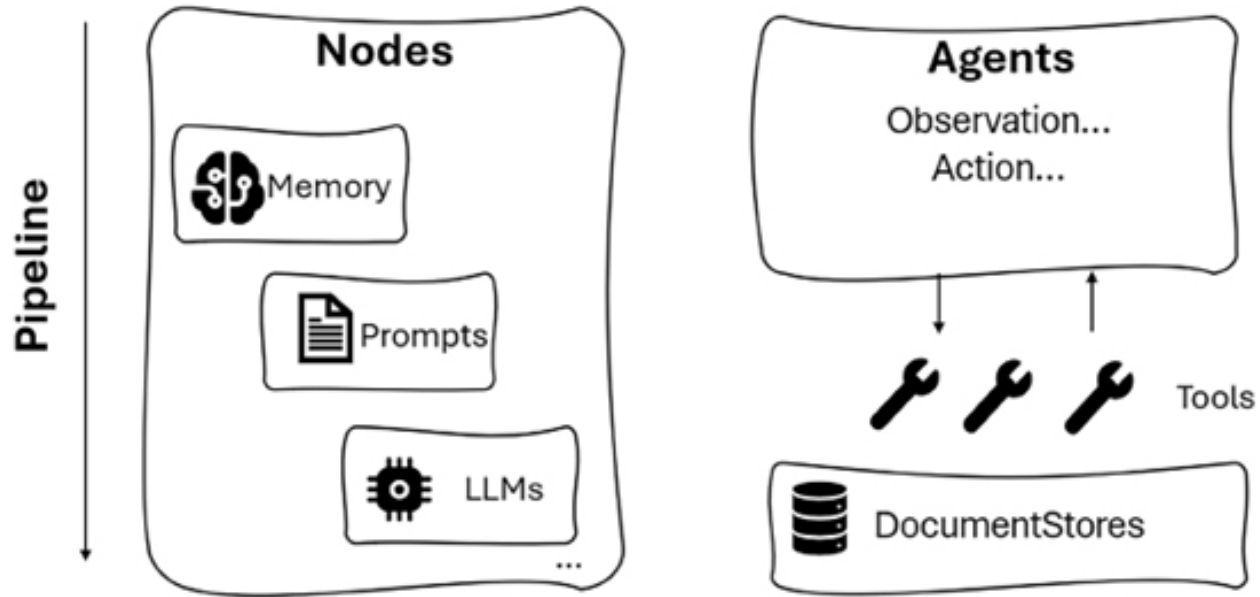
Framework	Best For	Deployment	Scalability	Current Adoption
LangChain & LangGraph	Enterprise AI apps, multi-agent workflows	Self-hosted, Cloud	High – widely used in production	Most adopted
Autogen & Semantic Kernel	Microsoft ecosystem, scalable AI apps	Azure, Self-hosted	High – enterprise-scale	Growing in enterprises
LlamaIndex	AI-powered search, RAG	Cloud, Self-hosted	High – efficient data processing	Strong in AI search
AutoGPT	No-code AI automation, continuous agents	Cloud-first, some self-hosted options	Medium – automation focused	Popular for prototyping
CrewAI	Workflow-based multi-agent AI apps	Cloud, Self-hosted	Medium – still early-stage	Fast-growing
PydanticAI	FastAPI & Pydantic-based AI apps	Self-hosted	Low – lightweight framework	Niche adoption
Spring AI	Java-based AI applications	Self-hosted, Cloud	Medium – depends on Java ecosystem	Popular in Java world
Haystack	LLM-powered search & RAG	Self-hosted, Cloud	High – built for production	Strong in search applications

# LangChain: a framework for building LLM-based Applications



Components of LangChain

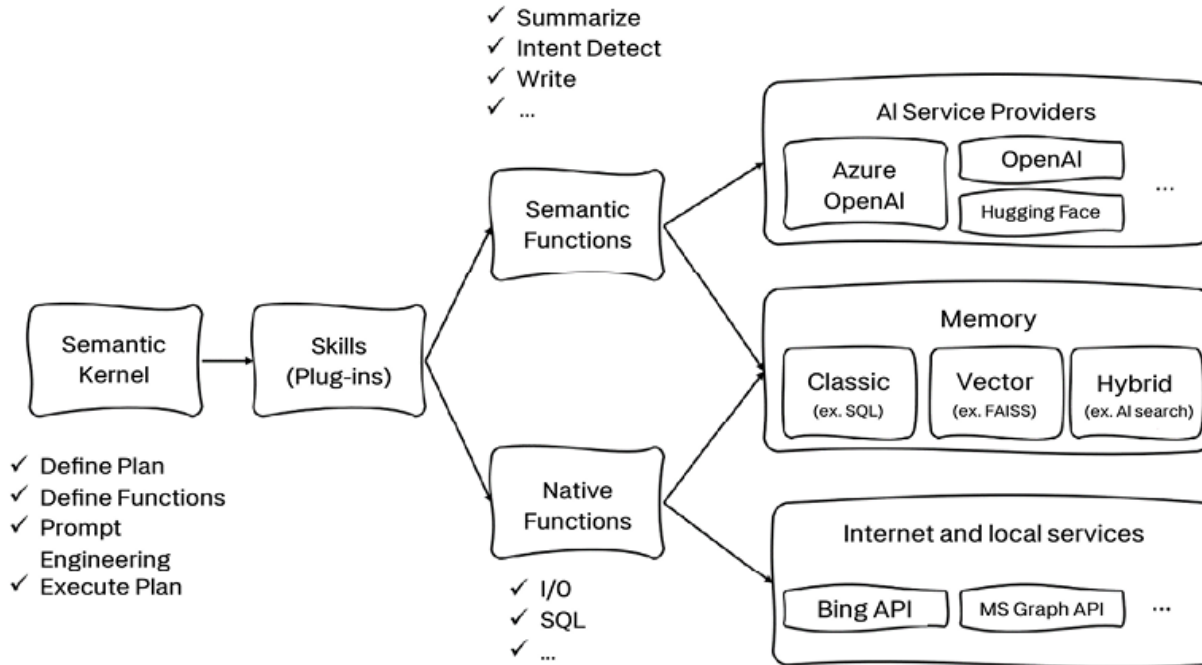
# Haystack: a framework for building LLM-based Applications



Haystack's components



# Semantic Kernel: a framework for building LLM-based Applications



## Anatomy of Semantic Kernel

# LangChain vs. Haystack vs. Semantic Kernel

Feature	LangChain	Haystack	Semantic Kernel
<b>LLM support</b>	Proprietary and open-source	Proprietary and open source	Proprietary and open source
<b>Supported languages</b>	Python and JS/TS	Python	C#, Java, and Python
<b>Process orchestration</b>	Chains	Pipelines of nodes	Pipelines of functions
<b>Deployment</b>	No REST API	REST API	No REST API
Feature	LangChain	Haystack	Semantic Kernel

*More on Semantic Kernel as well as other Agentic AI frameworks later...*

Table 2.1: Comparisons among the three AI orchestrators

# Developing LLM-based Applications with LangChain



# LangChain Overview

Open-source development framework for LLM applications

Python and JavaScript (TypeScript) packages

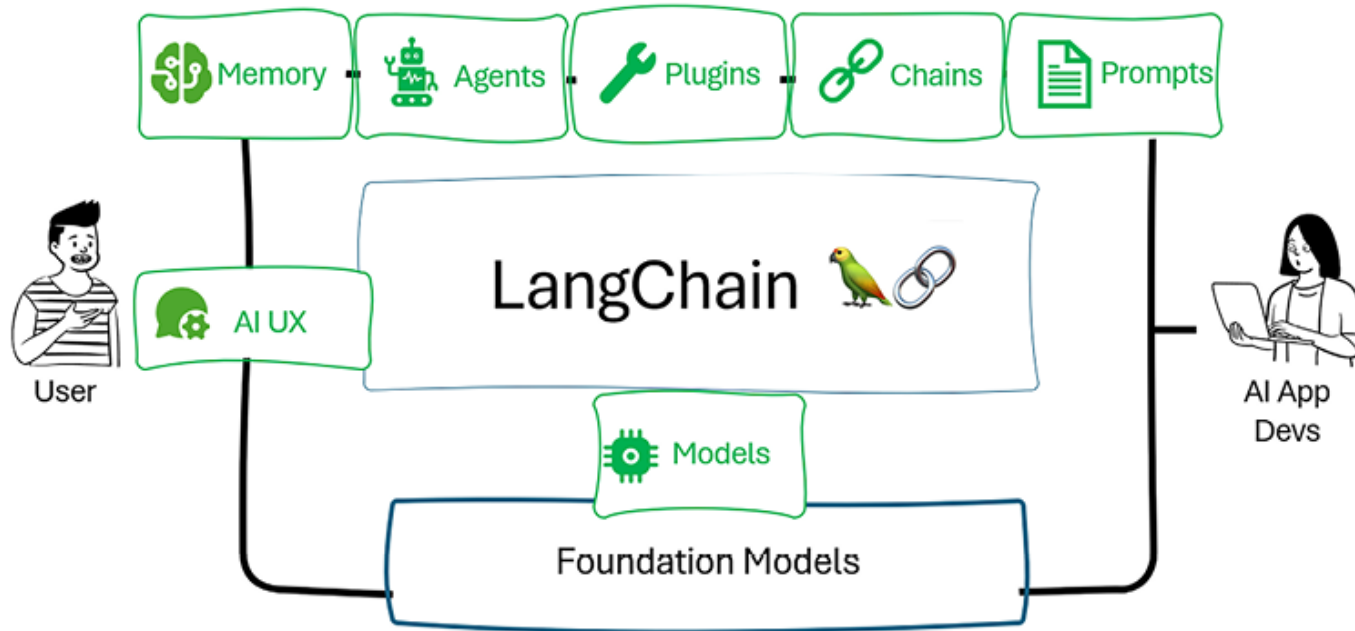
Focused on composition and modularity

Key value adds:

1. Modular components
2. Use cases: Common ways to combine components

*Source:* LangChain for LLM Application Development – a short course by Harrison Chase, Andrew Ng,  
<https://www.deeplearning.ai/short-courses/langchain-for-llm-application-development/>

# LangChain: a framework for building LLM-based Applications



Components of LangChain

# Components of LangChain

- **Models**

- LLMs: 20+ integrations
- Chat Models
- Text Embedding Models: 10+ integrations

- **Prompts**

- Prompt Templates
- Output Parsers: 5+ implementations
  - Retry/fixing logic
- Example Selectors: 5+ implementations

- **Indexes**

- Document Loaders: 50+ implementations
- Text Splitters: 10+ implementations
- Vector stores: 10+ integrations
- Retrievers: 5+ integrations/implementations

- **Chains**

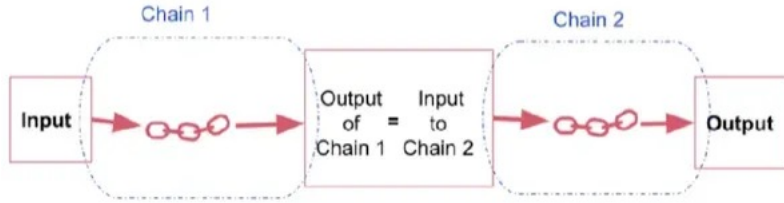
- Prompt + LLM + Output parsing
- Can be used as building blocks for longer chains
- More application specific chains: 20+ types

- **Agents**

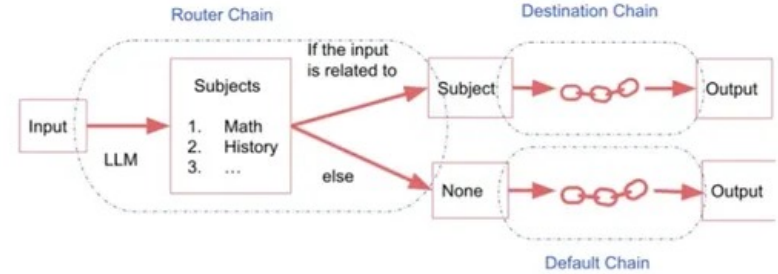
- Agent Types: 5+ types
  - Algorithms for getting LLMs to use tools
- Agent Toolkits: 10+ implementations
  - Agents armed with specific tools for a specific application

# Process Flow Orchestration with “Chains”

Simple Sequential Chain: One i/p, One o/p

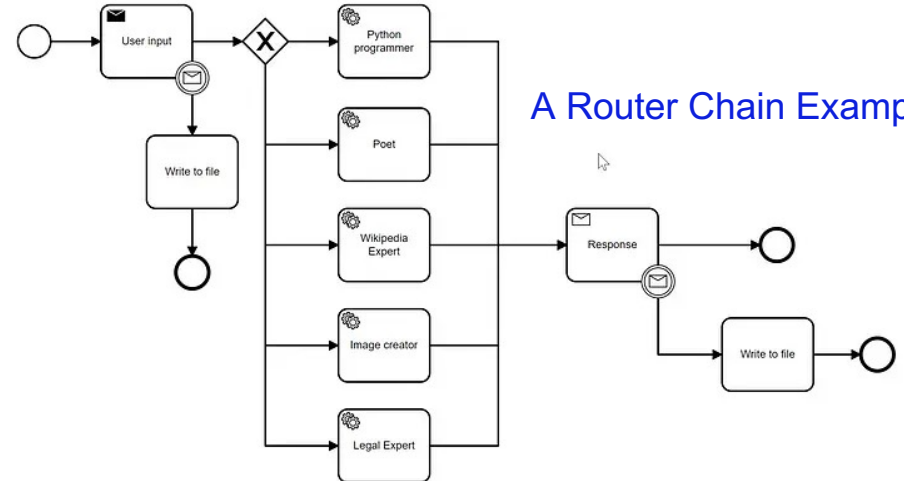
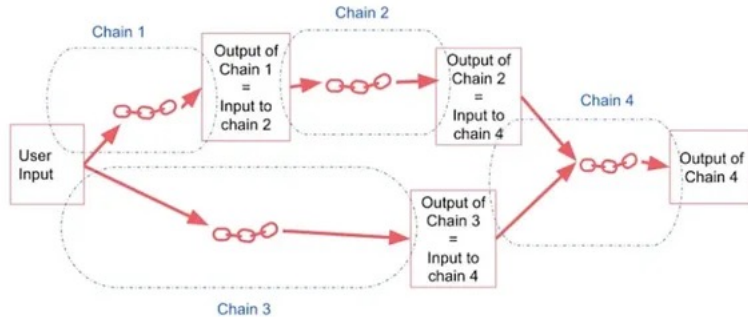


Router Chain:



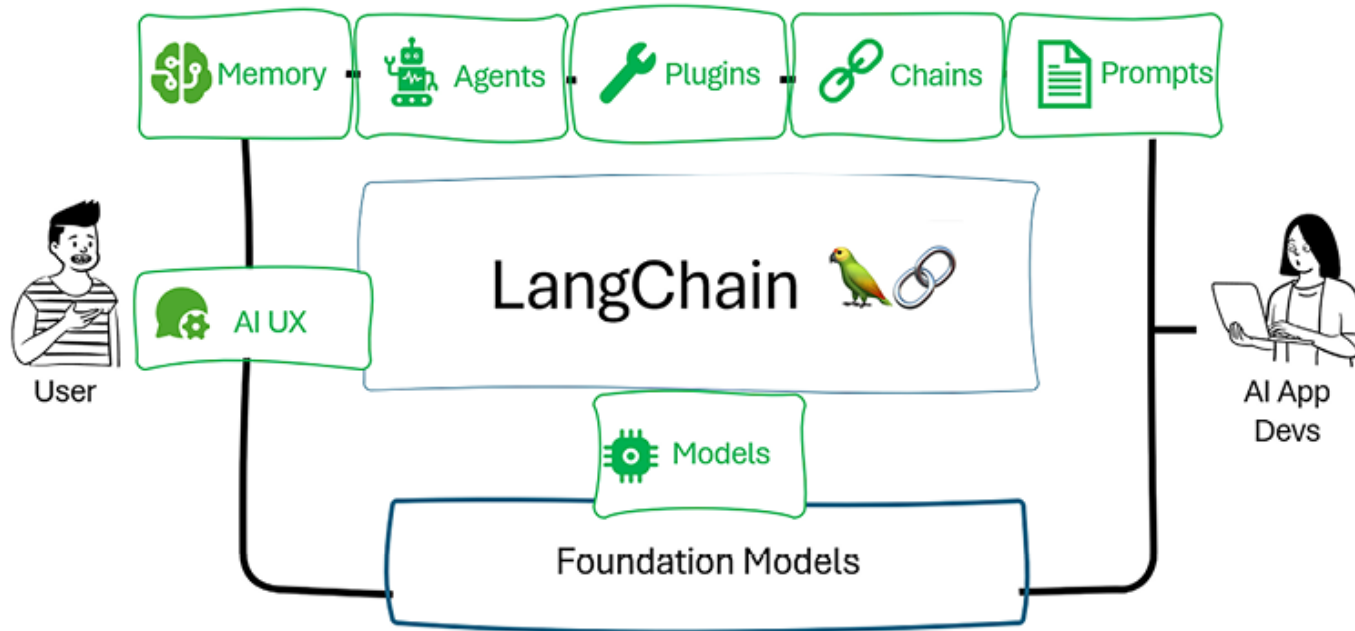
A Chain is sequence of calls. Those calls can be to LLMs, External Tools, or Data Processing steps

Sequential Chain: Multiple i/p, Multiple o/p



A Router Chain Example

# LangChain: a framework for building LLM-based Applications



Components of LangChain



# LangChain Memory Types

## ConversationBufferMemory

- This memory allows for storing of messages and then extracts the messages in a variable.

## ConversationBufferWindowMemory

- This memory keeps a list of the interactions of the conversation over time. It only uses the last K interactions.

## ConversationTokenBufferMemory

- This memory keeps a buffer of recent interactions in memory, and uses token length rather than number of interactions to determine when to flush interactions.

## ConversationSummaryMemory

- This memory creates a summary of the conversation over time.

# LangChain Memory Types

## Vector data memory

- Stores text (from conversation or elsewhere) in a vector database and retrieves the most relevant blocks of text.

## Entity memories

- Using an LLM, it remembers details about specific entities.

You can also use multiple memories at one time.

E.g., Conversation memory + Entity memory to recall individuals.

You can also store the conversation in a conventional database (such as key-value store or SQL)

# Document Loading with LangChain

- Loaders deal with the specifics of accessing and converting data

- o Accessing

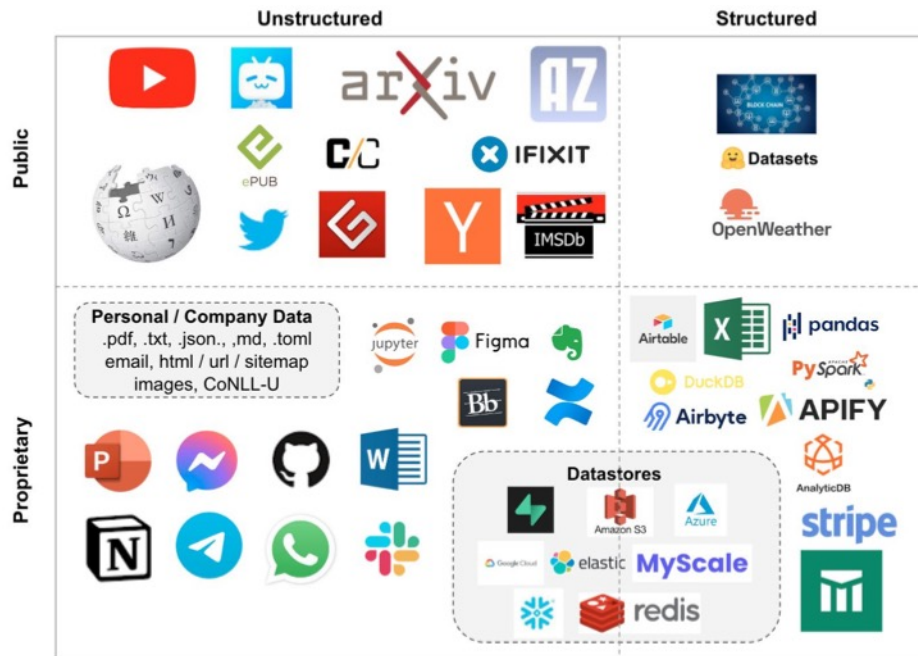
- Web Sites
- Data Bases
- YouTube
- arXiv
- ...

- o Data Types

- PDF
- HTML
- JSON
- Word, PowerPoint...

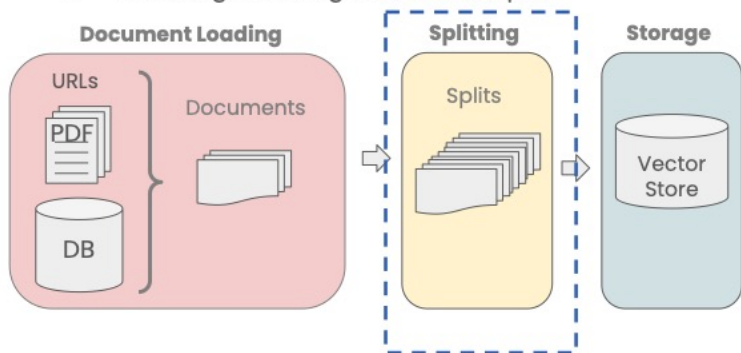
- Returns a list of `Document` objects:

```
[
Document(page_content='MachineLearning-Lecture01 \nInstructor (Andrew Ng): Okay.
Good morning. Welcome to CS229....',
metadata={'source': 'docs/cs229_lectures/MachineLearning-Lecture01.pdf', 'page': 0})
...
Document(page_content='[End of Audio] \nDuration: 69 minutes ',
metadata={'source': 'docs/cs229_lectures/MachineLearning-Lecture01.pdf', 'page': 21})
]
```



# Document Splitting with LangChain

- Splitting Documents into smaller chunks
  - Retaining meaningful relationships!



...  
on this model. The Toyota Camry has a head-snapping  
80 HP and an eight-speed automatic transmission that will

...

**Chunk 1:** on this model. The Toyota Camry has a head-snapping

**Chunk 2:** 80 HP and an eight-speed automatic transmission that will

**Question:** What are the specifications on the Camry?

## Example Splitter

```
langchain.text_splitter.CharacterTextSplitter(  
    separator: str = "\n\n"  
    chunk_size=4000,  
    chunk_overlap=200,  
    length_function=<builtin function len>,  
)
```

Methods:

**create\_documents()** - Create documents from a list of texts.

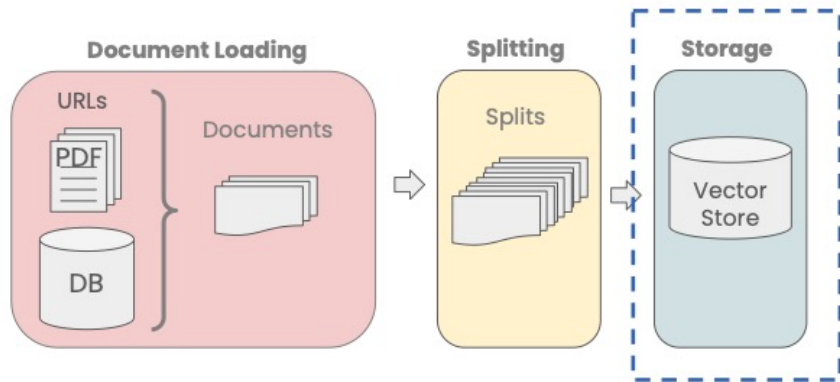
**split\_documents()** - Split documents.

# Different Types of Splitters in LangChain

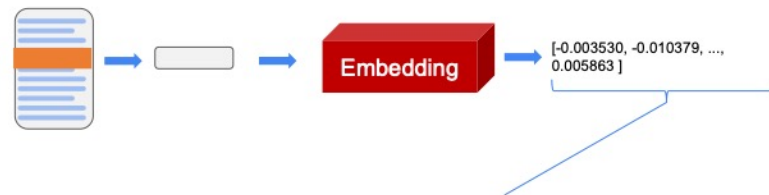
langchain.text\_splitter.

- **CharacterTextSplitter()**- Implementation of splitting text that looks at characters.
- **MarkdownHeaderTextSplitter()** - Implementation of splitting markdown files based on specified headers.
- **TokenTextSplitter()** - Implementation of splitting text that looks at tokens.
- **SentenceTransformersTokenTextSplitter()** - Implementation of splitting text that looks at tokens.
- **RecursiveCharacterTextSplitter()** - Implementation of splitting text that looks at characters. Recursively tries to split by different characters to find one that works.
- **Language()** – for CPP, Python, Ruby, Markdown etc
- **NLTKTextSplitter()** - Implementation of splitting text that looks at sentences using NLTK (Natural Language Tool Kit)
- **SpacyTextSplitter()** - Implementation of splitting text that looks at sentences using Spacy

# Support of Vector Store in LangChain

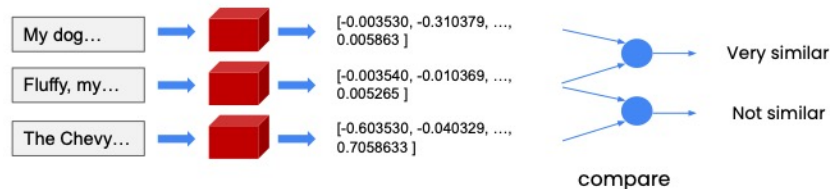


## Embeddings

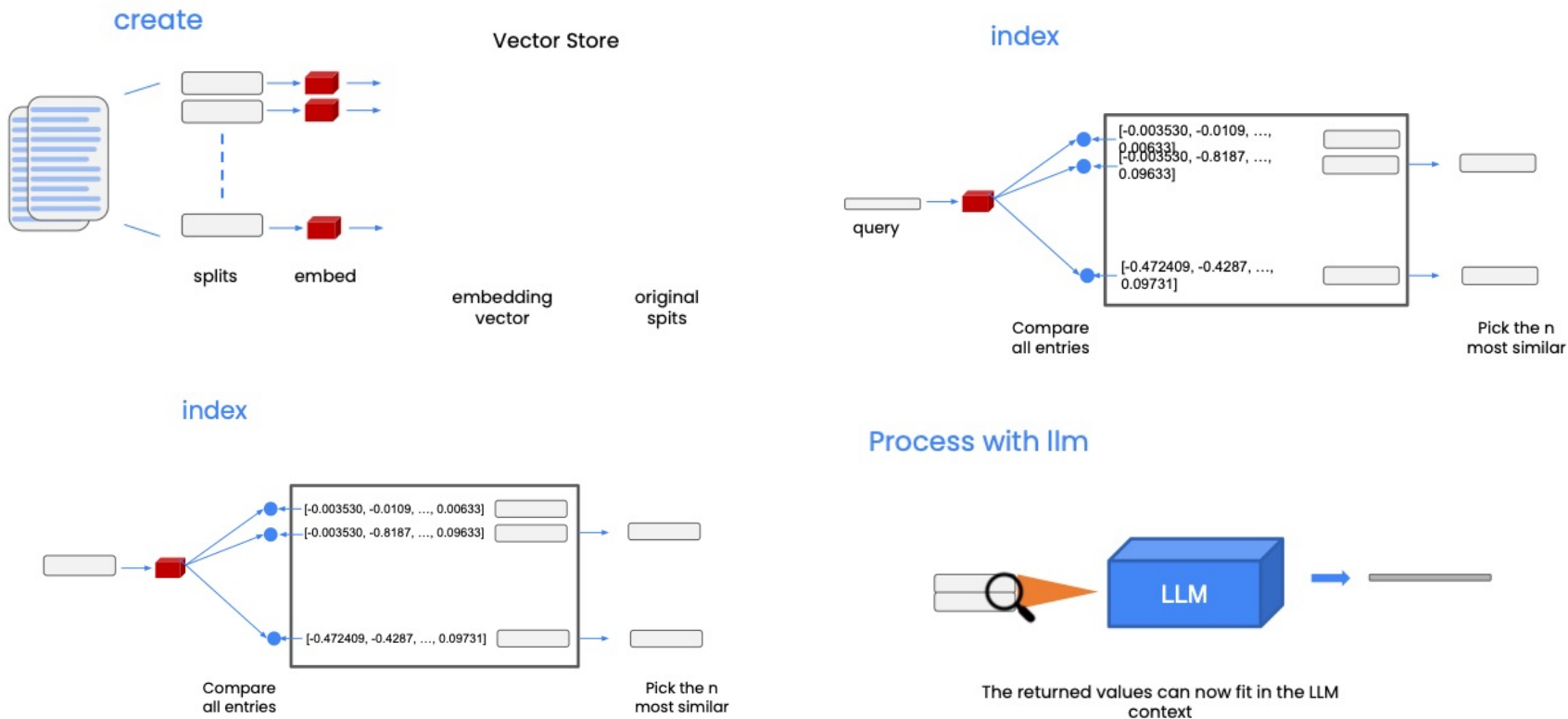


- Embedding vector captures content/meaning
- Text with similar content will have similar vectors

- 1) My dog Rover likes to chase squirrels.
- 2) Fluffy, my cat, refuses to eat from a can.
- 3) The Chevy Bolt accelerates to 60 mph in 6.7 seconds.

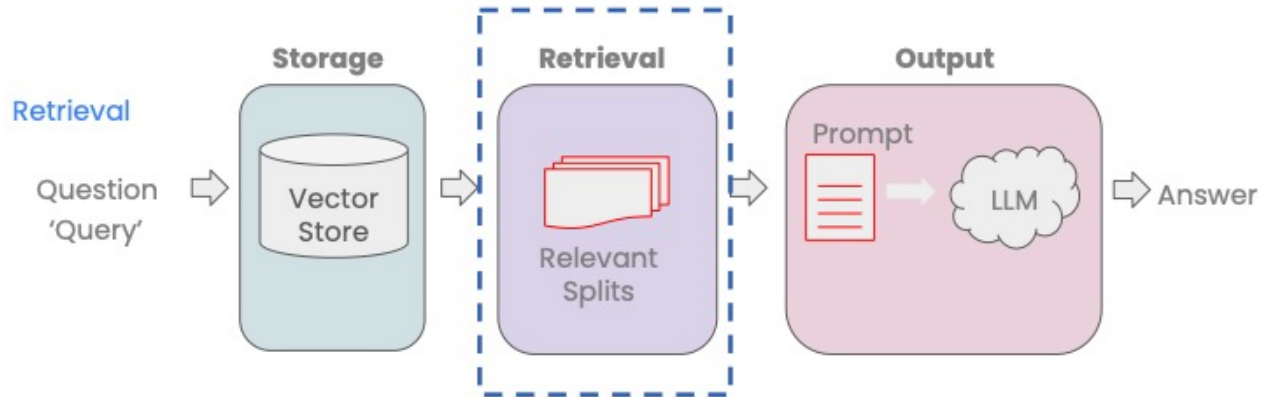


# Basic RAG support w/ Vector Store in LangChain





# Retrieval Support in LangChain



- Accessing/indexing the data in the vector store
  - Basic semantic similarity
  - Maximum marginal relevance
  - Including Metadata
- LLM Aided Retrieval

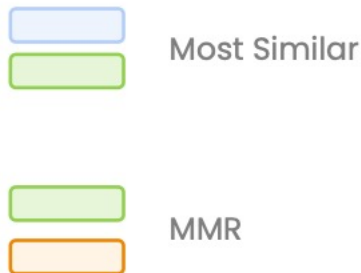
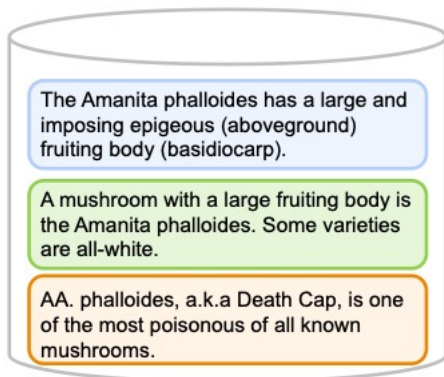


# Retrieval based on Maximum Marginal Relevance (MMR)

- You may not always want to choose the most similar responses

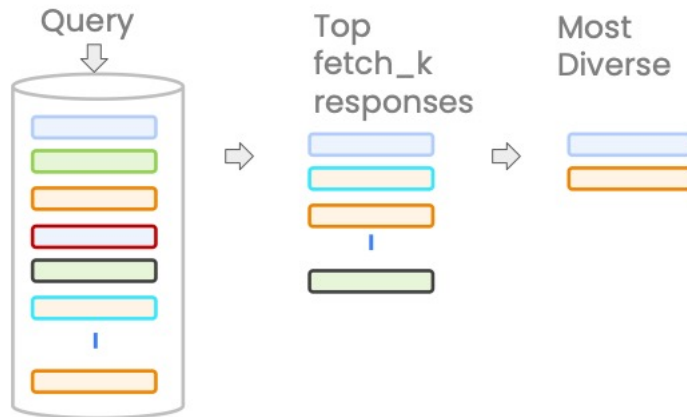


Tell me about all-white mushrooms with large fruiting bodies



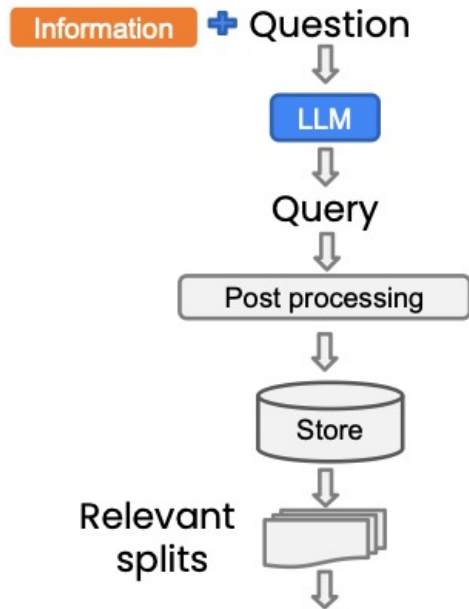
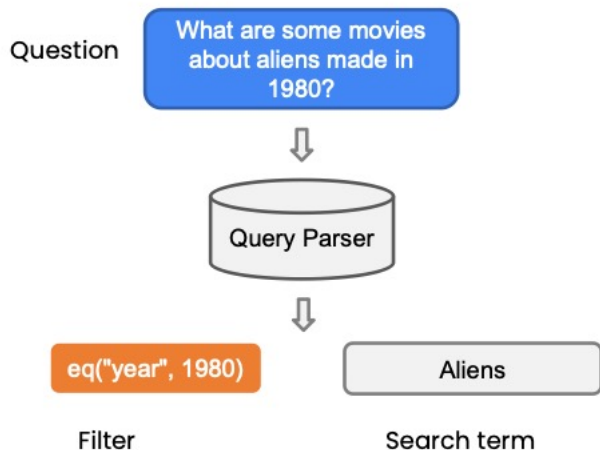
## MMR algorithm

- Query the Vector Store
- Choose the `fetch\_k` most similar responses
- Within those responses choose the `k` most diverse



# LLM-aided Retrieval

- There are several situations where the **Query** applied to the DB is more than just the **Question** asked.
- One is SelfQuery, where we use an LLM to convert the user question into a query



## Self-query

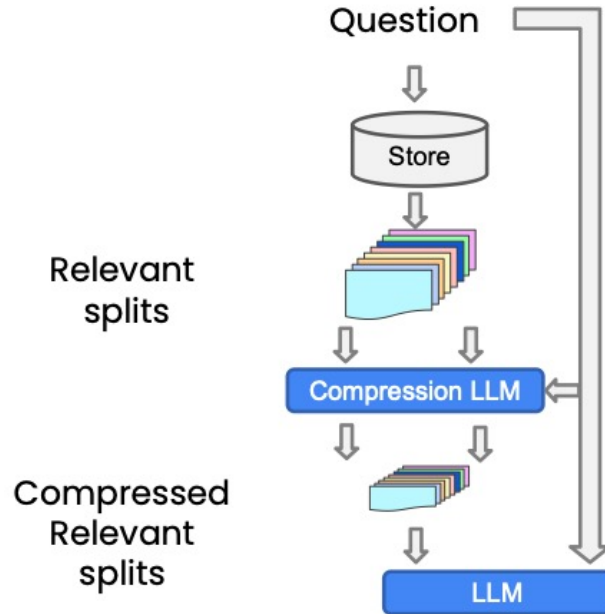
Information: Query format

Query: Question  
Filter: eq["section", "testing"]

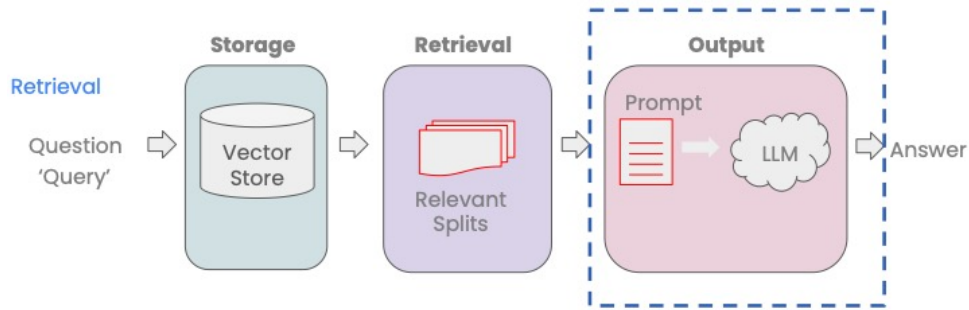
Query parser

# Compression Support

- Increase the number of results you can put in the context by shrinking the responses to only the relevant information.



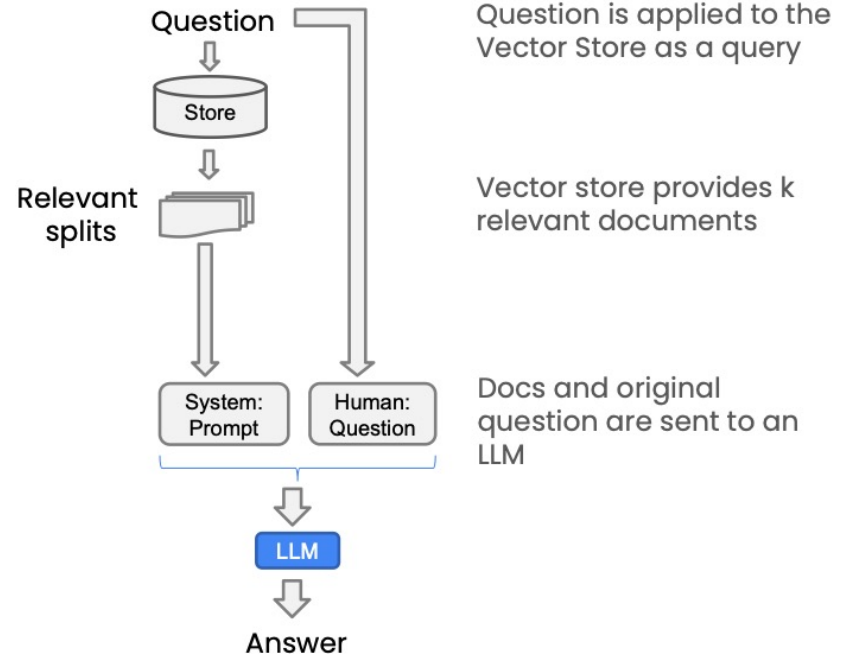
# Q&A Support via Retrieval



- Multiple relevant documents have been retrieved from the vector store
- Potentially compress the relevant splits to fit into the LLM context
- Send the information along with our question to an LLM to select and format an answer

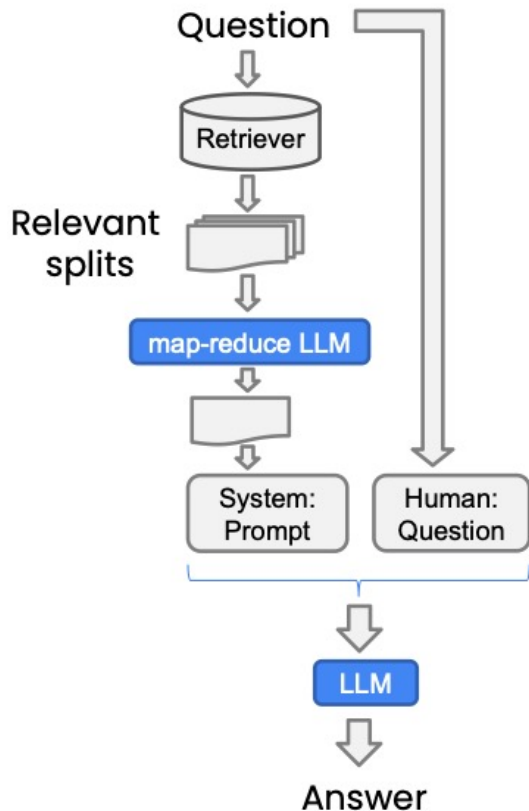
## RetrievalQA chain

`RetrievalQA.from_chain_type(, chain_type="stuff", ...)`



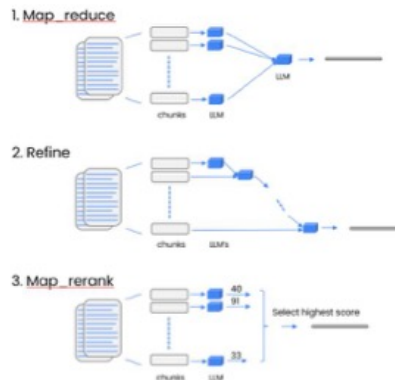
# Retrieval Chain with LLM Selection

```
RetrievalQA.from_chain_type(, chain_type="map_reduce",...)
```



You may have too many docs to fit into an LLM context. The solution is to use an LLM to select the 'most relevant' information

## 3 additional methods



# Facilitating Q&A over Documents with LangChain

## Stuff method



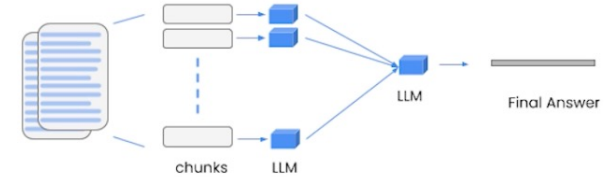
Stuffing is the simplest method. You simply stuff all data into the prompt as context to pass to the language model.

**Pros:** It makes a single call to the LLM. The LLM has access to all the data at once.

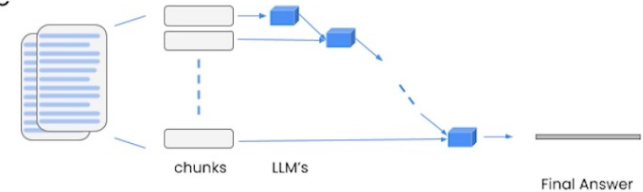
**Cons:** LLMs have a context length, and for large documents or many documents this will not work as it will result in a prompt larger than the context length.

## 3 additional methods

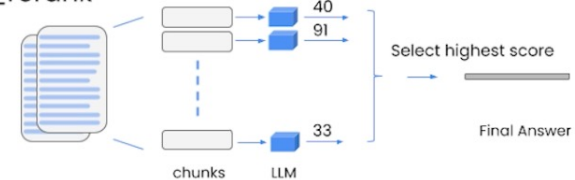
### 1. Map\_reduce



### 2. Refine



### 3. Map\_rerank



# Agent Support in LangChain



Agent refers to the idea of using large language models as reasoning engine to determine which actions to take and in what order.

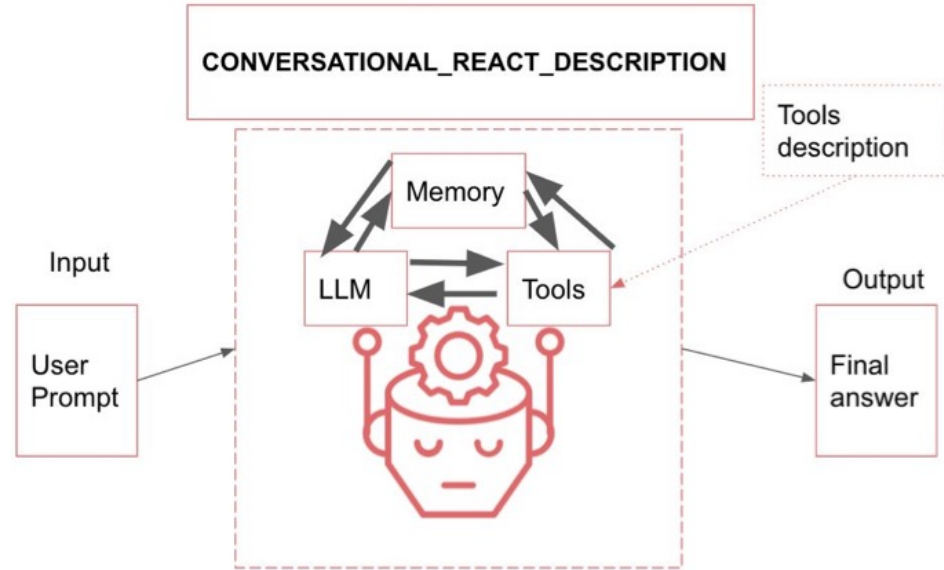
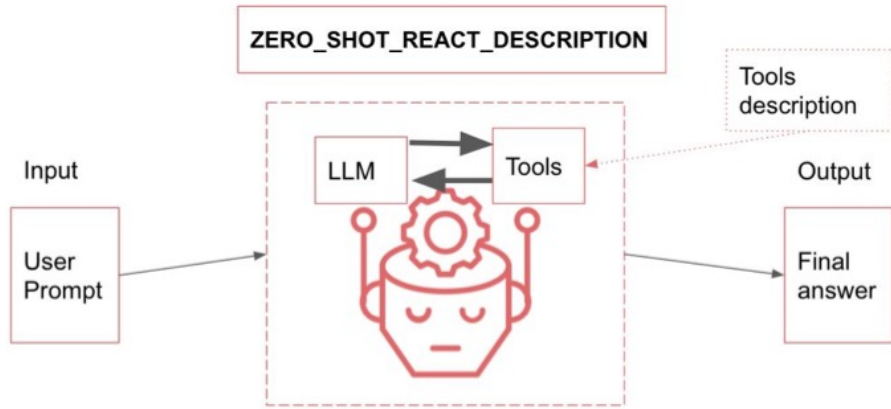
An action can be using a tool and observing its output and deciding what to return to the user.

```
from langchain.agents import  
initialize_agent, AgentType  
  
agent = initialize_agent(tools, llm,  
                        agent=AgentType.ZERO_SHOT_REACT_DESCRIP  
TION, verbose=True)
```

To construct a minimalist Agent, we need:

- **PromptTemplate:** this is responsible for taking the user input and previous steps and constructing a prompt to send to the language model
- **Language Model:** this takes the prompt constructed by the PromptTemplate and returns some output
- **Output Parser:** this takes the output of the Language Model and parses it into an AgentAction or AgentFinish object.

# More Agent Support in LangChain





# Implementation Examples with LangChain

*Source:* LangChain Chat with your Data – a short course by Harrison Chase, <https://www.deeplearning.ai/short-courses/langchain-chat-with-your-data/>

# Interfacing with existing LLMs via "Models"

```
In [ ]: # !pip install openai
```

```
In [1]: import openai
import os

from dotenv import load_dotenv, find_dotenv
_ = load_dotenv(find_dotenv()) # read local .env file
openai.api_key = os.environ['OPENAI_API_KEY']
```

```
In [2]: def get_completion(prompt, model="gpt-3.5-turbo"):
    messages = [{"role": "user", "content": prompt}]
    response = openai.ChatCompletion.create(
        model=model,
        messages=messages,
        temperature=0,
    )
    return response.choices[0].message["content"]
```

```
In [3]: get_completion("What is 1+1?")
```

```
Out[3]: 'As an AI language model, I can tell you that the answer
to 1+1 is 2.'
```

# Models, Prompts and Parsers

```
In [4]: customer_email = """
Arrr, I be fuming that me blender lid \
flew off and splattered me kitchen walls \
with smoothie! And to make matters worse, \
the warranty don't cover the cost of \
cleaning up me kitchen. I need yer help \
right now, matey!
"""
```

```
In [5]: style = """American English \
in a calm and respectful tone
"""
```

```
In [ ]: prompt = f"""Translate the text \
that is delimited by triple backticks
into a style that is {style}.
text: ```{customer_email}```
"""
print(prompt)
```

# Setting up/ Using Prompt Templates in LangChain

```
In [6]: prompt = f"""Translate the text \
that is delimited by triple backticks
into a style that is {style}.
text: ```{customer_email}```
"""
|
print(prompt)
```

```
Translate the text that is delimited by triple backticks
into a style that is American English in a calm and respect
ful tone
```

```
.
text: ```
Arrr, I be fuming that me blender lid flew off and splatte
red me kitchen walls with smoothie! And to make matters wo
rse, \
the warranty don't cover th
n. I need yer help right no
```
```

```
In [7]: response = get_completion(prompt)
```

```
In [8]: response
```

```
Out[8]: "I'm really frustrated that my blender lid flew off and ma
de a mess of my kitchen walls with smoothie! To add insult
to injury, the warranty doesn't cover the cost of cleaning
up my kitchen. I could really use your help right now, my
friend."
```

# Setting up/ Using Prompt Templates in LangChain

```
In [9]: from langchain.chat_models import ChatOpenAI
```

```
In [11]: chat = ChatOpenAI(temperature=0.0)
chat
```

```
Out[11]: ChatOpenAI(verbose=False, callbacks=None, callback_manager
=None, client=<class 'openai.api_resources.chat_completion
n.ChatCompletion'>, model_name='gpt-3.5-turbo', temperatur
e=0.0, model_kwargs={}, openai_api
se=None, openai_organization=None,
ax_retries=6, streaming=False, n=1
```

```
In [12]: template_string = """Translate the text \
that is delimited by triple backticks \
into a style that is {style}. \
text: ```{text}```
"""
```

```
In [13]: from langchain.prompts import ChatPromptTemplate

prompt_template = ChatPromptTemplate.from_template(template
```

```
In [14]: prompt_template.messages[0].prompt
```

```
Out[14]: PromptTemplate(input_variables=['style', 'text'], output_p
arser=None, partial_variables={}, template='Translate the
text that is delimited by triple backticks into a style th
at is {style}. text: ```{text}```\n', template_format='f-s
tring', validate_template=True)
```

# Setting up/ Using Prompt Templates in LangChain

```
In [16]: prompt_template.messages[0].prompt.input_variables
```

```
Out[16]: ['style', 'text']
```

```
In [17]: customer_style = """American English \
in a calm and respectful tone
"""
```

```
In [17]: customer_style = """American English \
in a calm and respectful tone
"""
```

```
In [ ]: customer_email = """
Arrr, I be fuming that me blender
lid flew off and splattered me kitchen
walls with smoothie! And to make matters
worse, the warranty don't cover the cost of
cleaning up me kitchen. I need yer help
right now, matey!
"""
```

```
In [18]: customer_email = """
Arrr, I be fuming that me blender lid \
flew off and splattered me kitchen walls \
with smoothie! And to make matters worse, \
the warranty don't cover the cost of \
cleaning up me kitchen. I need yer help \
right now, matey!
"""
```

```
In [19]: customer_messages = prompt_template.format_messages(
    style=customer_style,
    text=customer_email)
```

```
In [20]: print(type(customer_messages))
print(type(customer_messages[0]))
```

```
<class 'list'>
<class 'langchain.schema.HumanMessage'>
```

# Setting up/ Using Prompt Templates in LangChain

```
In [21]: print(customer_messages[0])
```

```
content="Translate the text that is delimited by triple ba  
ckticks into a style that is American English in a calm an  
d respectful tone\n. text: `blender lid flew off and spl  
smoothie! And to make matter ver the cost of cleaning up  
ght now, matey!\n```\n" addi
```

```
In [22]: customer_response = chat(cus
```

```
In [23]: print(customer_response.cont
```

```
I'm really frustrated that m  
e a mess of my kitchen walls  
rustration, the warranty doe  
g up my kitchen. Can you ple
```

```
In [23]: print(customer_response.content)
```

```
I'm really frustrated that my blender lid flew off and mad  
e a mess of my kitchen walls with smoothie. To add to my f  
rustration, the warranty doesn't cover the cost of cleanin  
g up my kitchen. Can you please help me out, friend?
```

```
In [24]: service_reply = """Hey there customer, \  
the warranty does not cover \  
cleaning expenses for your kitchen \  
because it's your fault that \  
you misused your blender \  
by forgetting to put the lid on before \  
starting the blender. \  
Tough luck! See ya!  
"""
```

```
In [25]: service_style_pirate = """\  
a polite tone \  
that speaks in English Pirate\  
"""
```



# Prompting and Parsing with LangChain

```
In [24]: service_reply = """Hey there customer, \
the warranty does not cover \
cleaning expenses for your kitchen \
because it's your fault that \
you misused your blender \
by forgetting to put the lid on before \
starting the blender. \
Tough luck! See ya!
"""
```

```
In [25]: service_style_pirate = """\
a polite tone \
that speaks in English Pirate\
"""
```

```
In [26]: service_messages = prompt_template.format_messages(
    style=service_style_pirate,
    text=service_reply)

print(service_messages[0].content)
```

```
Translate the text that is delimited by
nto a style that is a polite tone that
irate. text: ``Hey there customer, the
cover cleaning expenses for your kitch
fault that you misused your blender by
he lid on before starting the blender.
``
```

```
In [27]: service_response = chat(service_messages)
print(service_response.content)
```

```
Ahoy there, matey! I must kindly inform ye that the warran
ty be not coverin' the expenses o' cleaning yer galley, as
'tis yer own fault fer misusin' yer blender by forgettin'
to put the lid afore startin' it. Aye, tough luck! Fare
well, me hearty!
```



# Why LangChain introduces Prompt Templates ?

```
prompt = """
Your task is to determine if
the student's solution is
correct or not.
```

```
To solve the problem do the following:
- First, work out your own solution to the problem.
- Then compare your solution to the student's solution
and evaluate if the student's solution is correct or not.
...
Use the following format:
Question:
...
question here
...
Student's solution:
...
student's solution here
...
Actual solution:
...
...
steps to work out the solution and your solution here
...
Is the student's solution the same as actual solution \
just calculated:
...
yes or no
...
Student grade:
...
correct or incorrect
...

Question:
...
{question}
...
Student's solution:
...
{student_solution}
...
Actual solution:
...
***
```

Prompts can be long and detailed.

Reuse good prompts when you can!

LangChain also provides prompts for common operations.

# LangChain o/p Parsing works w/ Prompt Templates

```
EXAMPLES = ["""\nQuestion: What is the elevation range\nfor the area that the eastern sector\nof the Colorado orogeny extends into?\n\nThought: need to search Colorado orogeny, find\nthe area that the eastern sector of the Colorado\norogeny extends into, then find the elevation range\nof the area.\n\nAction: Search[Colorado orogeny]\n\nObservation: The Colorado orogeny was an\nepisode of mountain building (an orogeny) in\nColorado and surrounding areas.\n\nThought: It does not mention the eastern sector.\nSo I need to look up eastern sector.\nAction: Lookup[eastern sector]\n\n...\n\nThought: High Plains rise in elevation from\naround 1,800 to 7,000 ft, so the answer is 1,800 to\n7,000 ft.\n\nAction: Finish[1,800 to 7,000 ft]""",\n]
```

LangChain library  
functions parse the  
LLM's output  
assuming that it will  
use certain keywords.

Example here uses  
**Thought, Action,**  
**Observation** as  
keywords for Chain-  
of-Thought  
Reasoning. (ReAct)

# LangChain Output Parsing: Extraction and Formatting

```
In [28]: {
  "gift": False,
  "delivery_days": 5,
  "price_value": "pretty affordable!"
}

Out[28]: {'gift': False, 'delivery_days': 5, 'price_value': 'pretty
affordable!'}
```

```
In [ ]: customer_review = """\
This leaf blower is pretty amazing. It has four settings:\
candle blower, gentle breeze, windy city, and tornado. \
It arrived in two days, just in time for my wife's \
anniversary present. \
I think my wife liked it so much she was speechless. \
So far I've been the only one using it, and I've been \
using it every other morning to clear the leaves on our lawn. \
It's slightly more expensive than the other leaf blowers \
out there, but I think it's worth it for the extra features.
"""

review_template = """\
For the following text, extract the following information:

gift: Was the item purchased as a gift for someone else? Answer yes or no.
delivery_days: How many days did it take for the product to arrive?
price_value: Extract any sentences about the value or price, including
discounts or promotions.

Format the output as JSON with the following keys:
gift
delivery_days
price_value

text: {text}
"""
```

# LangChain Output Parsing: Extraction and Formatting

```
In [30]: from langchain.prompts import ChatPromptTemplate

prompt_template = ChatPromptTemplate.from_template(review_templ
print(prompt_template)
```

```
input_variables=['text'] output_parser=None partial_variables={} messages=[HumanMessagePromptTemplate(prompt=PromptTemplate(input_variables=['text'], output_parser=None, partial_variables={}, template='For the following text, extract the following information:\n\ngift: Was the item purchased as a gift for someone else? Answer True if yes, False if not or unknown.\ndelivery_days: How many days did it take for the product to arrive? If this information is not found, output -1.\nprice_value: Extract any sentences about the value or price, and output them as a list.\n\nFormat the output as a list of key-value pairs with the following keys:\ngift\ndelivery_days\nprice_value\n', template_format='f-string', val additional_kwargs={})]
```

```
In [31]: messages = prompt_template.format_messages(text=customer_review)
chat = ChatOpenAI(temperature=0.0)
response = chat(messages)
print(response.content)
```

```
{
  "gift": true,
  "delivery_days": 2,
  "price_value": ["It's slightly more expensive than the other leaf blowers out there, but I think it's worth it for the extra features."]
}
```

# LangChain Output Parsing: Extraction and Formatting

```
In [31]: messages = prompt_template.format_messages(text=customer_rev
chat = ChatOpenAI(temperature=0.0)
response = chat(messages)
print(response.content)
```

```
{
  "gift": true,
  "delivery_days": 2,
  "price_value": ["It's slight:
other leaf blowers out there, but
r the extra features."]
}
```

```
In [33]: response.content.get('gift')
```

```
-----
-----
AttributeError                                Traceback (most
recent call last)
Cell In[33], line 1
----> 1 response.content.get('gift')

AttributeError: 'str' object has no attribute 'get'
```

```
In [34]: from langchain.output_parsers import ResponseSchema
from langchain.output_parsers import StructuredOutputParser
```

```
In [ ]: gift_schema = ResponseSchema(name="gift", description="Was t
delivery_days_schema = ResponseSchema(name="delivery_days",
price_value_schema = ResponseSchema(name="price_value", desc
response_schemas = [gift_schema, delivery_days_schema, price_
```

# LangChain Output Parsing: Extraction and Formatting

```
In [34]: from langchain.output_parsers import ResponseSchema
         from langchain.output_parsers import StructuredOutputParser
```

```
In [ ]: gift_schema = ResponseSchema(name="gift", description="Was the item purchased as a gift for someone else? Answer True if yes, known. Otherwise, answer False.", output_key="gift")
         delivery_days_schema = ResponseSchema(name="delivery_days", description="How many days did it take for the product to arrive? If this information is not found, output -1.", output_key="delivery_days")
         price_value_schema = ResponseSchema(name="price_value", description="Extract any sentences about the value or price.", output_key="price_value")
         response_schemas = [gift_schema, delivery_days_schema, price_value_schema]
```

```
In [38]: format_instructions = output_parser.get_format_instructions()
```

```
In [39]: print(format_instructions)
```

The output should be a markdown code snippet following the following schema, including the leading and trailing backslashes: ````json` and `````:

```
```json
```

```
{
  "gift": string // Was the item purchased as a gift for someone else? Answer True if yes, known. Otherwise, answer False.
```

```
  "delivery_days": string // How many days did it take for the product to arrive? If this information is not found, output -1.
```

```
  "price_value": string // Extract any sentences about the value or price.
```

```
In [ ]: review_template_2 = """\
For the following text, extract the following information:

gift: Was the item purchased as a gift for someone else? Answer True if yes, known. Otherwise, answer False.
delivery_days: How many days did it take for the product to arrive? If this information is not found, output -1.
price_value: Extract any sentences about the value or price.

text: {text}

{format_instructions}
"""

prompt = ChatPromptTemplate.from_template(template=review_template_2)
messages = prompt.format_messages(text=customer_review,
                                  format_instructions=format_instructions)
```



# LangChain Output Parsing: Extraction and Formatting

```
In [41]: print(messages[0].content)
```

For the following text, extract the following information:

gift: Was the item purchased as a gift for someone else? A answer True if yes, False if not or unknown.  
delivery\_days: How many days did it take for the product to arrive? If this information is not found, output -1.  
price\_value: Extract any sentences about the value or price, and output them as a comma separated Python list.

text: This leaf blower is pretty amazing. It has four settings: candle blower, gentle breeze, windy city, and tornado. It arrived in two days, just in time for my wife's anniversary present. I think my wife liked it so much she was speechless. So far I've been the only one using it, and I've been using it every other morning to clear the leaves on our lawn. It's slightly more expensive than the other leaf blowers out there, but I think it's worth it for the extra features.

The output should be a markdown code snippet formatted in the following schema, including the leading and trailing ```json and ```:

```
```json
{
  "gift": string // Was the item purchased as a gift for someone else? Answer True if yes, False if not or unknown.
  "delivery_days": string // How many days did it take for the product to arrive? If this information is not found, output -1.
  "price_value": string // Extract any sentences about the value or price, and output them as a comma separated Python list.
}
```

```
In [42]: response = chat(messages)
```

```
In [43]: print(response.content)
```

```
```json
{
  "gift": true,
  "delivery_days": "2",
  "price_value": ["It's slightly more expensive than the other leaf blowers out there, but I think it's worth it for the extra features."]
}
```

```
In [44]: output_dict = output_parser.parse(response.content)
```

```
In [45]: output_dict
```

```
Out[45]: {'gift': True,
          'delivery_days': '2',
          'price_value': ["It's slightly more expensive than the other leaf blowers out there, but I think it's worth it for the extra features."]}
```

```
In [46]: type(output_dict)
```

```
Out[46]: dict
```

```
In [49]: output_dict.get('delivery_days')
```


```
Out[49]: '2'
```

# Empowering LLM-based Applications with External Actions



# How ChatGPT interacts w/ the Outside World via Plugins

**Prompt Engineering Strategies**

➤

**Use External Tools**  
Enhances model capabilities.  
Tactics include embeddings-based search, code execution, and access to specific functions.

Actions  
Memory

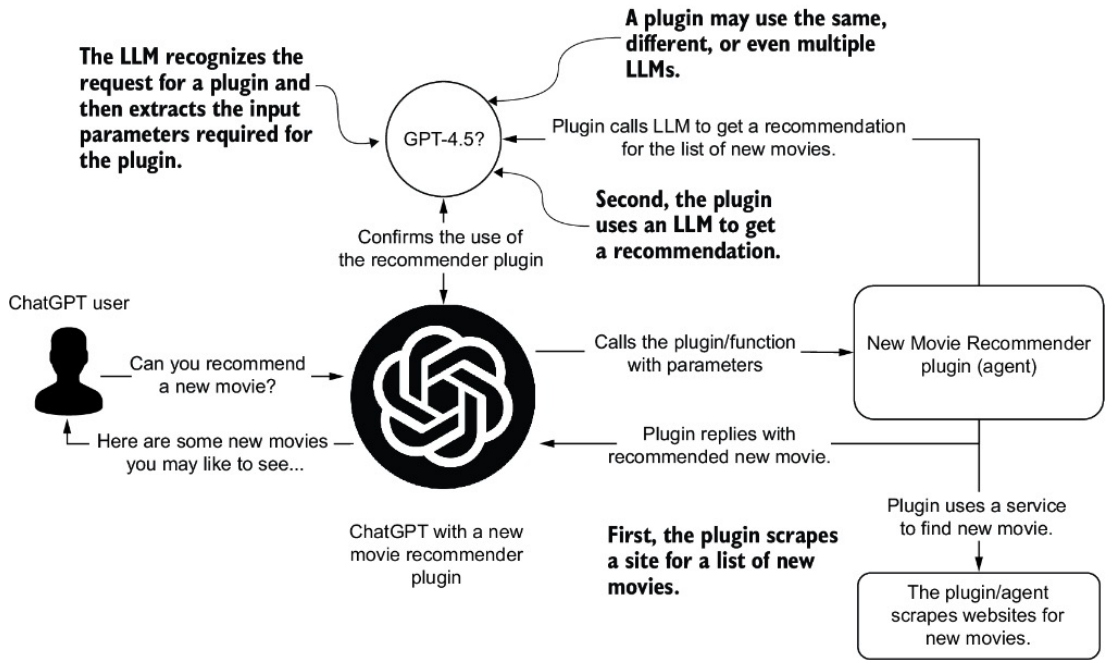
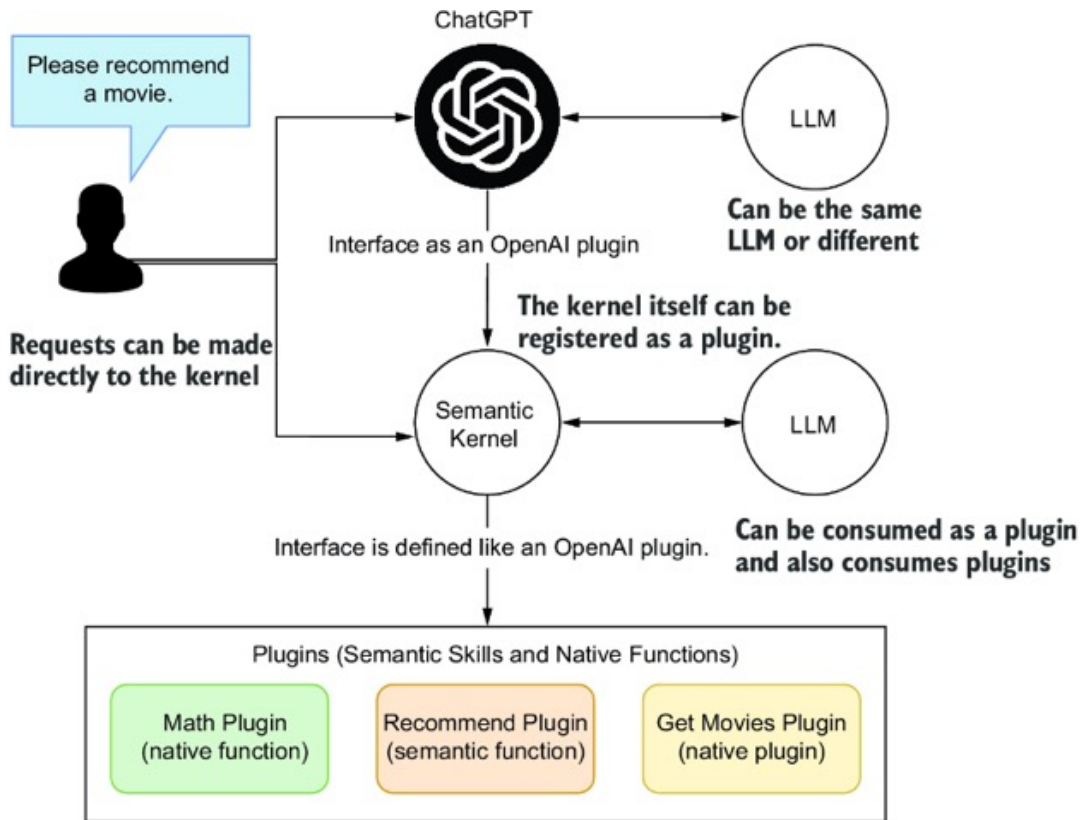


Figure 5.1 How a ChatGPT plugin operates and how plugins and other external tools (e.g., APIs) align with the Use External Tools prompt engineering strategy

# Microsoft's Semantic Kernel as a Plugin for ChatGPT



Semantic Kernel (SK) is another open source project from Microsoft intended to help build AI applications. At its core, the project is best used to define actions, or what the platform calls semantic plugins, which are wrappers for skills and functions.

This figure shows how the SK can be used as a plugin and a consumer of OpenAI plugins.

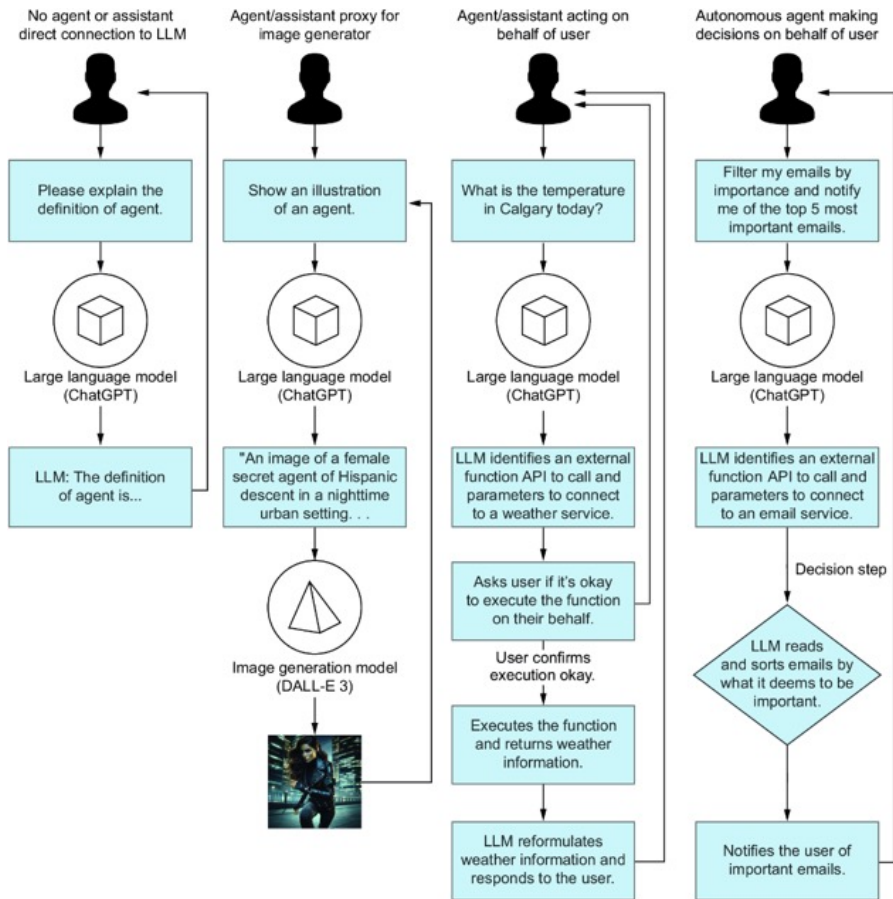
The SK relies on the OpenAI plugin definition to define a plugin. That way, it can consume and publish itself or other plugins to other systems.

An OpenAI plugin definition maps precisely to the function definitions. This means that SK is the orchestrator of API tool calls, aka plugins.

That also means that SK can help organize multiple plugins with a chat interface or an agent.

# From LLM-powered Application to AI Agent

# From “LLM-powered Application” to “AI Agent”



Agent == Act on “YOUR” behalf at various level of autonomy !

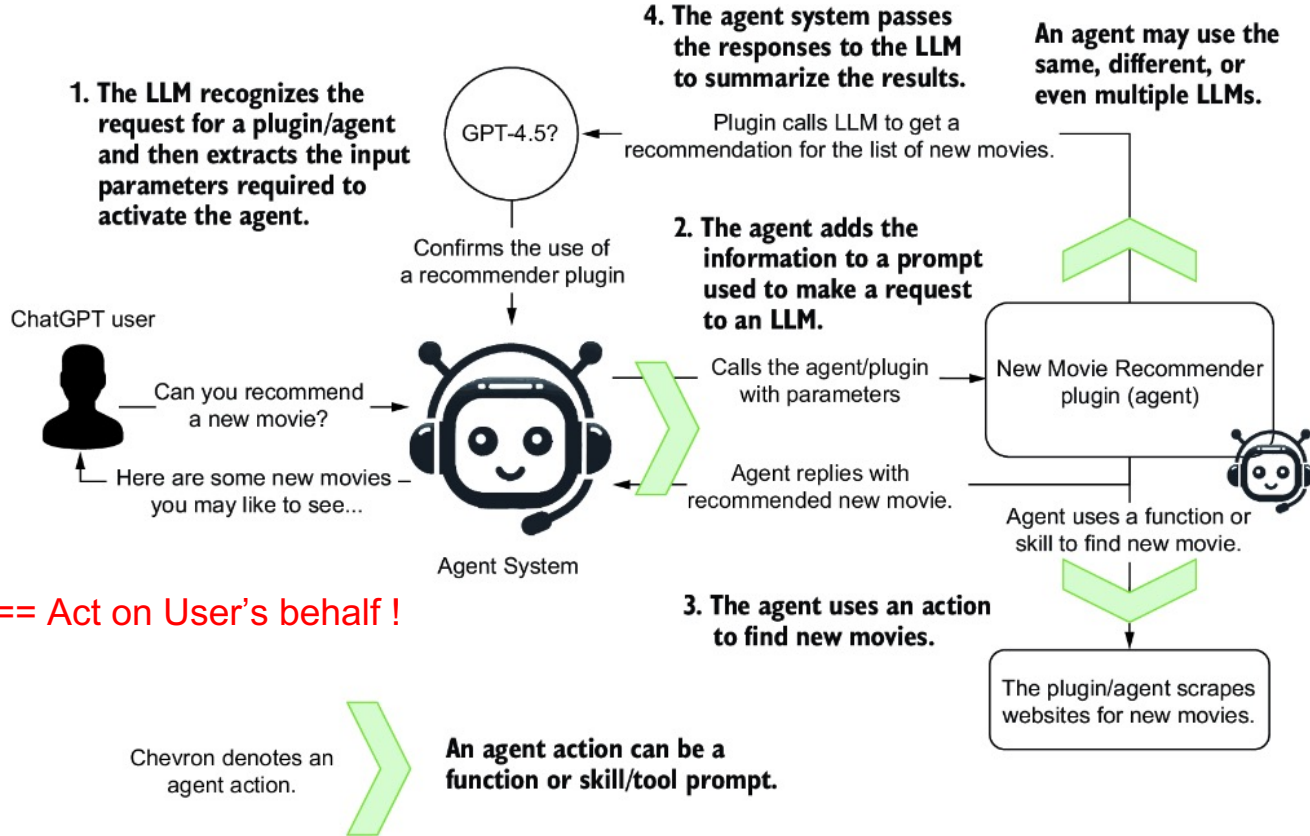
“YOUR” can be == Users

“YOUR” can also be an Ext. App

Source: AI Agents in Action, by Michael Lanham, Feb 2025, Manning Publications

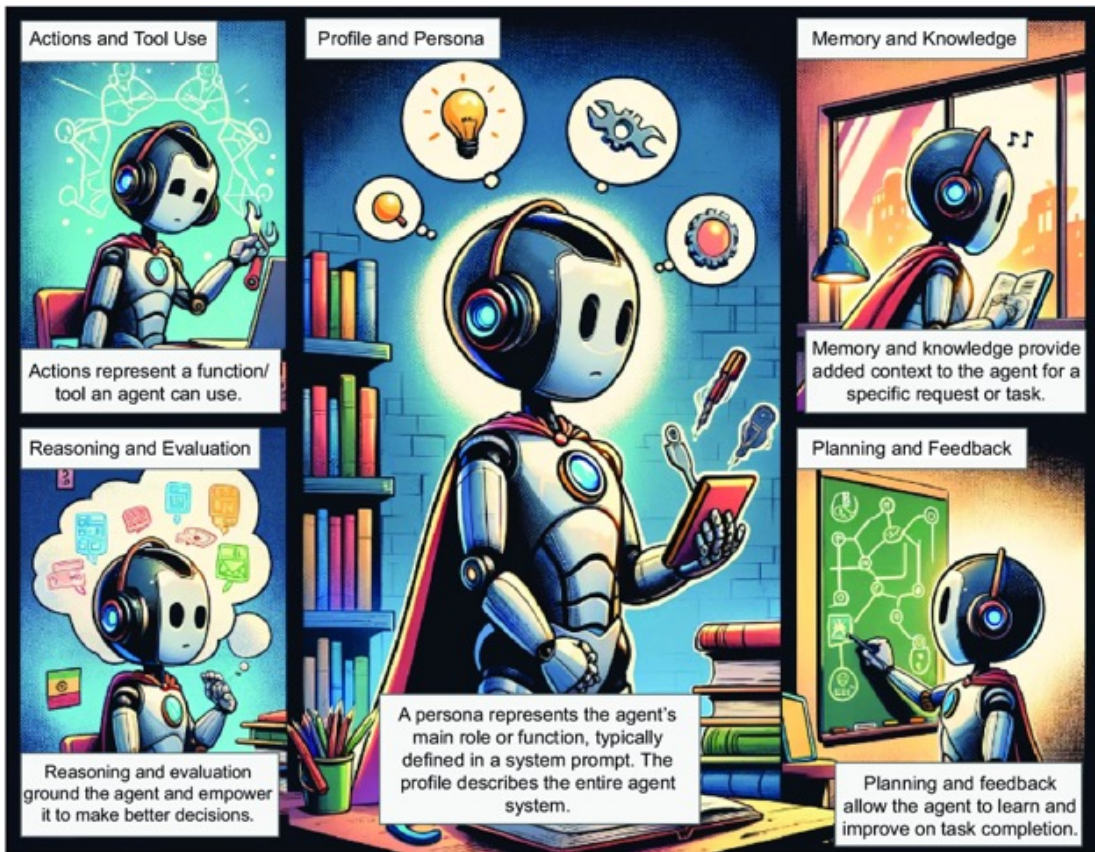
The differences between the LLM interactions from direct action compared to using proxy agents, agents, and autonomous agents

# How an Agent uses Actions to perform External Tasks



Agent == Act on User's behalf !

# Key Functional Components of a Single-Agent System



# In-Depth Look at the Key Functional Components of an Agent



# The Profile and Persona Component of an Agent



## Profile and Persona



### Profile Contents

Persona: Role, i.e., coder or tester  
Demographics: Sex, age, background

### Profile Generation

Handcrafted: Manually designed by humans  
LLM generated: Directed by human prompts  
Data generated: Constructed from data personas

**Agent persona:** We'll understand how to clearly define the persona, specifying their role and characteristics to guide the agent effectively.

**Agent role and demographics:** We'll see how relevant demographic and role details can provide agent context, such as age, gender, or background, for a more relevant interaction.

**Human vs. AI assistance for persona generation:** We'll highlight the role of human involvement in persona generation, whether it's entirely human driven or assisted by LLMs or other agents.

**Innovative persona techniques:** Prompts generated through data or other novel approaches such as evolutionary algorithms to enhance agent capabilities.



# Agent Action and Tool Use

**Action targets:** We'll learn the importance of defining action targets, whether for task completion, exploration, or communication, to clarify the agent's objectives.

**Action space and impact:** We'll learn the significance of understanding how actions affect task completion and their effect on the agent's environment, internal states, and self-knowledge.

**Action generation methods:** We'll see the various ways actions can be generated, such as manual creation, memory recollection, or plan following, to illustrate the diversity of agent behaviors.



# Agent Memory and Knowledge



## Memory and Knowledge



### Retrieval Structure

- Unified
- Hybrid

### Retrieval Formats

- Language
- Databases
- Embeddings
- Lists

### Retrieval Operation

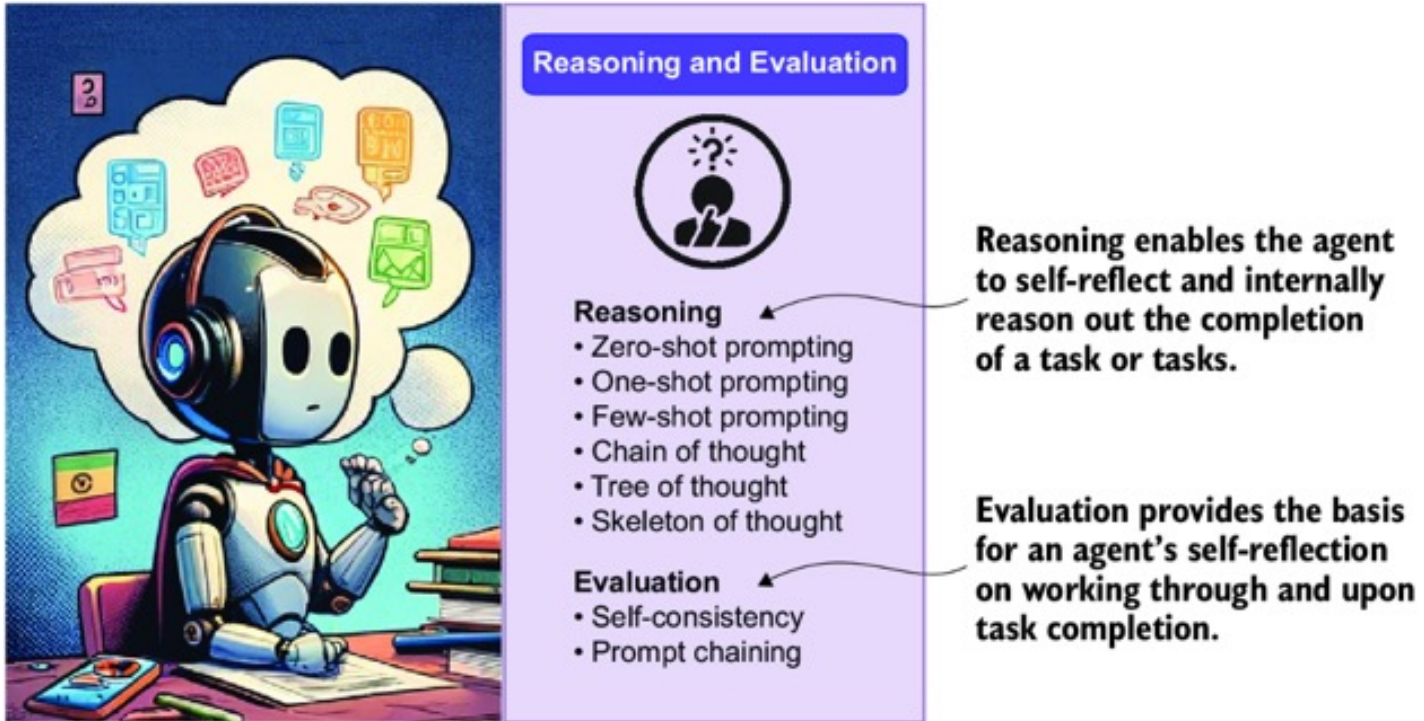
- Augmentation
- Semantic Extraction
- Compression

**Retrieval structure variety:** We'll learn about the diverse memory structures agents can employ, including unified and hybrid approaches, enabling flexibility in information storage.

**Retrieval formats:** We'll explore the various data sources for memory, such as language (e.g., PDF documents), databases (relational, object, or document), and embeddings, offering a rich pool of information to draw upon.

**Semantic similarity:** We'll learn how embeddings enable semantic similarity searches, facilitating efficient retrieval of relevant data and enhancing the agent's decision-making capabilities.

# Reasoning and Evaluation Functions of an Agent





# Planning and Feedback Functions of an Agent

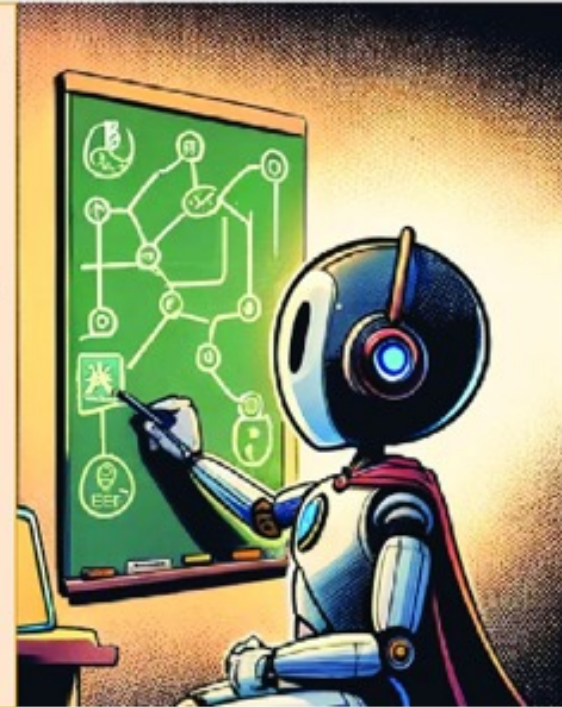
We'll look at various planning strategies with and without feedback—from basic and sequential planners to automatic tool use with reasoning.

Feedback may come from a variety of sources, such as environmental, human, and an LLM via various constructive feedback patterns.

## Planning and Feedback



- ▶ **Planning without feedback (autonomous)**
  - Basic planning
  - Automatic reasoning with tool use
  - Sequential planning
- ▶ **Planning with feedback**
  - Environmental feedback
  - Human feedback
  - LLM feedback
  - Adaptive constructive feedback



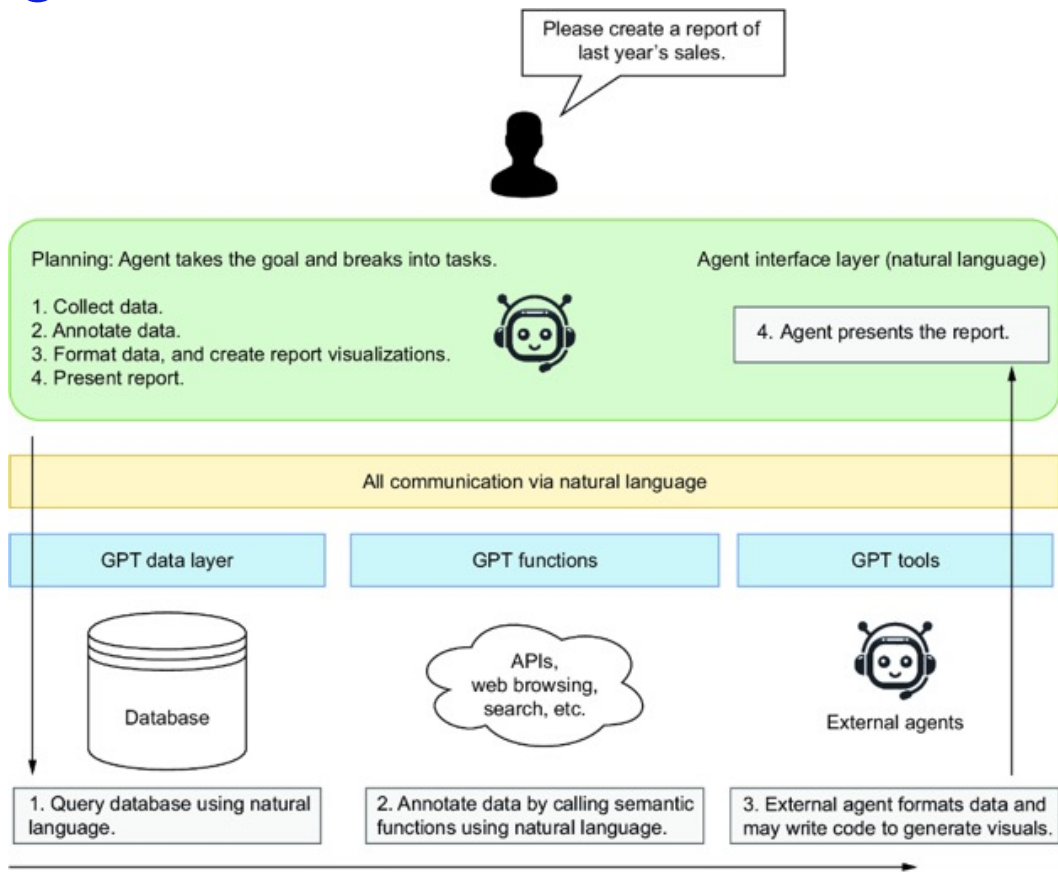
# Comparisons of Agentic AI Frameworks

Framework	Best For	Deployment	Scalability	Current Adoption
LangChain & LangGraph	Enterprise AI apps, multi-agent workflows	Self-hosted, Cloud	High – widely used in production	Most adopted
Autogen & Semantic Kernel	Microsoft ecosystem, scalable AI apps	Azure, Self-hosted	High – enterprise-scale	Growing in enterprises
LlamaIndex	AI-powered search, RAG	Cloud, Self-hosted	High – efficient data processing	Strong in AI search
AutoGPT	No-code AI automation, continuous agents	Cloud-first, some self-hosted options	Medium – automation focused	Popular for prototyping
CrewAI	Workflow-based multi-agent AI apps	Cloud, Self-hosted	Medium – still early-stage	Fast-growing
PydanticAI	FastAPI & Pydantic-based AI apps	Self-hosted	Low – lightweight framework	Niche adoption
Spring AI	Java-based AI applications	Self-hosted, Cloud	Medium – depends on Java ecosystem	Popular in Java world
Haystack	LLM-powered search & RAG	Self-hosted, Cloud	High – built for production	Strong in search applications

Source: <https://medium.com/@varmalearn/a-casual-guide-to-the-top-agentic-ai-frameworks-in-2025-9299dd94e034>



# A Vision on how an Agent Interacts with other Software Systems



# Implementing AI Agents with LangChain / LangGraph

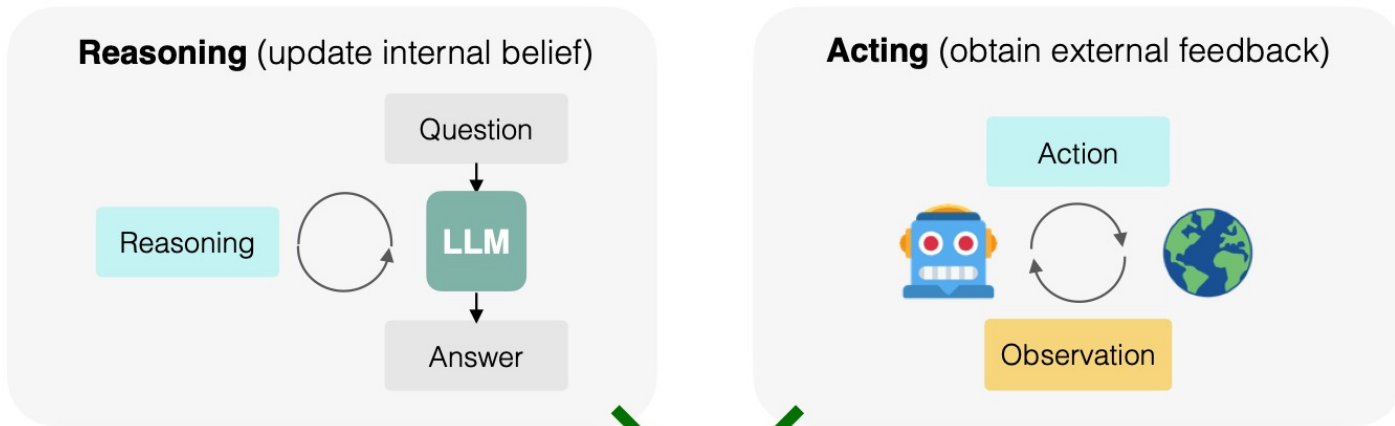


# Build an Agent from Scratch

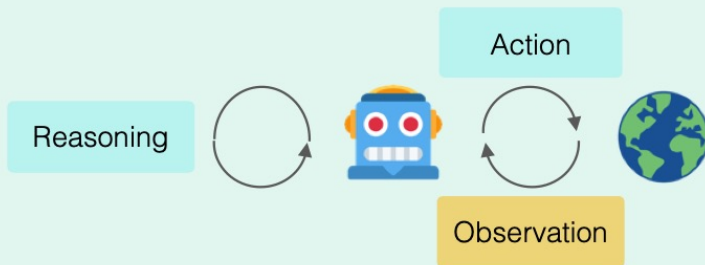


Use LangChain only, Not LangGraph !

# Recap: ReAct = Reasoning and Acting



**ReAct:** a new paradigm of agents that **reason and act**



- **Synergy** of reasoning and acting
- **Simple** and intuitive to use
- **General** across domains

# Example of ReAct w/ a Zero-Shot ReAct Prompt

You are an agent that answers questions by using two actions:

- search[query]: Google search the query. You can also use it to calculate math.
- finish[answer]: return the answer

Your generation should have the following format:

**Thought:** to analyze the observation and inform the action.

**Action:** your action

**Question:** If I have seven trillion dollars today, can I buy Apple, Nvidia, and Microsoft? If not, how much more money do I need?

# Example of ReAct

LLM

(GPT-4)

**Thought:** I need to find the current market capitalization of Apple, Nvidia, and Microsoft to understand how much it would cost to buy all of them.

**Action:** search["current market capitalization of Apple, Nvidia, and Microsoft"]

**Observation:** "Nvidia's market cap stood at \$1.832 trillion as of about 10 a.m. ET Wednesday, while Alphabet's was at \$1.814 trillion, bringing Nvidia behind just Microsoft (MSFT) at \$3.015 trillion and Apple (AAPL) at \$2.836 trillion."

**Thought:** Now that I have the market cap for all companies, let's add them together and see if \$7 trillion is enough.

**Action:** search["1.832 trillion + 3.015 trillion + 2.836 trillion"]

**Observation:** 7683000000000

**Thought:** The amount of money needed in addition to the \$7 trillion is \$683 billion.

**Action:** finish["You need an additional \$683 billion to be able to buy Apple, Nvidia, and Microsoft."] 



(Mar 2024)

# Example of ReAct

LLM

(GPT-4)

**Thought:** I need to find the current market capitalization of Apple, Nvidia, and Microsoft to understand how much it would cost to buy all of them.

**Action:** search["current market capitalization of Apple, Nvidia, and Microsoft"]

**Observation:** nothing found!

(My adversarial observation)

**Thought:** The search did not give any results. I will need to search for each company's market capitalization individually.

**Action:** search["current market capitalization of Apple"]

**Observation:** 'price': 180.75

(My adversarial observation)

**Thought:** This is the current price of Apple's stock, not its market capitalization. The market cap is calculated by multiplying the share price by the number of outstanding shares. Thus, I need to find the total number of Apple's outstanding shares to calculate its market cap.

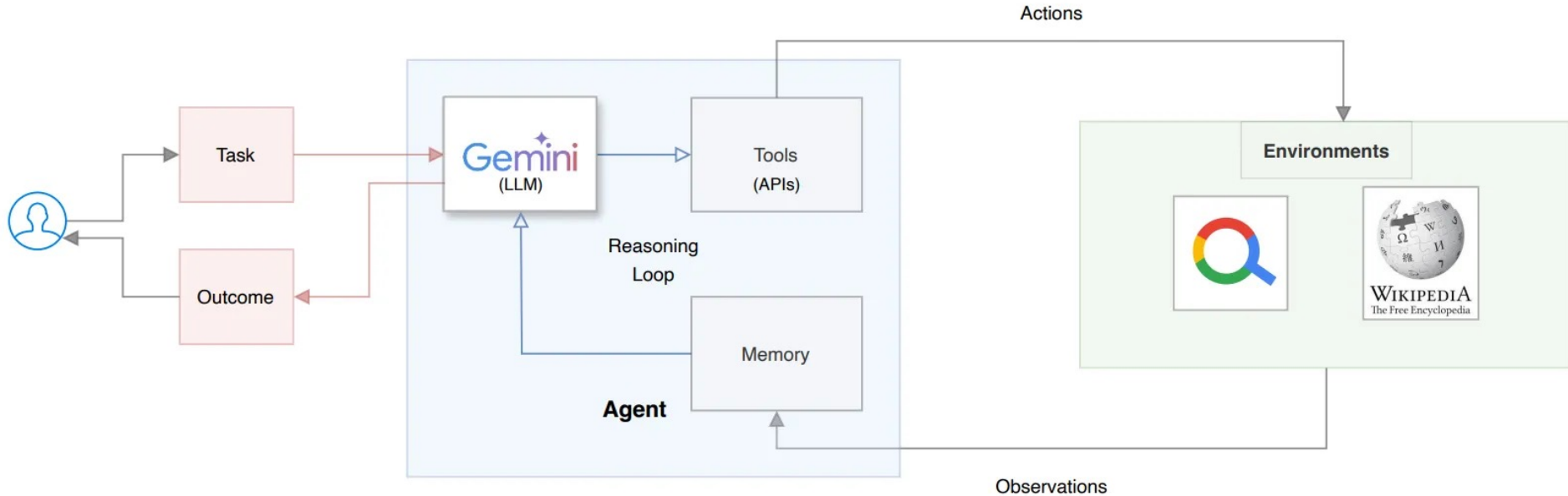
**Action:** search["total number of Apple's outstanding shares"]



(Continues to solve the task)

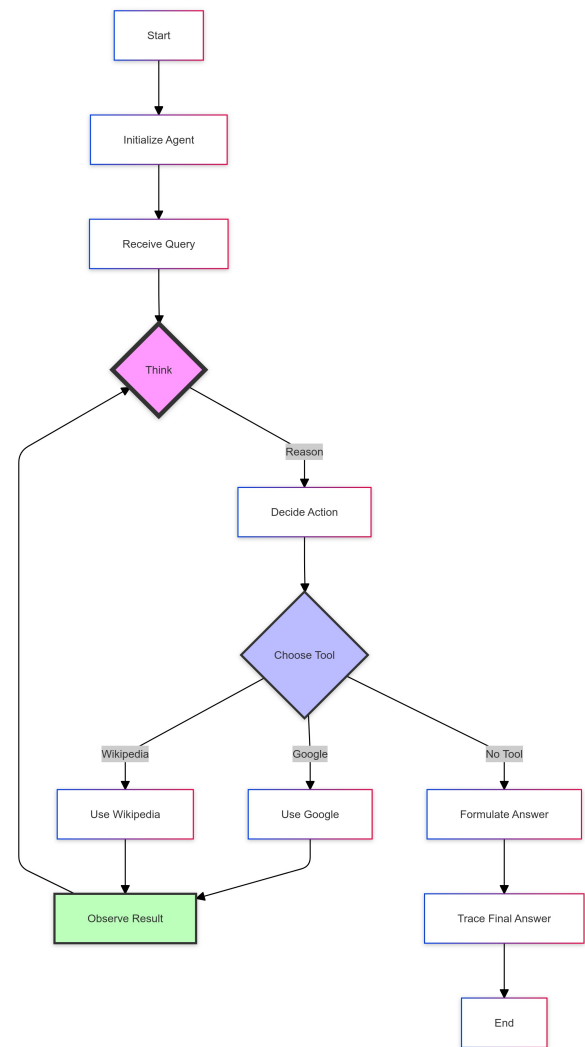
**Synergy:** acting support reasoning, reasoning guides acting

# High-level Architecture of a ReAct Agent using Gemini



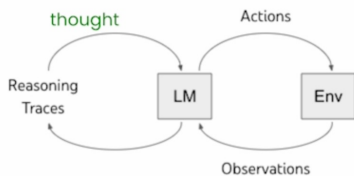
Source: <https://medium.com/google-cloud/building-react-agents-from-scratch-a-hands-on-guide-using-gemini-ffe4621d90ae>  
<https://github.com/arunpshankar/react-from-scratch>

# Process-Flow of a ReAct Agent





# Building a ReAct AI agent using LangChain Prompt Templates/ Library



ReAct (Reason + Act)

Published as a conference paper at ICLR 2023

## REACT: SYNERGIZING REASONING AND ACTING IN LANGUAGE MODELS

Shunyu Yao<sup>\*1</sup>, Jeffrey Zhao<sup>2</sup>, Dian Yu<sup>2</sup>, Nan Du<sup>2</sup>, Izhak Shafran<sup>2</sup>, Karthik Narasimhan<sup>1</sup>, Yuan Cao<sup>2</sup>

<sup>1</sup>Department of Computer Science, Princeton University

<sup>2</sup>Google Research, Brain team

<sup>\*</sup>{shunyuy, karthikh}@princeton.edu

```
prompt = """
You run in a loop of Thought, Action, PAUSE, Observation.
At the end of the loop you output an Answer
Use Thought to describe your thoughts about the question you have been a
Use Action to run one of the actions available to you - then return PAUS
Observation will be the result of running those actions.

Your available actions are:

calculate:
e.g. calculate: 4 * 7 / 3
Runs a calculation and returns the number - uses Python so be sure to us

average_dog_weight:
e.g. average_dog_weight: Collie
returns average weight of a dog when given the breed

Example session:

Question: How much does a Bulldog weigh?
Thought: I should look the dogs weight using average_dog_weight
Action: average_dog_weight: Bulldog
PAUSE

You will be called again with this:

Observation: A Bulldog weights 51 lbs

You then output:

Answer: A bulldog weights 51 lbs
""".strip()
```

# Building a ReAct AI agent using LangChain Prompt Templates/ Library

```
# based on https://til.simonwillison.net/llms/python-react-pattern
```

```
import openai
import re
import httpx
import os
from dotenv import load_dotenv, find_dotenv

_ = load_dotenv(find_dotenv())
from openai import OpenAI
```

```
client = OpenAI()
```

```
chat_completion = client.chat.completions.create(
    model="gpt-3.5-turbo",
    messages=[{"role": "user", "content": "Hello world"}]
)
```

```
chat_completion.choices[0].message.content
```

```
'Hello! How can I assist you today?'
```

```
class Agent:
    def __init__(self, system=""):
        self.system = system
        self.messages = []
        if self.system:
            self.messages.append({"role": "system", "content": system})

    def __call__(self, message):
        self.messages.append({"role": "user", "content": message})
        result = self.execute()
        self.messages.append({"role": "assistant", "content": result})
        return result
```

```
class Agent:
    def __init__(self, system=""):
        self.system = system
        self.messages = []
        if self.system:
            self.messages.append({"role": "system", "content": system})

    def __call__(self, message):
        self.messages.append({"role": "user", "content": message})
        result = self.execute()
        self.messages.append({"role": "assistant", "content": result})
        return result

    def execute(self):
        completion = client.chat.completions.create(
            model="gpt-4-0125-preview",
            temperature=0,
            messages=self.messages)
        return completion.choices[0].message.content
```

# Building a ReAct AI agent using LangChain Prompt Templates/ Library

```
prompt = """
You run in a loop of Thought, Action, PAUSE, Observation.
At the end of the loop you output an Answer
Use Thought to describe your thoughts about the question you have been a
Use Action to run one of the actions available to you - then return PAUSE
Observation will be the result of running those actions.

Your available actions are:

calculate:
e.g. calculate: 4 * 7 / 3
Runs a calculation and returns the number - uses Python so be sure to us

average_dog_weight:
e.g. average_dog_weight: Collie
returns average weight of a dog when given the breed

Example session:

Question: How much does a Bulldog weigh?
Thought: I should look the dogs weight using average_dog_weight
Action: average_dog_weight: Bulldog
PAUSE

You will be called again with this:

Observation: A Bulldog weights 51 lbs

You then output:

Answer: A bulldog weights 51 lbs
""".strip()
```

# Building a ReAct AI agent using LangChain Prompt Templates/ Library

```
def calculate(what):  
    return eval(what)  
  
def average_dog_weight(name):  
    if name in "Scottish Terrier":  
        return("Scottish Terriers average 20 lbs")  
    elif name in "Border Collie":  
        return("a Border Collies average weight is 37 lbs")  
    elif name in "Toy Poodle":  
        return("a toy poodles average weight is 7 lbs")  
    else:  
        return("An average dog weights 50 lbs")  
  
known_actions = {  
    "calculate": calculate,  
    "average_dog_weight": average_dog_weight  
}
```

```
abot = Agent(prompt)
```

```
result = abot("How much does a toy poodle weigh?")  
print(result)
```

```
Thought: I should look up the dog's weight using average_dog_weight for a Toy Poodle.
```

```
Action: average_dog_weight: Toy Poodle  
PAUSE
```

```
result = average_dog_weight("Toy Poodle")
```

```
result
```

```
'a toy poodles average weight is 7 lbs'
```

```
result
```

```
'a toy poodles average weight is 7 lbs'
```

```
next_prompt = "Observation: {}".format(result)
```

```
abot(next_prompt)
```

```
'Answer: A Toy Poodle weighs 7 lbs.'
```

```
abot.messages
```

```
[{'role': 'system',  
  'content': 'You run in a loop of Thought, Action, PAUSE, Observation. At the end of the loop you output an Answer. Use Thought to describe your thoughts about the question you have been asked. Use Action to run one of the actions available to you - then return PAUSE. Observation will be the result of running those actions. Your available actions are: calculate: e.g. calculate: 4 * 7 / 3. Returns a calculation and returns the number - uses Python so be sure to use floating point syntax if necessary. average_dog_weight: e.g. average_dog_weight: Collie. Returns average weight of a dog when given the breed. Example session: Question: How much does a Bulldog weigh? Thought: I should look the dogs weight using average_dog_weight. Action: average_dog_weight: Bulldog. PAUSE. You will be called again with this: Observation: A Bulldog weights 51 lbs. You then output: Answer: A bulldog weights 51 lbs'},  
 {'role': 'user', 'content': 'How much does a toy poodle weigh?'},  
 {'role': 'assistant',  
  'content': "Thought: I should look up the dog's weight using average_dog_weight for a Toy Poodle. Action: average_dog_weight: Toy Poodle. PAUSE"},  
 {'role': 'user',  
  'content': 'Observation: a toy poodles average weight is 7 lbs'},  
 {'role': 'assistant', 'content': 'Answer: A Toy Poodle weighs 7 lbs.'}]
```

# Building a ReAct AI agent using LangChain Prompt Templates/ Library

```
abot = Agent(prompt)
```

```
question = """I have 2 dogs, a border collie and a scottish terrier. \nWhat is their combined weight"""\nabot(question)
```

```
"Thought: To find the combined weight of a Border Collie and a Scottish Terrier, I need to first find the average weight of each breed and then add those weights together. I'll start by finding the average weight of a Border Collie.\n\nAction: average_dog_weight: Border Collie\nPAUSE"
```

```
next_prompt = "Observation: {}".format(average_dog_weight("Border Collie"))\nprint(next_prompt)
```

```
Observation: a Border Collies average weight is 37 lbs
```

```
abot(next_prompt)
```

```
"Thought: Now that I know a Border Collie's average weight is 37 lbs, I need to find the average weight of a Scottish Terrier to calculate the combined weight.\n\nAction: average_dog_weight: Scottish Terrier\nPAUSE"
```

```
next_prompt = "Observation: {}".format(average_dog_weight("Scottish Terrier"))\nprint(next_prompt)
```

```
Observation: Scottish Terriers average 20 lbs
```

```
abot(next_prompt)
```

```
'Thought: With the average weight of a Border Collie being 37 lbs and a Scottish Terrier being 20 lbs, I can now calculate their combined weight.\n\nAction: calculate: 37 + 20\nPAUSE'
```

```
Observation: 57
```

```
abot(next_prompt)
```

```
'Answer: The combined weight of a Border Collie and a Scottish Terrier is 57 lbs.'
```

---

# Building a ReAct AI agent using LangChain Prompt Templates/ Library

```
action_re = re.compile('^Action: (\w+): (.*)$')
```

```
def query(question, max_turns=5):
    i = 0
    bot = Agent(prompt)
    next_prompt = question
    while i < max_turns:
        i += 1
        result = bot(next_prompt)
        print(result)
        actions = [
            action_re.match(a)
            for a in result.split('\n')
            if action_re.match(a)
        ]
        if actions:
            # There is an action to run
            action, action_input = actions[0].groups()
            if action not in known_actions:
                raise Exception("Unknown action: {}: {}".format(action,
                    print("-- running {} {}".format(action, action_input))
                    observation = known_actions[action](action_input)
                    print("Observation:", observation)
                    next_prompt = "Observation: {}".format(observation)
            else:
                return
```

```
question = """I have 2 dogs, a border collie and a scottish terrier. \
What is their combined weight"""
query(question)]
```

```
question = """I have 2 dogs, a border collie and a scottish terrier. \
What is their combined weight"""
query(question)
```

Thought: To find the combined weight of a Border Collie and a Scottish Terrier, I need to first find the average weight of each breed and then add those weights together. I will start by finding the average weight of a Border Collie.

Action: average\_dog\_weight: Border Collie

PAUSE

-- running average\_dog\_weight Border Collie

Observation: a Border Collies average weight is 37 lbs

Thought: Now that I know the average weight of a Border Collie is 37 lbs, I need to find the average weight of a Scottish Terrier to calculate the combined weight of both dogs.

Action: average\_dog\_weight: Scottish Terrier

PAUSE

-- running average\_dog\_weight Scottish Terrier

Observation: Scottish Terriers average 20 lbs

Thought: With the average weight of a Border Collie being 37 lbs and a Scottish Terrier being 20 lbs, I can now calculate their combined weight by adding these two numbers together.

Action: calculate: 37 + 20

PAUSE

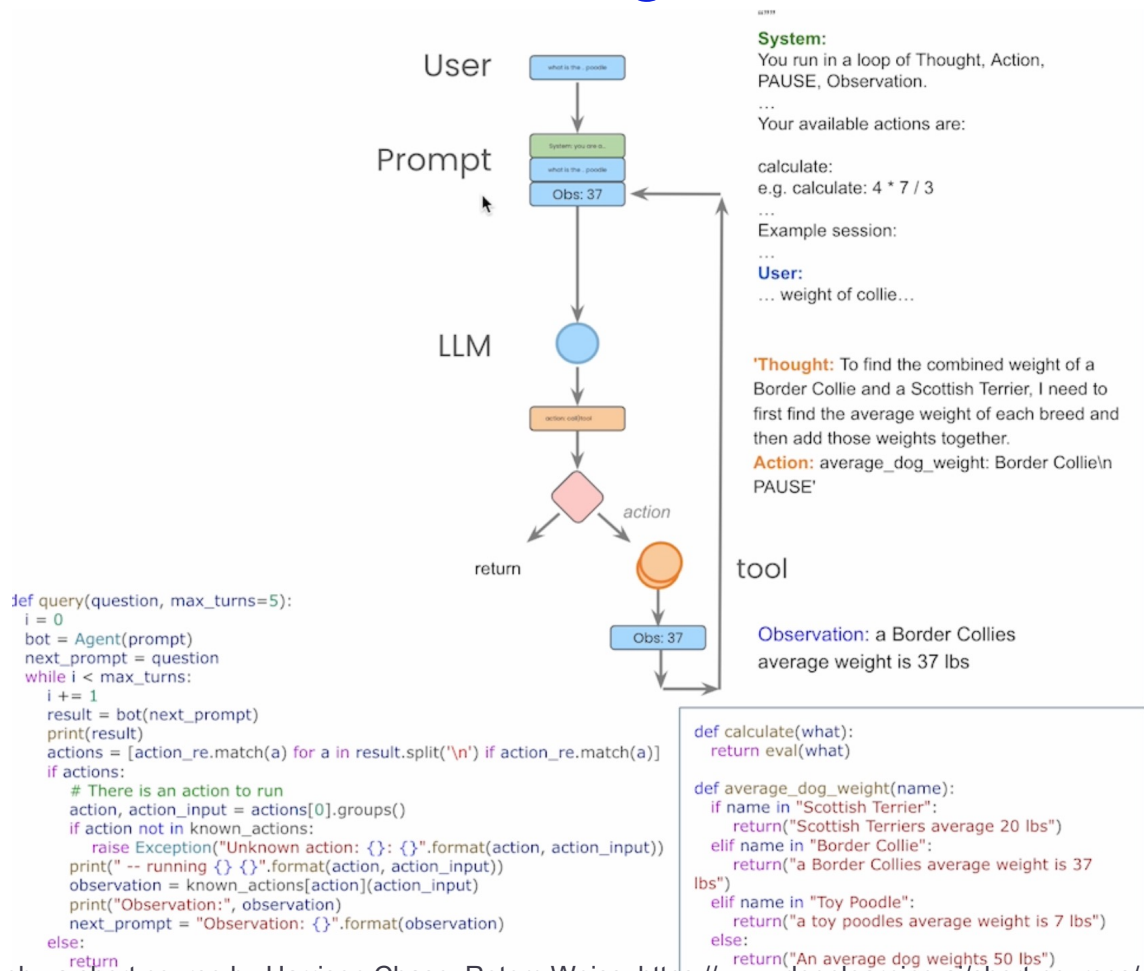
-- running calculate 37 + 20

Observation: 57

Answer: The combined weight of a Border Collie and a Scottish Terrier is 57 lbs.



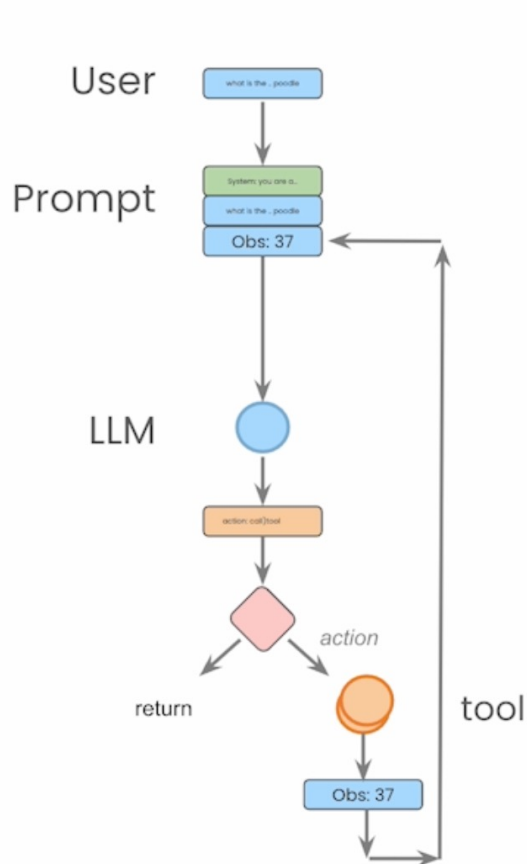
# Building a ReAct-based AI Agent w/ vanilla LangChain





# AI Agents in LangGraph

# Recap: LangChain Prompt Templates



Prompt templates allow reusable prompts

```
from langchain.prompts import PromptTemplate
prompt_template = PromptTemplate.from_template(
    "Tell me a {adjective} joke about {content}."
```

There are also prompts for agents available in the hub:

```
prompt = hub.pull("hwchase17/react")
```

```
https://smith.langchain.com/hub/hwchase17/react
```

# LangChain Prompt Templates

The image shows a screenshot of the LangChain Prompt Templates repository on GitHub. The main focus is on a specific prompt template named "PromptTemplate".

**Repository Overview:**

- Repository: `hwchase17/react` (Public)
- Updated: 7 months ago
- Stats: 15 likes, 24.2k views, 663k downloads

**PromptTemplate Details:**

- Action: the action to take, should be one of `[[tool_names]]`
- Action Input: the input to the action
- Observation: the result of the action
- ... (this Thought/Action/Action Input/Observation can repeat N times)
- Thought: I now know the final answer
- Final Answer: the final answer to the original input question

**Example Usage:**

```
Begin!
```

```
Question: {input}
Thought: {agent_scratchpad}
```

**Code Snippet:**

```
from langchain import hub
prompt = hub.pull("hwchase17/react")
```

**Comments:**

- Matteo Notaro (2 months ago): i've a problem with this prompt. If used with a multi-input tools it does pass only one input. It's by design this way?

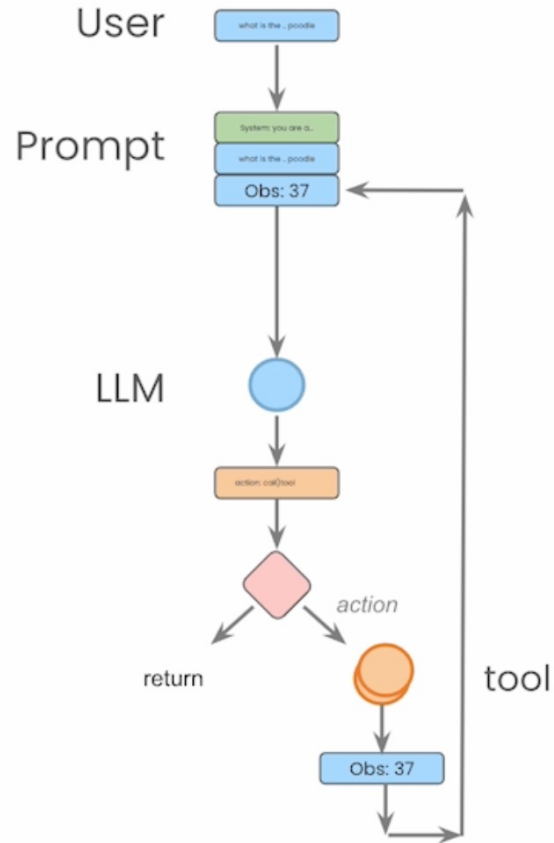
**Other Templates:**

- `homanp/superagent`: This prompt ads sequential function calling to models other than GPT-0613. (118 likes, 50.3k views, 27.4k downloads)
- `hardkothari/prompt-maker`: Convert your small and lazy prompt into a detailed and better prompts with this template. (106 likes, 29k views, 8.78k downloads)
- `rlm/rag-prompt`: This is a prompt for retrieval-augmented-generation. It is useful for chat, QA, or other applications that rely on passing context to an LLM. (91 likes, 92.7k views, 954k downloads)
- `smithing-gold/assumption-checker`: Assert whether assumptions are made in a user's query and provide follow up questions to debunk their claims. (63 likes, 15.3k views, 2.12k downloads)

**Navigation and Filtering:**

- Search: Search for prompts, use cases, models...
- Filters: Use Cases (Agent simulations, Agents, Autonomous agents, Chatbots, Classification, Code understanding, Code writing, Evaluation, Extraction, Interacting with AP..., Multi-modal, QA over documents, Self-checking, SQL, Summarization, Tagging), Type (ChatPromptTem..., StringPromptTem..., StructuredPrompt), Language (Chinese, English, French, German, Russian, Spanish), Model (anthropic:claude-2).
- Sort: Top Favored, Top Viewed, Top Downloaded, Recently Updated.

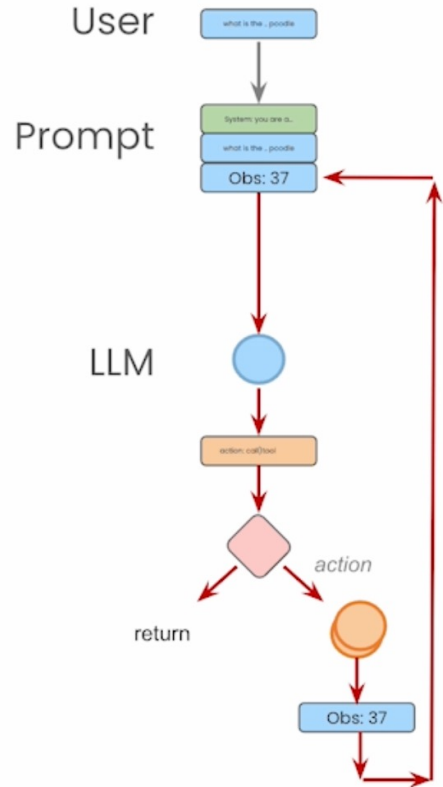
# LangChain Tools



```
# get a tool from the library
from langchain_community.tools.tavily_search import \
    TavilySearchResults
tool = TavilySearchResults(max_results=2)

self.tools = {t.name: t for t in tools}
self.model = model.bind_tools([tool])
```

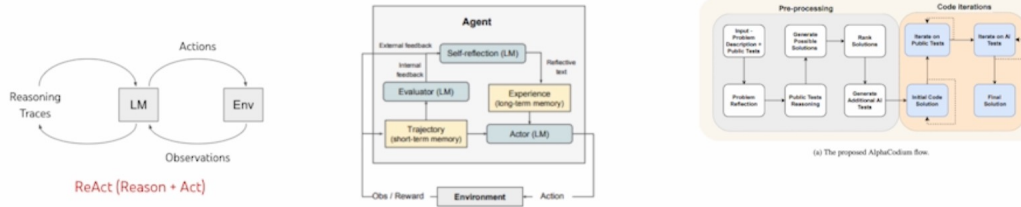
# What's New in LangGraph for building AI Agents ?



- Cyclic Graphs
- Persistence
- Human-in-the-loop

# What's New in LangGraph for building AI Agents ?

## Graphs



- LangGraph is an extension of LangChain that supports graphs.
- Single and Multi-agent flows are described and represented as graphs.
- Allows for extremely controlled “flows”
- Built-in persistence allows for human-in-the-loop workflows

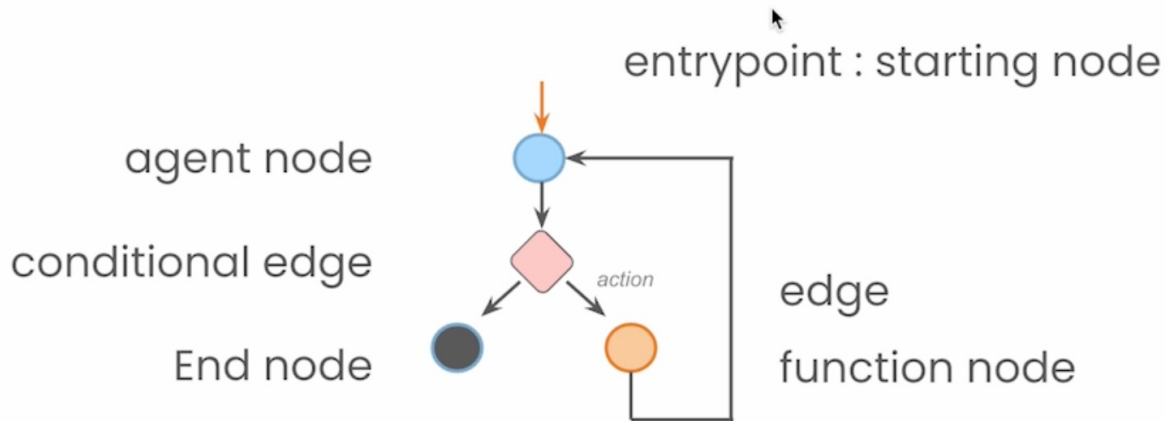
# What's New in LangGraph for building AI Agents ?

## Graphs

  **Nodes:** Agents or functions

 **Edges:** connect nodes

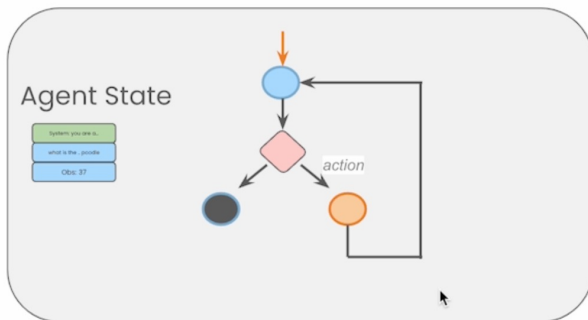
 **Conditional edges:** decisions





# What's New in LangGraph for building AI Agents ?

## Data/State



- Agent State is accessible to all parts of the graph
- It is local to the graph
- Can be stored in a persistence layer

## Simple

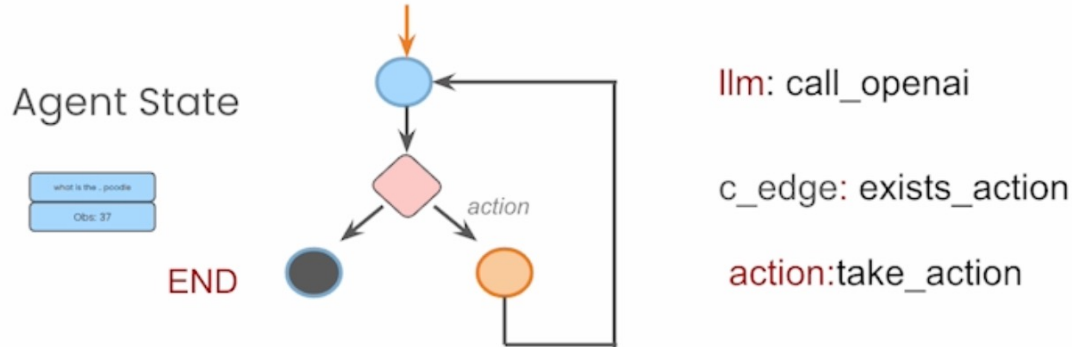
```
class AgentState(TypedDict):  
    messages: Annotated[Sequence[BaseMessage], operator.add]
```

## Complex

```
class AgentState(TypedDict):  
    input: str  
    chat_history: list[BaseMessage]  
    agent_outcome: Union[AgentAction, AgentFinish, None]  
    intermediate_steps: Annotated[list[tuple[AgentAction, str]], operator.add]
```

# What's New in LangGraph for building AI Agents ?

## CODE



## State

```
class AgentState(TypedDict):  
    messages: Annotated[list[AnyMessage], operator.add]
```

# AI Agents in LangGraph

```
from dotenv import load_dotenv, find_dotenv
```

```
_ = load_dotenv(find_dotenv())
```

```
from langgraph.graph import StateGraph, END
from typing import TypedDict, Annotated
import operator
from langchain_core.messages import AnyMessage, SystemMessage, HumanMessage
from langchain_openai import ChatOpenAI
from langchain_community.tools.tavily_search import TavilySearchResults
```

```
tool = TavilySearchResults(max_results=2)
print(type(tool))
print(tool.name)
```

```
<class 'langchain_community.tools.tavily_search.tool.TavilySearchResults'>
tavily_search_results_json
```

```
class AgentState(TypedDict):
    messages: Annotated[list[AnyMessage], operator.add]
```

```
class Agent:

    def __init__(self, model, tools, system=""):
        self.system = system
        graph = StateGraph(AgentState)
        graph.add_node("llm", ...)
        graph.add_node("action", ...)
        graph.add_conditional_edges(
            "llm",
            ...,
            {True: "action", False: END}
        )
        graph.add_edge("action", "llm")
        graph.set_entry_point("llm")
```

```
class Agent:

    def __init__(self, model, tools, system=""):
        self.system = system
        graph = StateGraph(AgentState)
        graph.add_node("llm", self.call_openai)
        graph.add_node("action", ...)
        graph.add_conditional_edges(
            "llm",
            ...,
            {True: "action", False: END}
        )
        graph.add_edge("action", "llm")
        graph.set_entry_point("llm")
        self.graph = graph.compile()
        self.tools = {t.name: t for t in tools}
        self.model = model.bind_tools(tools)

    def call_openai(self, state: AgentState):
        messages = state['messages']
        if self.system:
            messages = [SystemMessage(content=self.system)] + messages
        message = self.model.invoke(messages)
        return {'messages': [message]}

    def take_action(self, state: AgentState):
        tool_calls = state['messages'][-1].tool_calls
        results = []
        for t in tool_calls:
            print(f"Calling: {t}")
            result = self.tools[t['name']].invoke(t['args'])
            results.append(ToolMessage(tool_call_id=t['id'], name=t['name'], content=result))
        print("Back to the model!")
        return {'messages': results}
```

# AI Agents in LangGraph

```
class Agent:
    def __init__(self, model, tools, system=""):
        self.system = system
        graph = StateGraph(AgentState)
        graph.add_node("llm", self.call_openai)
        graph.add_node("action", self.take_action)
        graph.add_conditional_edges(
            "llm",
            self.exists_action,
            {True: "action", False: END}
        )
        graph.add_edge("action", "llm")
        graph.set_entry_point("llm")
        self.graph = graph.compile()
        self.tools = {t.name: t for t in tools}
        self.model = model.bind_tools(tools)

    def exists_action(self, state: AgentState):
        result = state['messages'][-1]
        return len(result.tool_calls) > 0

    def call_openai(self, state: AgentState):
        messages = state['messages']
        if self.system:
            messages = [SystemMessage(content=self.system)] + messages
        message = self.model.invoke(messages)
        return {'messages': [message]}

    def take_action(self, state: AgentState):
        tool_calls = state['messages'][-1].tool_calls
        results = []
        for t in tool_calls:
            print(f"Calling: {t}")
            result = self.tools[t['name']].invoke(t['args'])
            results.append(ToolMessage(tool_call_id=t['id'], name=t['name'],
                                      content=result))
        print("Back to the model!")
        return {'messages': results}
```

```
return len(result.tool_calls) > 0
```

```
def call_openai(self, state: AgentState):
    messages = state['messages']
    if self.system:
        messages = [SystemMessage(content=self.system)] + messages
    message = self.model.invoke(messages)
    return {'messages': [message]}

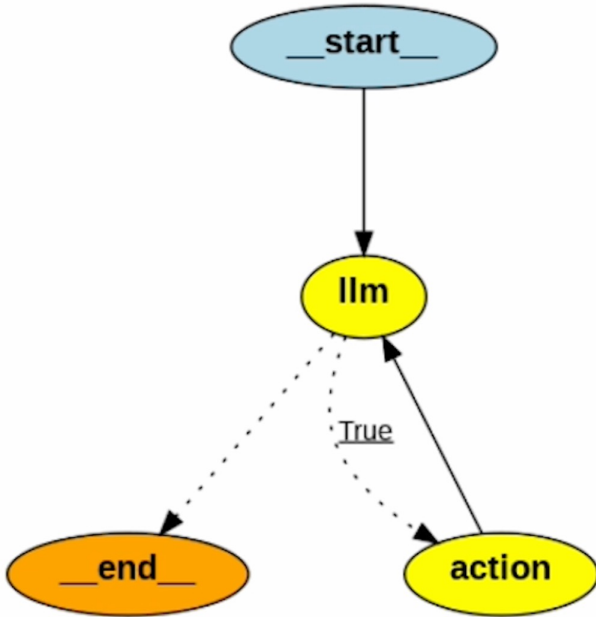
def take_action(self, state: AgentState):
    tool_calls = state['messages'][-1].tool_calls
    results = []
    for t in tool_calls:
        print(f"Calling: {t}")
        result = self.tools[t['name']].invoke(t['args'])
        results.append(ToolMessage(tool_call_id=t['id'], name=t['name'],
                                    content=result))
    print("Back to the model!")
    return {'messages': results}
```

```
prompt = """You are a smart research assistant. Use the search engine to
You are allowed to make multiple calls (either together or in sequence).
Only look up information when you are sure of what you want. \
If you need to look up some information before asking a follow up questi
"""
model = ChatOpenAI(model="gpt-4-turbo")
abot = Agent(model, [tool], system=prompt)
```

# AI Agents in LangGraph

```
from IPython.display import Image
```

```
Image(abot.graph.get_graph().draw_png())
```



```
messages = [HumanMessage(content="What is the weather in sf?")]  
result = abot.graph.invoke({"messages": messages})
```

```
Calling: {'name': 'tavily_search_results_json', 'args': {'query': 'current weather in San Francisco'}, 'id': 'call_3aeXDjFZbI2DQodx9490uQdk'}  
Back to the model!
```

```
result
```



# AI Agents in LangGraph

result

```
{'messages': [HumanMessage(content='What is the weather in sf?'),
  AIMessage(content='', additional_kwargs={'tool_calls': [{'id': 'call_3aeXDJfZbI2DQodx9490uQdk', 'function': {'arguments': '{"query": "current weather in San Francisco"}', 'name': 'tavily_search_results_json', 'type': 'function'}]}], response_metadata={'token_usage': {'completion_tokens': 22, 'prompt_tokens': 153, 'total_tokens': 175}, 'model_name': 'gpt-4-turbo', 'system_fingerprint': 'fp_76f018034d', 'finish_reason': 'tool_calls', 'logprobs': None, id='run-d20dbbbb-f650-428d-a0b0-2f676bd34da8-0', tool_calls=[{'name': 'tavily_search_results_json', 'args': {'query': 'current weather in San Francisco'}, 'id': 'call_3aeXDJfZbI2DQodx9490uQdk'}]}],
```

```
  ToolMessage(content='{\'url\': \'https://www.weatherapi.com/\', \'content\': \"\{\'location\': {\\'name\': \'San Francisco\', \'region\': \'California\', \'country\': \'United States of America\', \'lat\': 37.78, \'lon\': -122.42, \'tz_id\': \'America/Los_Angeles\', \'localtime_epoch\': 1712954805, \'localtime\': \'2024-04-12 13:46\', \'current\': {\\'ast_updated_epoch\': 1712954700, \'last_updated\': \'2024-04-12 13:45\', \'temp_c\': 15.6, \'temp_f\': 60.1, \'is_day\': 1, \'condition\': {\\'text\': \'Partly cloudy\', \'icon\': \'//cdn.weatherapi.com/weather/64x64/day/116.png\', \'code\': 1003}, \'wind_mph\': 15.0, \'wind_kph\': 24.1, \'wind_degree\': 270, \'wind_dir\': \'W\', \'pressure_mb\': 1009.0, \'pressure_in\': 29.8, \'precip_mm\': 0.0, \'precip_in\': 0.0, \'humidity\': 60, \'cloud\': 50, \'feelslike_c\': 15.6, \'feelslike_f\': 60.1, \'vis_km\': 16.0, \'vis_miles\': 9.0, \'uv\': 4.0, \'gust_mph\': 19.5, \'gust_kph\': 31.3}}\", \'url\': \'https://www.wunderground.com/hourly/us/ca/san-francisco/94129/date/2024-04-12\', \'content\': \'San Francisco Weather Forecasts. Weather Underground provides local & long-range weather forecasts, weatherreports, maps & tropical weather conditions for the San Francisco area. ... Friday 04/12 ...\'}, name='tavily_search_results_json', tool_call_id='call_3aeXDJfZbI2DQodx9490uQdk'),
```

```
  AIMessage(content='The current weather in San Francisco is partly cloudy with a temperature of 60.1°F (15.6°C). The wind is coming from the west at 15.0 mph (24.1 kph), and the atmospheric pressure is 1009.0 mb. The humidity level is at 60%, and visibility is approximately 9.0 miles (16.0 km).', response_metadata={'token_usage': {'completion_tokens': 80, 'prompt_tokens': 607, 'total_tokens': 687}, 'model_name': 'gpt-4-turbo'})
```

```
result['messages'][-1].content
```

'The current weather in San Francisco is partly cloudy with a temperature of 60.1°F (15.6°C). The wind is coming from the west at 15.0 mph (24.1 kph), and the atmospheric pressure is 1009.0 mb. The humidity level is at 60%, and visibility is approximately 9.0 miles (16.0 km).'

```
messages = [HumanMessage(content="What is the weather in SF and LA?")]
result = abot.graph.invoke({"messages": messages})
```

```
Calling: {'name': 'tavily_search_results_json', 'args': {'query': 'current weather in San Francisco'}, 'id': 'call_3aeXDJfZbI2DQodx9490uQdk'}
Calling: {'name': 'tavily_search_results_json', 'args': {'query': 'current weather in Los Angeles'}, 'id': 'call_w7IBLdVHXmdUIpHgybMaE5Vw'}
Back to the model!
```

```
result['messages'][-1].content
```

'Here is the current weather information for San Francisco and Los Angeles:  
San Francisco: Temperature: 15.6°C (60.1°F) Condition: Partly cloudy  
Wind: 15.0 mph from the west Humidity: 60% Visibility: 16 km (9 miles)  
Los Angeles: Temperature: 16.1°C (61.0°F) Condition: Overcast  
Wind: 4.3 mph from the southeast Humidity: 67% Visibility: 16 km (9 miles)'





# AI Agents in LangGraph

```
query = "Who won the super bowl in 2024? What is the GDP of that state?"  
messages = [HumanMessage(content=query)]  
result = abot.graph.invoke({"messages": messages})
```

```
Calling: {'name': 'tavily_search_results_json', 'args': {'query': '2024  
Super Bowl winner'}, 'id': 'call_8KigXEXM5GShFhcaLjJLMe7Z'}
```

Back to the model!

```
Calling: {'name': 'tavily_search_results_json', 'args': {'query': 'GDP  
of Missouri 2023'}, 'id': 'call_q3ynNhiv6p0mae07KaN6tuxG'}
```

Back to the model!

```
result['messages'][-1].content
```

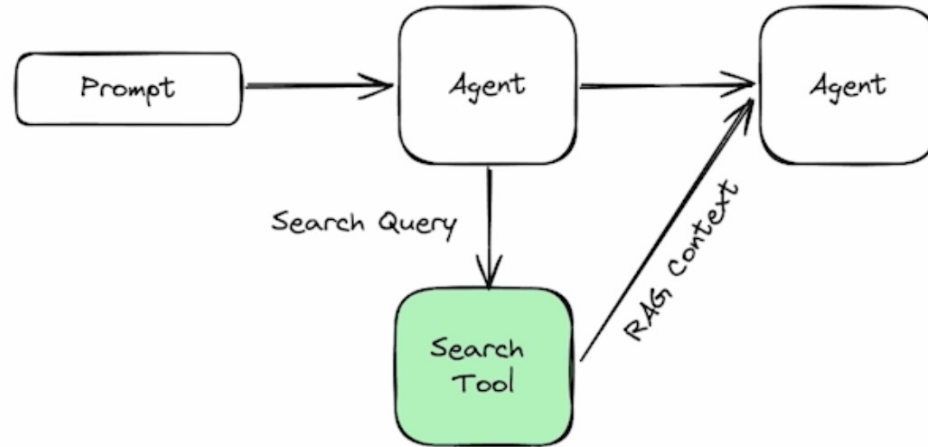
```
'The Kansas City Chiefs won the Super Bowl in 2024.\n\nThe GDP of Missouri, where the Kansas City Chiefs are based, was $423.6 billion in the 3rd quarter of 2023.'
```



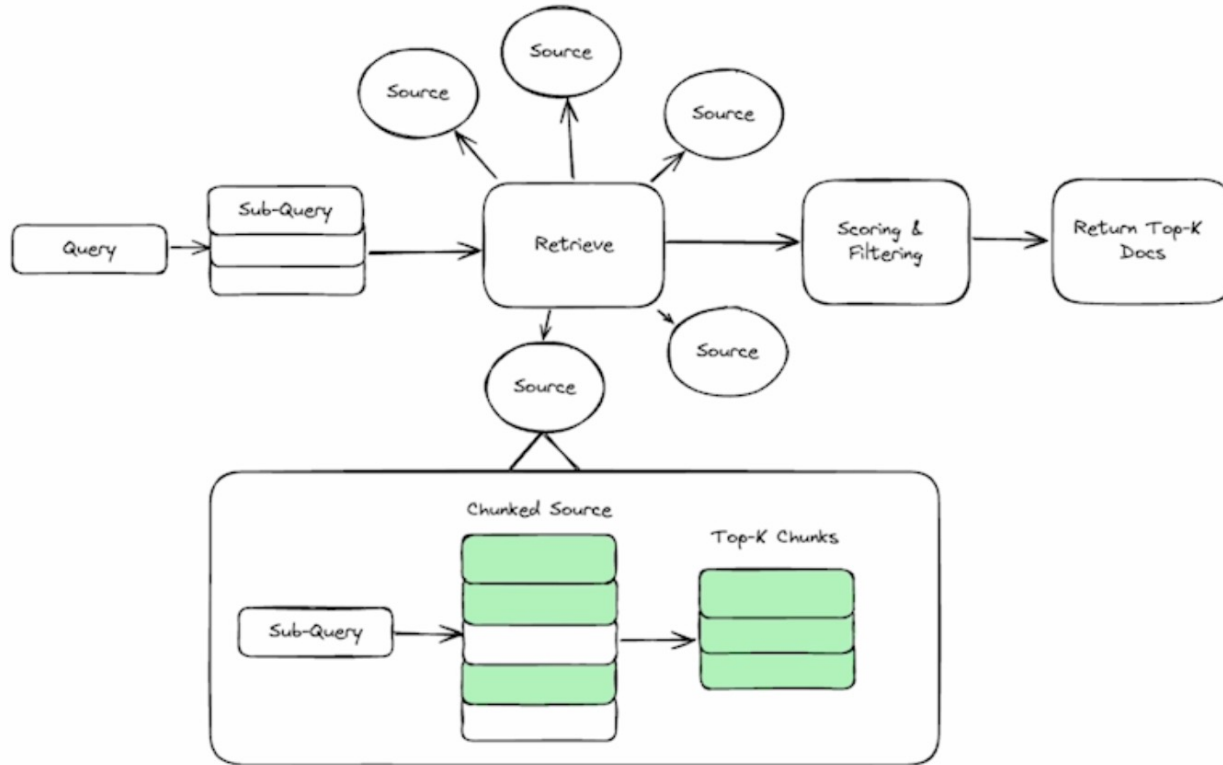


# AI Agents in LangGraph

## Why Search Tool



# AI Agents in LangGraph



# AI Agents in LangGraph

```
from dotenv import load_dotenv
import os
from tavily import TavilyClient

# load environment variables from .env file
_ = load_dotenv()

# connect
client = TavilyClient(api_key=os.environ.get("TAVILY_API_KEY"))

# run search
result = client.search("What is in Nvidia's new Blackwell GPU?",
                       include_answer=True)

# print the answer
result["answer"]
```

"The Nvidia Blackwell GPU features 208 billion transistors and connects two chip parts with a super-fast link between them for a cohesive architecture. It also includes an Advanced Transformer Engine, an AI engine designed to double Blackwell's processing power for more efficient task handling."

# AI Agents in LangGraph

```
# choose location (try to change to your own city!)
```

```
city = "San Francisco"
```

```
query = f"""
```

```
    what is the current weather in {city}?
```

```
    Should I travel there today?
```

```
    """
```

```
import requests
from bs4 import BeautifulSoup
from duckduckgo_search import DDGS
import re
```

```
ddg = DDGS()
```

```
def search(query, max_results=3):
    results = ddg.text(query, max_results=max_results)
    return [i["href"] for i in results]
```

```
for i in search(query):
    print(i)
```

```
https://www.accuweather.com/en/us/san-francisco/94103/current-weather/347629
```

```
https://www.accuweather.com/en/us/san-francisco/94103/weather-forecast/347629
```

```
https://www.accuweather.com/en/us/san-francisco/94103/air-travel-weather/347629
```

```
def scrape_weather_info(url):
    """Scrape content from the given URL"""
    if not url:
        return "Weather information could not be found."
```

```
# fetch data
```

```
headers = {'User-Agent': 'Mozilla/5.0'}
```

```
response = requests.get(url, headers=headers)
```

```
if response.status_code != 200:
```

```
    return "Failed to retrieve the webpage."
```

```
# parse result
```

```
soup = BeautifulSoup(response.text, 'html.parser')
```

```
return soup
```

```
# use DuckDuckGo to find a websites and take the first result
```

```
url = search(query)[0]
```

```
# scrape first wesbsite
```

```
soup = scrape_weather_info(url)
```

```
print(f"Website: {url}\n\n")
```

```
print(soup)
```

```
Website: https://www.accuweather.com/en/us/san-francisco/94103/current-weather/347629
```

```
<!DOCTYPE html>
```

```
<html class="accuweather" lang="en-us">
```

```
<head>
```

```
<meta content="IE=edge,chrome=1" http-equiv="X-UA-Compatible"/>
```



# AI Agents in LangGraph

```
weather_data = []
for tag in soup.find_all(['h1', 'h2', 'h3', 'p']):
    text = tag.get_text(" ", strip=True)
    weather_data.append(text)

# combine all elements into a single string
weather_data = "\n".join(weather_data)

# remove all spaces from the combined text
weather_data = re.sub(r'\s+', ' ', weather_data)

print(f"Website: {url}\n\n")
print(weather_data)
```

Website: <https://www.accuweather.com/en/us/san-francisco/94103/current-weather/347629>

San Francisco, CA San Francisco California Around the Globe Around the  
Globe Hurricane Tracker Severe Weather Radar & Maps News & Features Ast  
ronomy Business Climate Health Recreation Sports Travel Video Winter Ce  
nter Current Weather 2:24 PM Day Max UV Index 7 High Wind WSW 13 mph Wi  
nd Gusts 25 mph Probability of Precipitation 25% Probability of Thunder  
storms 3% Precipitation 0.00 in Cloud Cover 69% Morning Afternoon Night  
Wind 5 14 mph Wind Gusts 17 mph Probability of Precipitation 91% Probab  
ility of Thunderstorms 17% Precipitation 0.33 in Rain 0.33 in Hours of  
Precipitation 5 Hours of Rain 5 Cloud Cover 93% Evening Overnight Sun &  
Moon Temperature History Further Ahead Further Ahead Hourly Daily Month  
ly Around the Globe Around the Globe Hurricane Tracker Severe Weather R  
adar & Maps News Video Winter Center Top Stories Severe Weather Powerfu  
l storms unleash severe flooding, tornadoes from Texas to Maine 1 hour  
ago Weather Forecasts Chilly winds to follow storms, flooding rain in e  
astern US 2 hours ago Severe Weather Severe storms to threaten weekend  
plans from Midwest to Northeast 5 hours ago Winter Weather Mid-April st  
orm to eye California with rain, mountain snow 2 hours ago Severe Weath  
er Severe weather, tornado threat to build on Plains, Mississippi Valle  
y 2 hours ago Featured Stories Weather News How wet weather can be dang  
erous to your pets Climate March 2024 hottest on record, EU climate se

```
import json
from pygments import highlight, lexers, formatters

# parse JSON
parsed_json = json.loads(data.replace("'", ""))

# pretty print JSON with syntax highlighting
formatted_json = json.dumps(parsed_json, indent=4)
colorful_json = highlight(formatted_json,
                           lexers.JsonLexer(),
                           formatters.TerminalFormatter())

print(colorful_json)
```

```
{
  "location": {
    "name": "San Francisco",
    "region": "California",
    "country": "United States of America",
    "lat": 37.78,
    "lon": -122.42,
    "tz_id": "America/Los_Angeles",
    "localtime_epoch": 1712957187,
    "localtime": "2024-04-12 14:26"
  },
  "current": {
    "last_updated_epoch": 1712956500,
    "last_updated": "2024-04-12 14:15",
    "temp_c": 15.0,
    "temp_f": 59.0,
    "is_day": 1,
    "condition": {
      "text": "Partly cloudy",
      "icon": "https://cdn.weatherapi.com/weather/64x64/day/116.png",
      "code": 1003
    }
  }
}
```

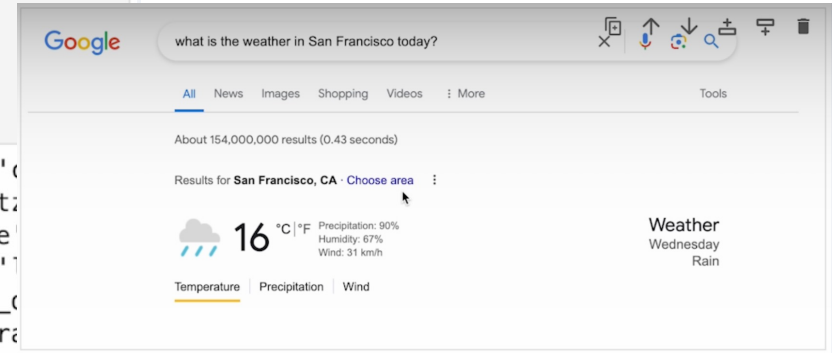
# AI Agents in LangGraph

```
# run search
result = client.search(query, max_results=1)

# print first result
data = result["results"][0]["content"]

print(data)

{'location': {'name': 'San Francisco', 'region': 'California', 'country': 'United States of America', 'lat': 37.78, 'lon': -122.42, 'timezone': 'America/Los_Angeles', 'localtime_epoch': 1712957187, 'localtime': '2024-04-12 14:26'}, 'current': {'last_updated_epoch': 1712956500, 'last_updated': '2024-04-12 14:15', 'temp_c': 15.0, 'temp_f': 59.0, 'is_day': 1, 'condition': {'text': 'Partly cloudy', 'icon': '//cdn.weather.com/weather/64x64/day/116.png', 'code': 1003}, 'wind_mph': 16.1, 'wind_kph': 25.9, 'wind_degree': 210, 'wind_dir': 'SSW', 'pressure_mb': 1009.0, 'pressure_in': 29.8, 'precip_mm': 0.0, 'precip_in': 0.0, 'humidity': 60, 'cloud': 75, 'feelslike_c': 13.6, 'feelslike_f': 56.5, 'vis_km': 16.0, 'vis_miles': 9.0, 'uv': 4.0, 'gust_mph': 20.6, 'gust_kph': 33.1}}
```





# AI Agents in LangGraph

```
from dotenv import load_dotenv, find_dotenv
```

```
_ = load_dotenv(find_dotenv())
```

```
from langgraph.graph import StateGraph, END
from typing import TypedDict, Annotated
import operator
from langchain_core.messages import AnyMessage, SystemMessage, HumanMessage
from langchain_openai import ChatOpenAI
from langchain_community.tools.tavily_search import TavilySearchResults
```

```
tool = TavilySearchResults(max_results=2)
```

```
class AgentState(TypedDict):
    messages: Annotated[list[AnyMessage], operator.add]
```

```
prompt = """You are a smart research assistant. Use the search engine to
You are allowed to make multiple calls (either together or in sequence).
Only look up information when you are sure of what you want. \
If you need to look up some information before asking a follow up question
"""
```

```
model = ChatOpenAI(model="gpt-4-turbo")
abot = Agent(model, [tool], system=prompt, checkpointer=memory)
```

```
from langgraph.checkpoint.sqlite import SqliteSaver
```

```
memory = SqliteSaver.from_conn_string(":memory:")
```

```
class Agent:
    def __init__(self, model, tools, checkpointer, system=""):
        self.system = system
        graph = StateGraph(AgentState)
        graph.add_node("llm", self.call_openai)
        graph.add_node("action", self.take_action)
        graph.add_conditional_edges("llm", self.exists_action, {True: "action", False: "llm"})
        graph.add_edge("action", "llm")
        graph.set_entry_point("llm")
        self.graph = graph.compile(checkpointer=checkpointer)
        self.tools = {t.name: t for t in tools}
        self.model = model.bind_tools(tools)

    def call_openai(self, state: AgentState):
        messages = state['messages']
        if self.system:
            messages = [SystemMessage(content=self.system)] + messages
        message = self.model.invoke(messages)
        return {'messages': [message]}

    def exists_action(self, state: AgentState):
        result = state['messages'][-1]
        return len(result.tool_calls) > 0

    def take_action(self, state: AgentState):
        tool_calls = state['messages'][-1].tool_calls
        results = []
        for t in tool_calls:
            print(f"Calling: {t}")
            result = self.tools[t['name']].invoke(t['args'])
            results.append(ToolMessage(tool_call_id=t['id'], name=t['name'], content=result)))
```



# AI Agents in LangGraph

```
messages = [HumanMessage(content="What is the weather in sf?")]
```

```
thread = {"configurable": {"thread_id": "1"}}
```

```
for event in abot.graph.stream({"messages": messages}, thread):  
    for v in event.values():  
        print(v['messages'])
```

```
[AIMessage(content='', additional_kwargs={'tool_calls': [{'id': 'call_v0EHAAC8cE33PI4vQUaLUyMB', 'function': {'arguments': '{"query": "current weather in San Francisco"}', 'name': 'tavily_search_results_json'}, 'type': 'function'}]}, response_metadata={'token_usage': {'completion_tokens': 22, 'prompt_tokens': 153, 'total_tokens': 175}, 'model_name': 'gpt-4-turbo', 'system_fingerprint': 'fp_76f018034d', 'finish_reason': 'tool_calls', 'logprobs': None}, id='run-07978a9a-6798-480b-a8f4-b99c439e6aef-0', tool_calls=[{'name': 'tavily_search_results_json', 'args': {'query': 'current weather in San Francisco'}, 'id': 'call_v0EHAAC8cE33PI4vQUaLUyMB'}]))  
Calling: {'name': 'tavily_search_results_json', 'args': {'query': 'current weather in San Francisco'}, 'id': 'call_v0EHAAC8cE33PI4vQUaLUyMB'}  
Back to the model!
```

```
[ToolMessage(content='[{"url": "https://www.weatherapi.com/", "content": "{\\"location\\": {\\"name\\": \\"San Francisco\\", \\"region\\": \\"California\\", \\"country\\": \\"United States of America\\", \\"lat\\": 37.78, \\"lon\\": -122.42, \\"tz_id\\": \\"America/Los_Angeles\\", \\"localtime_epoch\\": 1712959432, \\"localtime\\": \\"2024-04-12 15:03\\", \\"current\\": {\\"last_updated_epoch\\": 1712959200, \\"last_updated\\": \\"2024-04-12 15:00\\", \\"temp_c\\": 15.0, \\"temp_f\\": 59.0, \\"is_day\\": 1, \\"condition\\": \\"Partly cloudy\\", \\"icon\\": "\\/cdn.weatherapi.com/weather/64x64/day/116.png\\", \\"code\\": 1003}, \\"wind_mph\\": 16.1, \\"wind_kph\\": 25.9, \\"wind_degree\\": 210, \\"wind_dir\\": \\"SSW\\", \\"pressure_mb\\": 1009.0, \\"pressure_in\\": 29.8, \\"precip_mm\\": 0.0, \\"precip_in\\": 0.0, \\"humidity\\": 60, \\"cloud\\": 75, \\"feelslike_c\\": 13.6, \\"feelslike_f\\": 56.4, \\"vis_km\\": 16.0, \\"vis_miles\\": 9.0, \\"uv\\": 4.0, \\"gust_mph\\": 20.6, \\"gust_kph\\": 33.1}}}', {"url": "https://www.weathertab.com/en/c/e/04/united-states/california/san-francisco/", "content": "\\"Explore comprehensive April 2024 weather forecasts for San Francisco, including daily high and low temperatures, precipitation risks, and monthly temperature trends. Featuring detailed day-by-day forecasts, dynamic graphs of daily rain probabilities, and temperature trends to help you plan ahead. ... 12 64°F 49°F 18°C 9°C 49% 13 62°F 49°F ...\\"}], name='tavily_search_results_json', tool_call_id='call_v0EHAAC8cE33PI4vQUaLUyMB')]
```

```
pe: function]], response_metadata={'token_usage': {'completion_tokens': 22, 'prompt_tokens': 153, 'total_tokens': 175}, 'model_name': 'gpt-4-turbo', 'system_fingerprint': 'fp_76f018034d', 'finish_reason': 'tool_calls', 'logprobs': None}, id='run-07978a9a-6798-480b-a8f4-b99c439e6aef-0', tool_calls=[{'name': 'tavily_search_results_json', 'args': {'query': 'current weather in San Francisco'}, 'id': 'call_v0EHAAC8cE33PI4vQUaLUyMB'}]))  
Calling: {'name': 'tavily_search_results_json', 'args': {'query': 'current weather in San Francisco'}, 'id': 'call_v0EHAAC8cE33PI4vQUaLUyMB'}  
Back to the model!
```

```
[ToolMessage(content='[{"url": "https://www.weatherapi.com/", "content": "{\\"location\\": {\\"name\\": \\"San Francisco\\", \\"region\\": \\"California\\", \\"country\\": \\"United States of America\\", \\"lat\\": 37.78, \\"lon\\": -122.42, \\"tz_id\\": \\"America/Los_Angeles\\", \\"localtime_epoch\\": 1712959432, \\"localtime\\": \\"2024-04-12 15:03\\", \\"current\\": {\\"last_updated_epoch\\": 1712959200, \\"last_updated\\": \\"2024-04-12 15:00\\", \\"temp_c\\": 15.0, \\"temp_f\\": 59.0, \\"is_day\\": 1, \\"condition\\": \\"Partly cloudy\\", \\"icon\\": "\\/cdn.weatherapi.com/weather/64x64/day/116.png\\", \\"code\\": 1003}, \\"wind_mph\\": 16.1, \\"wind_kph\\": 25.9, \\"wind_degree\\": 210, \\"wind_dir\\": \\"SSW\\", \\"pressure_mb\\": 1009.0, \\"pressure_in\\": 29.8, \\"precip_mm\\": 0.0, \\"precip_in\\": 0.0, \\"humidity\\": 60, \\"cloud\\": 75, \\"feelslike_c\\": 13.6, \\"feelslike_f\\": 56.4, \\"vis_km\\": 16.0, \\"vis_miles\\": 9.0, \\"uv\\": 4.0, \\"gust_mph\\": 20.6, \\"gust_kph\\": 33.1}}}', {"url": "https://www.weathertab.com/en/c/e/04/united-states/california/san-francisco/", "content": "\\"Explore comprehensive April 2024 weather forecasts for San Francisco, including daily high and low temperatures, precipitation risks, and monthly temperature trends. Featuring detailed day-by-day forecasts, dynamic graphs of daily rain probabilities, and temperature trends to help you plan ahead. ... 12 64°F 49°F 18°C 9°C 49% 13 62°F 49°F ...\\"}], name='tavily_search_results_json', tool_call_id='call_v0EHAAC8cE33PI4vQUaLUyMB')]
```

```
[AIMessage(content='The current weather in San Francisco is partly cloudy with a temperature of 59°F (15°C). The wind is blowing from the south-southwest at 16.1 mph (25.9 kph). The humidity is at 60%, and the visibility is 9 miles (16 km).', response_metadata={'token_usage': {'completion_tokens': 62, 'prompt_tokens': 643, 'total_tokens': 705}, 'model_name': 'gpt-4-turbo', 'system_fingerprint': 'fp_76f018034d', 'finish_reason': 'stop', 'logprobs': None}, id='run-89f58e44-dfa6-4145-93f2-3fdb5763a8-0')]
```

# AI Agents in LangGraph

```
messages = [HumanMessage(content="What about in la?")]
thread = {"configurable": {"thread_id": "1"}}
for event in abot.graph.stream({"messages": messages}, thread):
    for v in event.values():
        print(v)

{'messages': [AIMessage(content='', additional_kwargs={'tool_calls':
[{'id': 'call_xg0cm52ju8a3Ib6g4dYGzk29', 'function': {'arguments': {'q
uery': 'current weather in Los Angeles'}}, 'name': 'tavily_search_result
s_json'}, 'type': 'function'}}], response_metadata={'token_usage': {'co
mpletion_tokens': 22, 'prompt_tokens': 717, 'total_tokens': 739}, 'mode
l_name': 'gpt-4-turbo', 'system_fingerprint': 'fp_76f018034d', 'finish
reason': 'tool_calls', 'logprobs': None}, id='run-aafe90d2-d6c0-4697-80
8d-5a83d187b050-0', tool_calls=[{'name': 'tavily_search_results_json',
'args': {'query': 'current weather in Los Angeles'}, 'id': 'call_xg0cm5
2ju8a3Ib6g4dYGzk29'}})]}]
Calling: {'name': 'tavily_search_results_json', 'args': {'query': 'curr
ent weather in Los Angeles'}, 'id': 'call_xg0cm52ju8a3Ib6g4dYGzk29'}
Back to the model!
{'messages': [ToolMessage(content='[{"url": "https://www.weatherapi.
com/", "content": "{\\"location\\": {\\"name\\": \\"Los Angeles\\", \\"regi
on\\": \\"California\\", \\"country\\": \\"United States of America\\", \\"lat
\\": 34.05, \\"lon\\": -118.24, \\"tz_id\\": \\"America/Los_Angeles\\", \\"loca
ltime_epoch\\": 1712959468, \\"localtime\\": \\"2024-04-12 15:04\\", \\"curr
ent\\": {\\"last_updated_epoch\\": 1712959200, \\"last_updated\\": \\"2024-04
-12 15:00\\", \\"temp_c\\": 16.7, \\"temp_f\\": 62.1, \\"is_day\\": 1, \\"condi
tion\\": {\\"text\\": \\"Overcast\\", \\"icon\\": \\"//cdn.weatherapi.com/weath
er/64x64/day/122.png\\", \\"code\\": 1009}, \\"wind_mph\\": 2.2, \\"wind_kph
\\": 3.6, \\"wind_degree\\": 177, \\"wind_dir\\": \\"S\\", \\"pressure_mb\\": 10
15.0, \\"pressure_in\\": 29.96, \\"precip_mm\\": 0.0, \\"precip_in\\": 0.0,
\\\"humidity\\": 67, \\"cloud\\": 100, \\"feelslike_c\\": 16.7, \\"feelslike_f
\\": 62.1, \\"vis_km\\": 16.0, \\"vis_miles\\": 9.0, \\"uv\\": 5.0, \\"gust_mph
\\": 13.1, \\"gust_kph\\": 21.1}}"}], {"url": "https://www.latimes.com/c
alifornia/story/2024-04-12/cold-weather-is-coming-los-angeles", "cont
ent": '\\"Jireh (they/them) is the 2023-24 fellow at the Los Angeles Tim
es and a queer Asian American writer and filmmaker born and raised in t
he San Gabriel Valley. More From the Los Angeles Times Business\\"}],
name='tavily_search_results_json', tool_call_id='call_xg0cm52ju8a3Ib6g4
dYGzk29'}})]}]
{'messages': [AIMessage(content='The current weather in Los Angeles is
overcast with a temperature of 62.1°F (16.7°C). The wind is blowing fro
m the south at a speed of 2.2 mph (3.6 kph). The humidity level is 67%,
and visibility is 9 miles (16 km).', response_metadata={'token_usage':
{'completion_tokens': 65, 'prompt_tokens': 1178, 'total_tokens': 1243},
'model_name': 'gpt-4-turbo', 'system_fingerprint': 'fp_a39722e138', 'fi
nish_reason': 'stop', 'logprobs': None}, id='run-0d2ea079-0efe-4a5c-8c6
4-d1f4a52cab3-0')}}]
```

```
l_name': 'gpt-4-turbo', 'system_fingerprint': 'fp_76f018034d', 'finish_
reason': 'tool_calls', 'logprobs': None}, id='run-aafe90d2-d6c0-4697-80
8d-5a83d187b050-0', tool_calls=[{'name': 'tavily_search_results_json',
'args': {'query': 'current weather in Los Angeles'}, 'id': 'call_xg0cm5
2ju8a3Ib6g4dYGzk29'}})]}]
Calling: {'name': 'tavily_search_results_json', 'args': {'query': 'curr
ent weather in Los Angeles'}, 'id': 'call_xg0cm52ju8a3Ib6g4dYGzk29'}
Back to the model!
{'messages': [ToolMessage(content='[{"url": "https://www.weatherapi.
com/", "content": "{\\"location\\": {\\"name\\": \\"Los Angeles\\", \\"regi
on\\": \\"California\\", \\"country\\": \\"United States of America\\", \\"lat
\\": 34.05, \\"lon\\": -118.24, \\"tz_id\\": \\"America/Los_Angeles\\", \\"loca
ltime_epoch\\": 1712959468, \\"localtime\\": \\"2024-04-12 15:04\\", \\"curr
ent\\": {\\"last_updated_epoch\\": 1712959200, \\"last_updated\\": \\"2024-04
-12 15:00\\", \\"temp_c\\": 16.7, \\"temp_f\\": 62.1, \\"is_day\\": 1, \\"condi
tion\\": {\\"text\\": \\"Overcast\\", \\"icon\\": \\"//cdn.weatherapi.com/weath
er/64x64/day/122.png\\", \\"code\\": 1009}, \\"wind_mph\\": 2.2, \\"wind_kph
\\": 3.6, \\"wind_degree\\": 177, \\"wind_dir\\": \\"S\\", \\"pressure_mb\\": 10
15.0, \\"pressure_in\\": 29.96, \\"precip_mm\\": 0.0, \\"precip_in\\": 0.0,
\\\"humidity\\": 67, \\"cloud\\": 100, \\"feelslike_c\\": 16.7, \\"feelslike_f
\\": 62.1, \\"vis_km\\": 16.0, \\"vis_miles\\": 9.0, \\"uv\\": 5.0, \\"gust_mph
\\": 13.1, \\"gust_kph\\": 21.1}}"}], {"url": "https://www.latimes.com/c
alifornia/story/2024-04-12/cold-weather-is-coming-los-angeles", "cont
ent": '\\"Jireh (they/them) is the 2023-24 fellow at the Los Angeles Tim
es and a queer Asian American writer and filmmaker born and raised in t
he San Gabriel Valley. More From the Los Angeles Times Business\\"}],
name='tavily_search_results_json', tool_call_id='call_xg0cm52ju8a3Ib6g4
dYGzk29'}})]}]
{'messages': [AIMessage(content='The current weather in Los Angeles is
overcast with a temperature of 62.1°F (16.7°C). The wind is blowing fro
m the south at a speed of 2.2 mph (3.6 kph). The humidity level is 67%,
and visibility is 9 miles (16 km).', response_metadata={'token_usage':
{'completion_tokens': 65, 'prompt_tokens': 1178, 'total_tokens': 1243},
'model_name': 'gpt-4-turbo', 'system_fingerprint': 'fp_a39722e138', 'fi
nish_reason': 'stop', 'logprobs': None}, id='run-0d2ea079-0efe-4a5c-8c6
4-d1f4a52cab3-0')}}]
```



# AI Agents in LangGraph

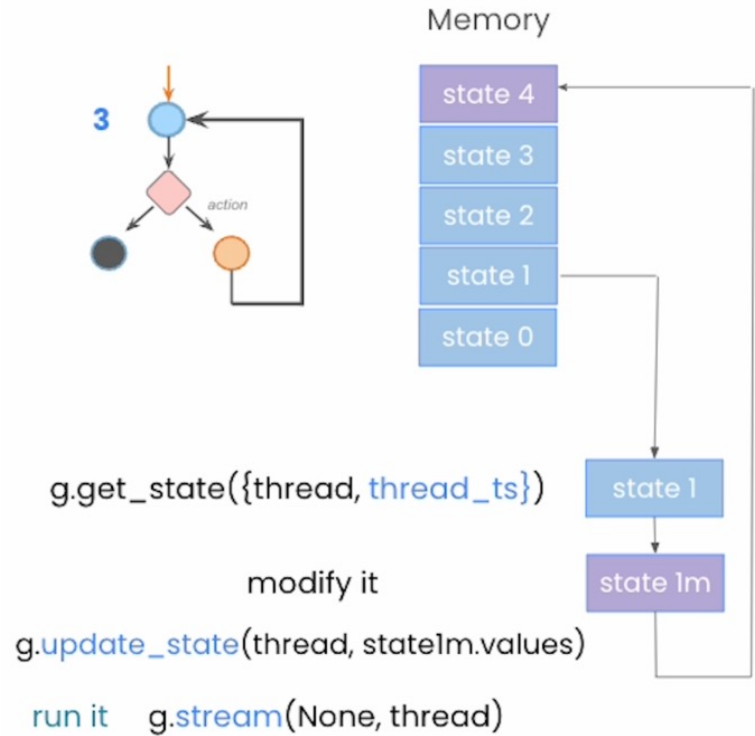
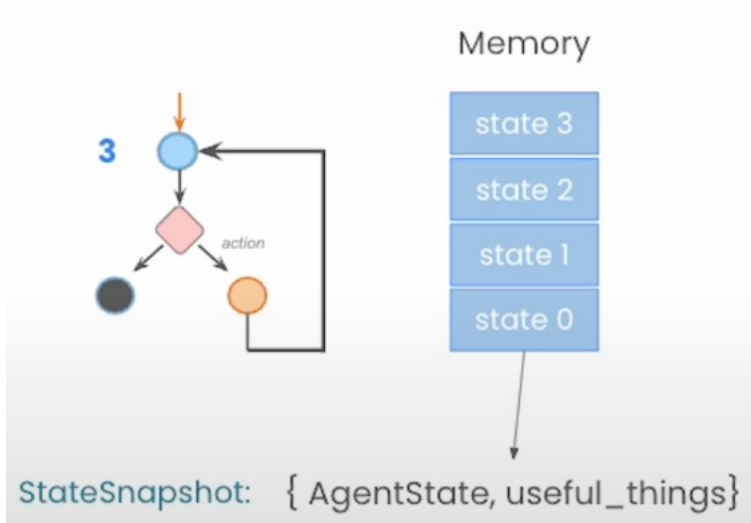
```
messages = [HumanMessage(content="Which one is warmer?")]
thread = {"configurable": {"thread_id": "1"}}
for event in abot.graph.stream({"messages": messages}, thread):
    for v in event.values():
        print(v)
```

```
{'messages': [AIMessage(content='Los Angeles is currently warmer with a temperature of 62.1°F (16.7°C) compared to San Francisco, which has a temperature of 59°F (15°C).', response_metadata={'token_usage': {'completion_tokens': 38, 'prompt_tokens': 1255, 'total_tokens': 1293}, 'model_name': 'gpt-4-turbo', 'system_fingerprint': 'fp_67e6987839', 'finish_reason': 'stop', 'logprobs': None}, id='run-810ffd03-8f21-48a5-ac30-ac5ac358f465-0')]}
```

```
messages = [HumanMessage(content="Which one is warmer?")]
thread = {"configurable": {"thread_id": "2"}}
for event in abot.graph.stream({"messages": messages}, thread):
    for v in event.values():
        print(v)
```

```
{'messages': [AIMessage(content='Could you please specify the two or more items you are comparing to determine which one is warmer?', response_metadata={'token_usage': {'completion_tokens': 20, 'prompt_tokens': 151, 'total_tokens': 171}, 'model_name': 'gpt-4-turbo', 'system_fingerprint': 'fp_a39722e138', 'finish_reason': 'stop', 'logprobs': None}, id='run-d84c5fe0-8ab5-4ca6-89e6-575467ed1ded-0')]}
```

# Tracking and Manipulating State Memory for Agents in LangGraph



# AI Agents in LangGraph

```
from langgraph.checkpoint.aiosqlite import AsyncSqliteSaver
```

```
memory = AsyncSqliteSaver.from_conn_string(":memory:")  
abot = Agent(model, [tool], system=prompt, checkpointer=memory)
```

```
messages = [HumanMessage(content="What is the weather in SF?")]  
thread = {"configurable": {"thread_id": "4"}}  
async for event in abot.graph.astream_events({"messages": messages}, thread_id="4",  
kind = event["event"]):  
    if kind == "on_chat_model_stream":  
        content = event["data"]["chunk"].content  
        if content:  
            # Empty content in the context of OpenAI means  
            # that the model is asking for a tool to be invoked  
            # So we only print non-empty content  
            print(content, end="|")
```

```
from langgraph.checkpoint.aiosqlite import AsyncSqliteSaver  
memory = AsyncSqliteSaver.from_conn_string(":memory:")  
abot = Agent(model, [tool], system=prompt, checkpointer=memory)
```

```
/Users/dlai.filming/Documents/GitHub/SC-Langchain-LangGraph-C5/.venv_11_1_LC5/lib/python3.11/site-packages/langchain_core/_api/beta_decorator.py:87: LangChainBetaWarning: This API is in beta and may change in the future.
```

```
warn_beta(  
Calling: {'name': 'tavily_search_results_json', 'args': {'query': 'current weather in San Francisco'}, 'id': 'call_AIkbXdDIVJcjReXmEGIvnUKp'}  
Back to the model!
```

```
The| current| weather| in| San| Francisco| is| partly| cloudy| with| a|  
temperature| of| 59|°F| (15|°C)|.| Winds| are| coming| from| the| south-|west| at| 16|.1| mph| (25|.9| k|ph)|.| The| humidity| is|  
| at| 60|%,| and| the| atmospheric| pressure| is| 100|9| mb|.| Visibility| is| approximately| 9| miles| (16| km)|.
```

# Human in the Loop w/ LangGraph



```
class Agent:
    def __init__(self, model, tools, system="", checkpointer=None):
        self.system = system
        graph = StateGraph(AgentState)
        graph.add_node("llm", self.call_openai)
        graph.add_node("action", self.take_action)
        graph.add_conditional_edges("llm", self.exists_action, {True: "a
        graph.add_edge("action", "llm")
        graph.set_entry_point("llm")
        self.graph = graph.compile(
            checkpointer=checkpointer,
            interrupt_before=["action"]
        )
        self.tools = {t.name: t for t in tools}
        self.model = model.bind_tools([tool])

    def call_openai(self, state: AgentState):
        messages = state['messages']
        if self.system:
            messages = [SystemMessage(content=self.system)] + messages
        message = self.model.invoke(messages)
        return {'messages': [message]}

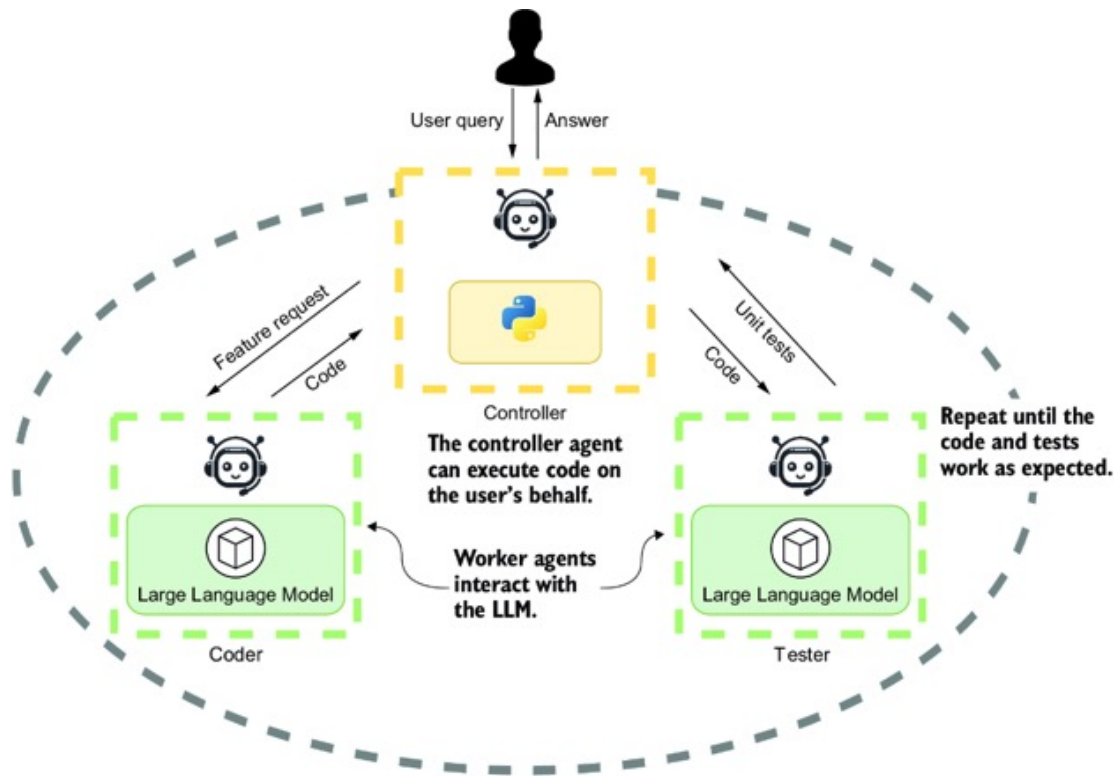
    def exists_action(self, state: AgentState):
        print(state)
        result = state['messages'][-1]
        return len(result.tool_calls) > 0

    def take_action(self, state: AgentState):
        tool_calls = state['messages'][-1].tool_calls
        results = []
        for t in tool_calls:
            print(f"Calling: {t}")
            result = self.tools[t['name']].invoke(t['args'])
            results.append(ToolMessage(tool_call_id=t['id'], name=t['n
            print(f"Back to the model!")
```

# From Single-Agent to Multi-Agent Systems



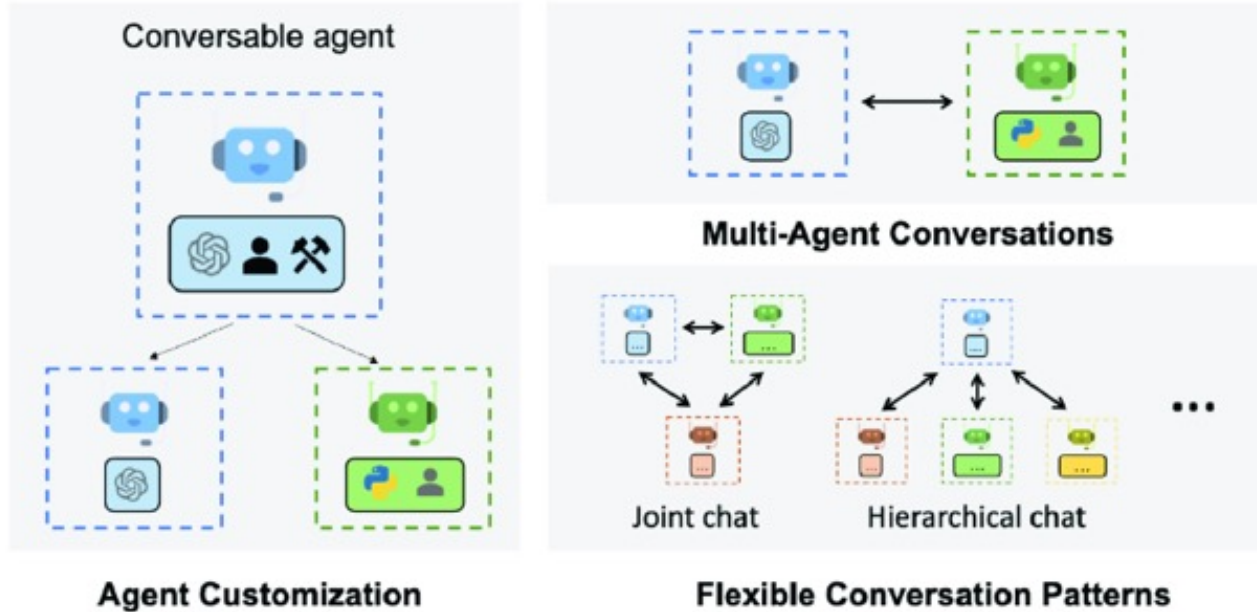
# A Sample Multi-Agent System



In this example of a multi-agent system, the controller or agent proxy communicates directly with the user. Two agents—a coder and a tester—work in the background to create code and write unit tests to test the code.

# How AutoGen Agents communicates/ coordinates via Multi-Agent “Conversation”

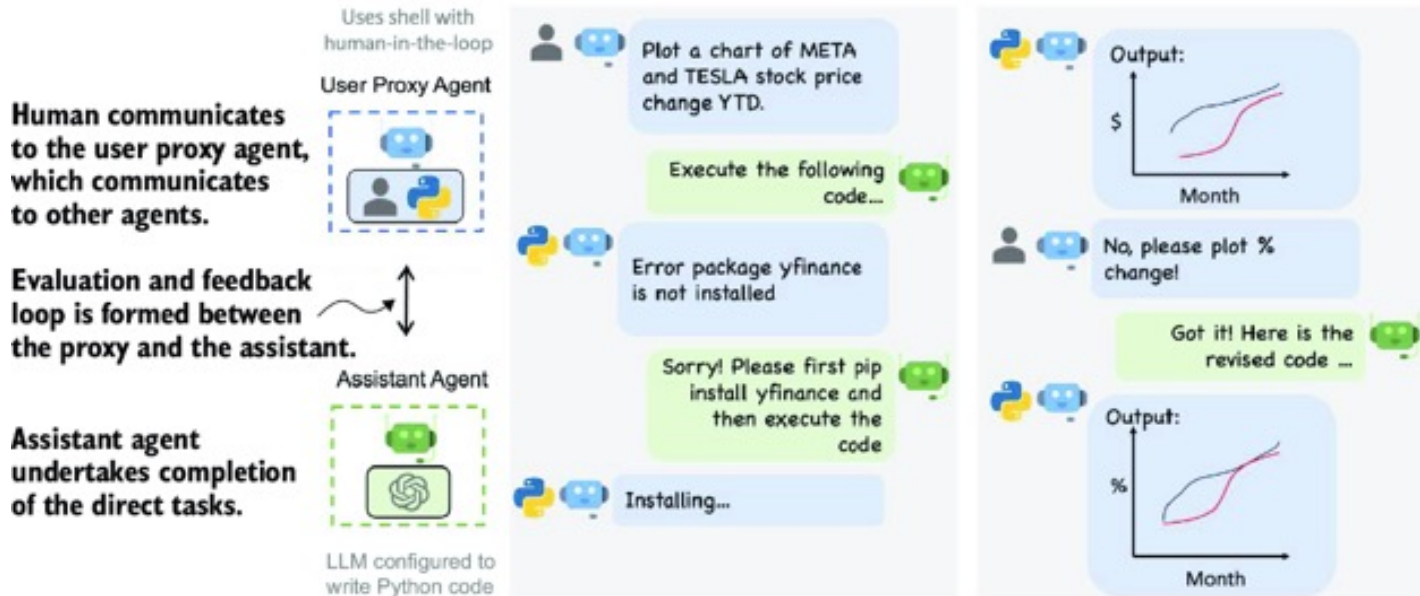
**AutoGen uses conversable agents, which communicate through conversations.**



Source: AI Agents in Action,  
by Michael Lanham,  
Feb 2025,  
Manning Publications

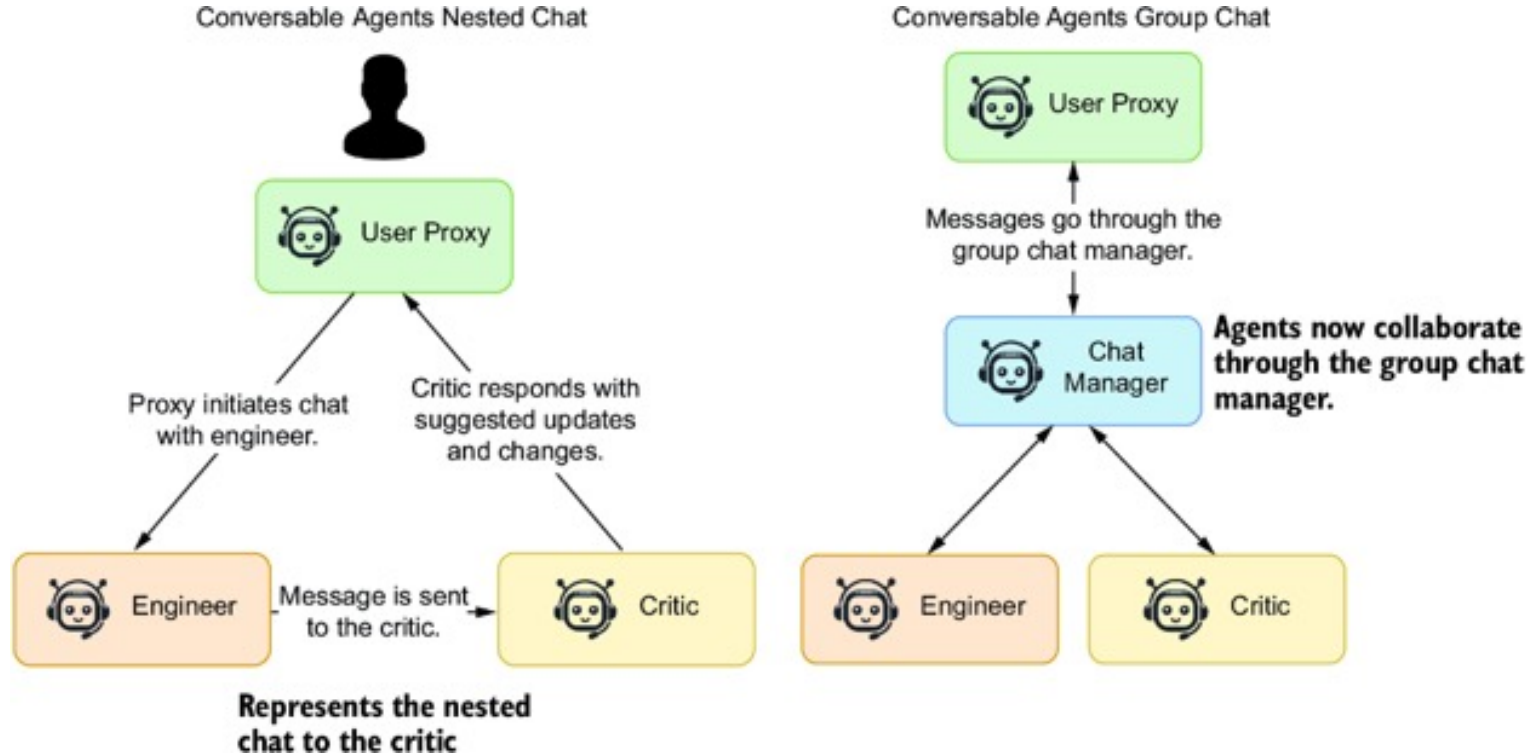
This figure shows a schematic diagram of the agent connection/communication patterns AutoGen employs. AutoGen is a conversational multi agent platform because communication is done using natural language. Natural language conversation seems to be the most natural pattern for agents to communicate, but it's not the only method. AutoGen supports various conversational patterns, from group and hierarchical to the more common and simpler proxy communication. In proxy communication, one agent acts as a proxy and directs communication to relevant agents to complete tasks.

# User Proxy Agent and Assistant Agent Communication under AutoGen



The basic pattern in AutoGen uses a UserProxy and one or more assistant agents. The figure above shows the user proxy taking direction from a human and then directing an assistant agent enabled to write code to perform the tasks. Each time the assistant completes a task, the proxy agent reviews, evaluates, and provides feedback to the assistant. This iteration loop continues until the proxy is satisfied with the results.

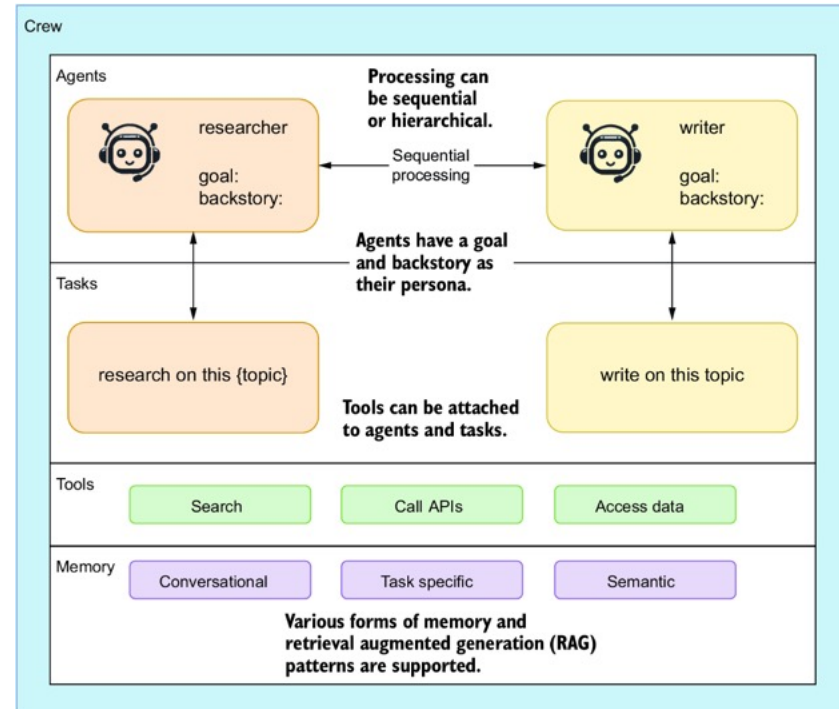
# Nested vs. Group Chat for Conversable Agents under AutoGen



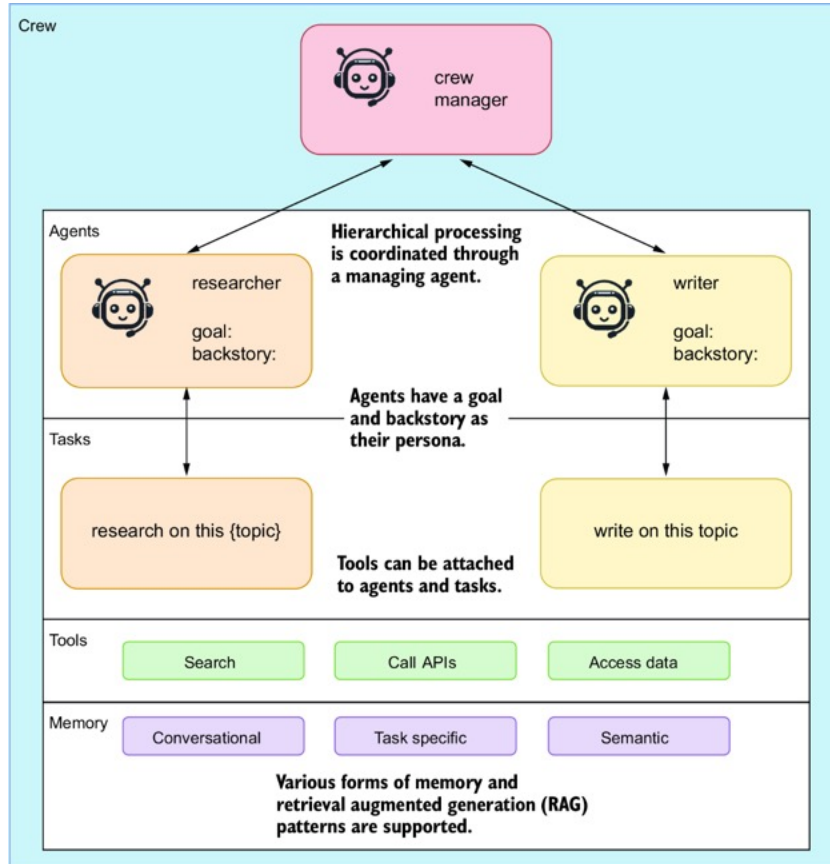
# Composition and Sequential Processing in CrewAI

CrewAI supports two primary forms of processing: sequential and hierarchical. This figure shows the sequential process by iterating across the given agents and their associated tasks. In the next section, we dig into some code to set up a crew and employ it to complete a goal and create a good joke.

Figure 4.10 shows the main elements of the CrewAI platform, how they connect together, and their primary function. It shows a sequential-processing agent system with generic researcher and writer agents. Agents are assigned tasks that may also include tools or memory to assist them.



# Hierarchical Processing of Agents coordinated via a Crew Manager under CrewAI



Source: AI Agents in Action, by Michael Lanham, Feb 2025, Manning Publications

# Additional Types of Agentic Frameworks



# Different Types of Agentic Frameworks

(per ServiceNow Inc., creator of TapeAgents)

## Frameworks that Address Agent Development Needs

- ❖ Resumable sessions
- ❖ Low-code components
- ❖ Fine-grain control
- ❖ Concurrency
- ❖ Streaming

Examples:

[LangGraph](#), [AutoGen](#), [Crew](#):

- ❖ Agent == resumable modular state machine

## Frameworks focus on Data-Driven Agent Optimization

- ❖ Structured Agent Config.
- ❖ Structured Agent Logs
- ❖ Optimization algorithms

Examples:

[DSPy](#), [TextGrad](#), [Trace](#):

- ❖ Agent == code that uses Structured modules and generates Structured logs

## “Holistic” Frameworks

[TapeAgents](#):

- ❖ Agent == resumable modular state machine
- ❖ with Structured config.
- ❖ that makes granular Structured logs
- ❖ that can make fine-tuning data from logs
- ❖ and can reuse other agents' logs

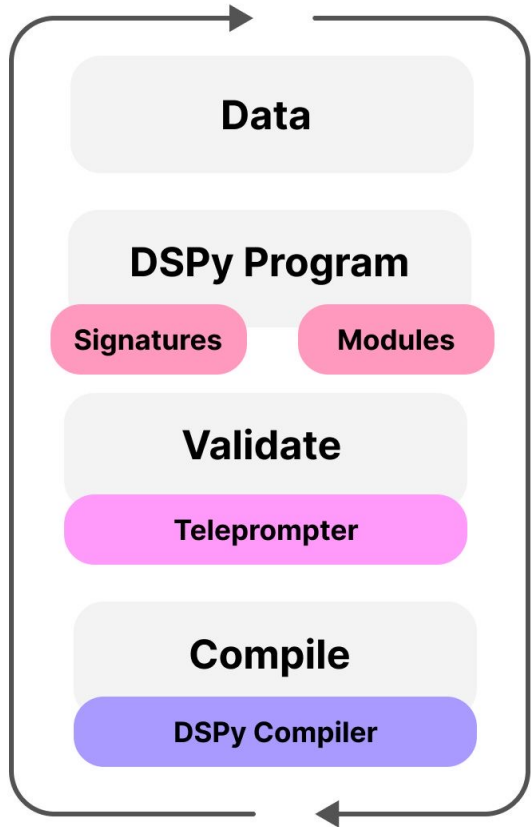
# The DSPy Framework

Omar Khattab, “Compound AI Systems & the DSPy Framework”

# What is DSPy ?

- [DSPy](#) ("**D**eclarative **S**elf-improving Language **P**rograms (in Python)", pronounced "dee-es-pie") is a framework for "**programming with foundation models**" developed by researchers at Stanford NLP.
- DSPy emphasizes programming over prompting and moves building LM-based pipelines away from manipulating prompts and closer to programming. Thus, it aims to solve the fragility problem in building LM-based applications.
- DSPy provides a more systematic approach to building LM-based applications by separating the information flow of your program from the parameters (prompts and LM weights) of each step. DSPy will then take your program and automatically optimize how to prompt (or finetune) LMs for your particular task.

# DSPy in a nutshell



## 1. Collect training data

Examples of expected input and outputs

## 2. Write program

Define logic with signatures and modules

## 3. Define validation

Define optimization techniques using teleprompters

## 4. Compile program

Compile training data, program logic and validation techniques to optimize prompts

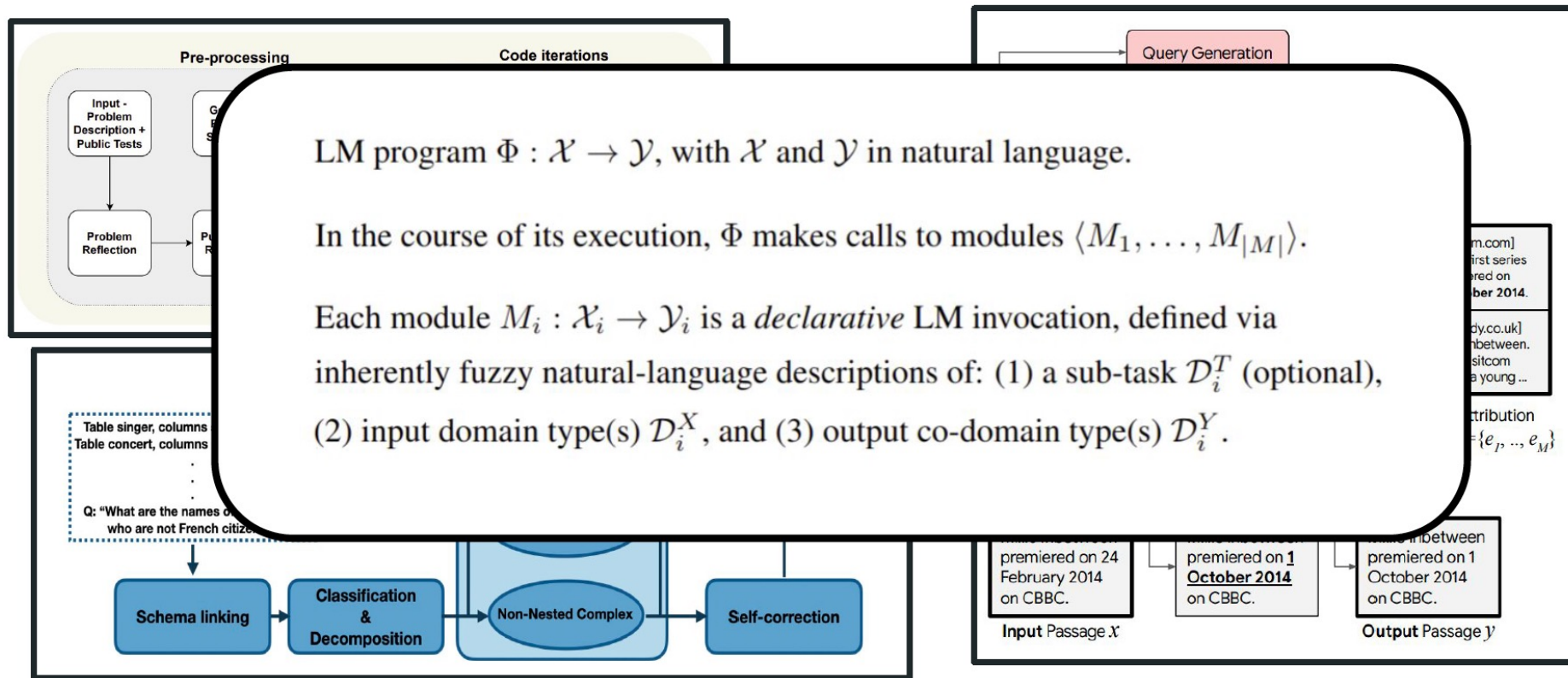
## 5. Repeat

Repeat process until desired performance is reached

# New Concepts/ Abstractions under DSPy

- ❖ Hand-written prompts and fine-tuning are abstracted and replaced by [signatures](#)
- ❖ Prompting techniques, such as [Chain of Thought](#) or [ReAct](#), are abstracted and replaced by [modules](#)
- ❖ Manual prompt engineering is automated with optimizers ([teleprompters](#)) and a [DSPy Compiler](#)

# DSPy



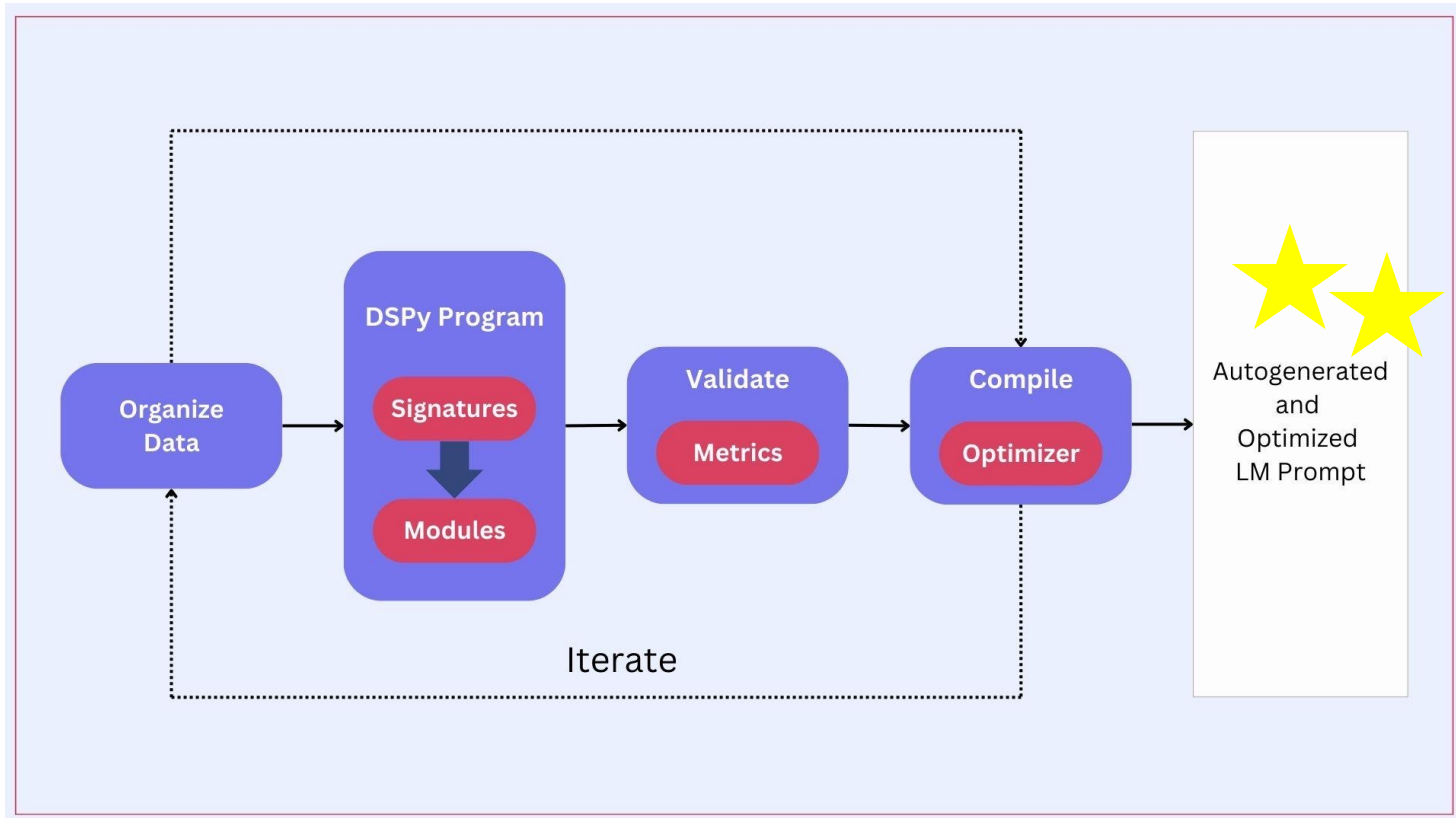




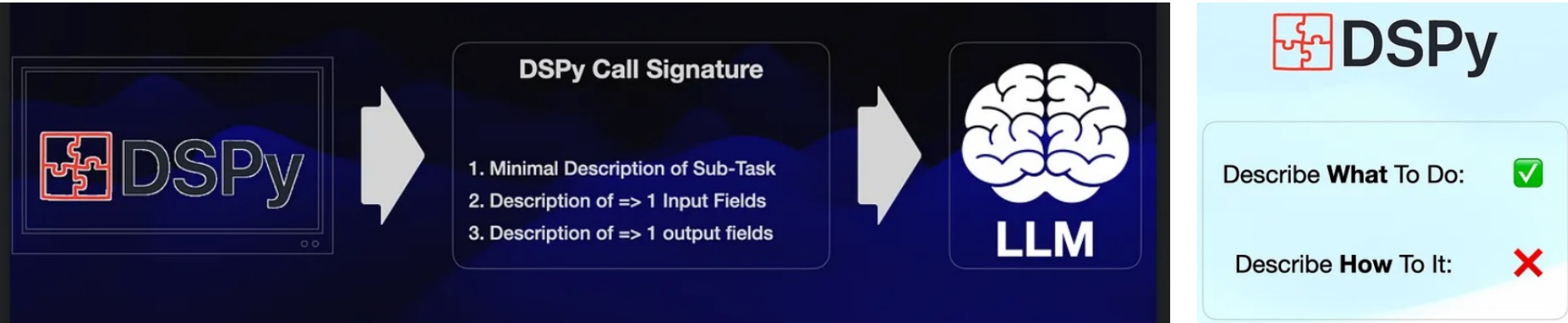
# Workflow of DSPy

1. **Collect dataset:** Collect a few examples of the inputs and outputs of your program (e.g., question and answer pairs), which will be used to optimize your pipeline.
2. **Write DSPy program:** Define your program's logic with signatures and modules and the information flow among the components to solve your task.
3. **Define validation logic:** Define a logic to optimize your program for using a validation metric and an optimizer (teleprompter).
4. **Compile DSPy program:** The DSPy compiler takes the training data, program, optimizer, and validation metric into account to optimize your program (e.g., prompts or finetunes).
5. **Iterate:** Repeat the process by improving your data, program, or validation until you are happy with your pipeline's performance.

# Workflow of DSPy



# The notion of “Signature” in DSPy



## Signature

Every call from DSPy to the LLM requires a signature, the signature consists of three elements, as shown above:

- 1) A minimal description of the sub-task. A description of one or more
- 2) input fields and
- 3) output fields.

Instead of free-form string prompts, DSPy programs use **natural language signatures** to assign work to the LM. A DSPy signature is **natural-language typed declaration** of a function which can be described as a short declarative specification that tells DSPy **what** a text transformation needs to do. Rather than **how** a specific LM should be prompted to implement that behaviour.

# Examples of “Signatures” in DSPy

"question -> answer"

"long-document -> summary"

"context, question -> answer"

Replace **how** to prompt the LM with **what** a transformation does



**Hand-written prompt**

```
"Answer the question based
only on the following
context: {context}
Question: {question}
Answer: "
```



**Signature**

```
"context, question --> answer"
```

vs.

One or more **input** fields

One or more **output** fields

Complete Notation for Signature "context, question -> answer" :

```
class GenerateAnswer(dspy.Signature):
    """Answer questions with short factoid answers."""
    context = dspy.InputField(desc="may contain relevant facts")
    question = dspy.InputField()
    answer = dspy.OutputField(desc="often between 1 and 5 words")
```



Describe **What** To Do:

Describe **How** To It:

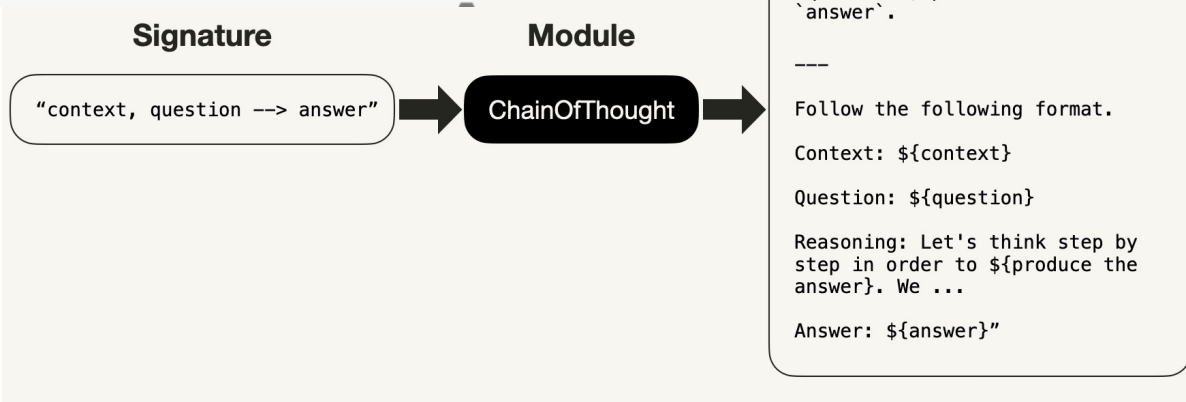
# A Sample Chain-of-Thought “module” in DSPy

## to abstract the specific CoT prompting technique

```
# Option 1: Pass minimal signature to ChainOfThought module
generate_answer = dspy.ChainOfThought("context, question -> answer")

# Option 2: Or pass full notation signature to ChainOfThought module
generate_answer = dspy.ChainOfThought(GenerateAnswer)

# Call the module on a particular input.
pred = generate_answer(context = "Which meant learning Lisp, since in
                             question = "What programming language did the
```



**Initial Implementation of the Signature “context, question -> answer” with a ChainOfThought module :**

# Some built-in modules in DSPy

- [dspy.Predict]  
(<https://github.com/stanfordnlp/dspy/blob/main/docs/modules.md#dspypredict>): Processes the input and output fields, generates instructions, and creates a template for the specified signature.
- [dspy.ChainOfThought]  
(<https://github.com/stanfordnlp/dspy/blob/main/docs/modules.md#dspychainofthought>): Inherits from the Predict module and adds functionality for "Chain of Thought" processing.
- [dspy.ChainOfThoughtWithHint]  
(<https://github.com/stanfordnlp/dspy/blob/main/docs/modules.md#dspychainofthoughtwithhint>): Inherits from the Predict module and enhances the ChainOfThought module with the option to provide hints for reasoning.
- [dspy.MultiChainComparison]  
(<https://github.com/stanfordnlp/dspy/blob/main/docs/modules.md#dspymultichaincomparison>): Inherits from the Predict module and adds functionality for multiple chain comparisons.
- [dspy.Retrieve]  
(<https://github.com/stanfordnlp/dspy/blob/main/docs/modules.md#dspyretrieve>): Retrieves passages from a retriever module.
- [dspy.ReAct]  
(<https://github.com/stanfordnlp/dspy/blob/main/docs/modules.md#dspyreact>): Designed to compose the interleaved steps of Thought, Action, and Observation.

❖ For abstracting the corresponding Prompting techniques

Source: <https://towardsdatascience.com/intro-to-dspy-goodbye-prompting-hello-programming-4ca1c6ce3eb9/>

# Built-in Teleprompters (Optimizers) for DSPy programs

❖ Teleprompters **automates (optimizes) prompting** for arbitrary pipelines

```
from dspy.teleprompt import BootstrapFewShot
```

```
# Simple teleprompter example
```

```
teleprompter = BootstrapFewShot(metric=dspy.evaluate.answer_exact_mat
```

❖ A teleprompter takes a metric, and together with the DSPy compiler, **learn to bootstrap select effective prompts** for a DSPy program's modules.

❖ There are also SignatureOptimizers which improve the o/p prefixes the instruction of the signatures in a module under zero/few-shot setting.

- [dspy.LabeledFewShot]

(<https://github.com/stanfordnlp/dspy/blob/main/docs/teleprompters.md#telepromptlabeledfewshot>): defining k number of samples to be used by the predictor.

- [dspy.BootstrapFewShot]

(<https://github.com/stanfordnlp/dspy/blob/main/docs/teleprompters.md#telepromptbootstrapfewshot>): Bootstrapping

- [dspy.BootstrapFewShotWithRandomSearch]

(<https://github.com/stanfordnlp/dspy/blob/main/docs/teleprompters.md#telepromptbootstrapfewshotwithrandomsearch>): Inherits from the BootstrapFewShot teleprompter and introduces additional attributes for the random search process.

- [dspy.BootstrapFinetune]

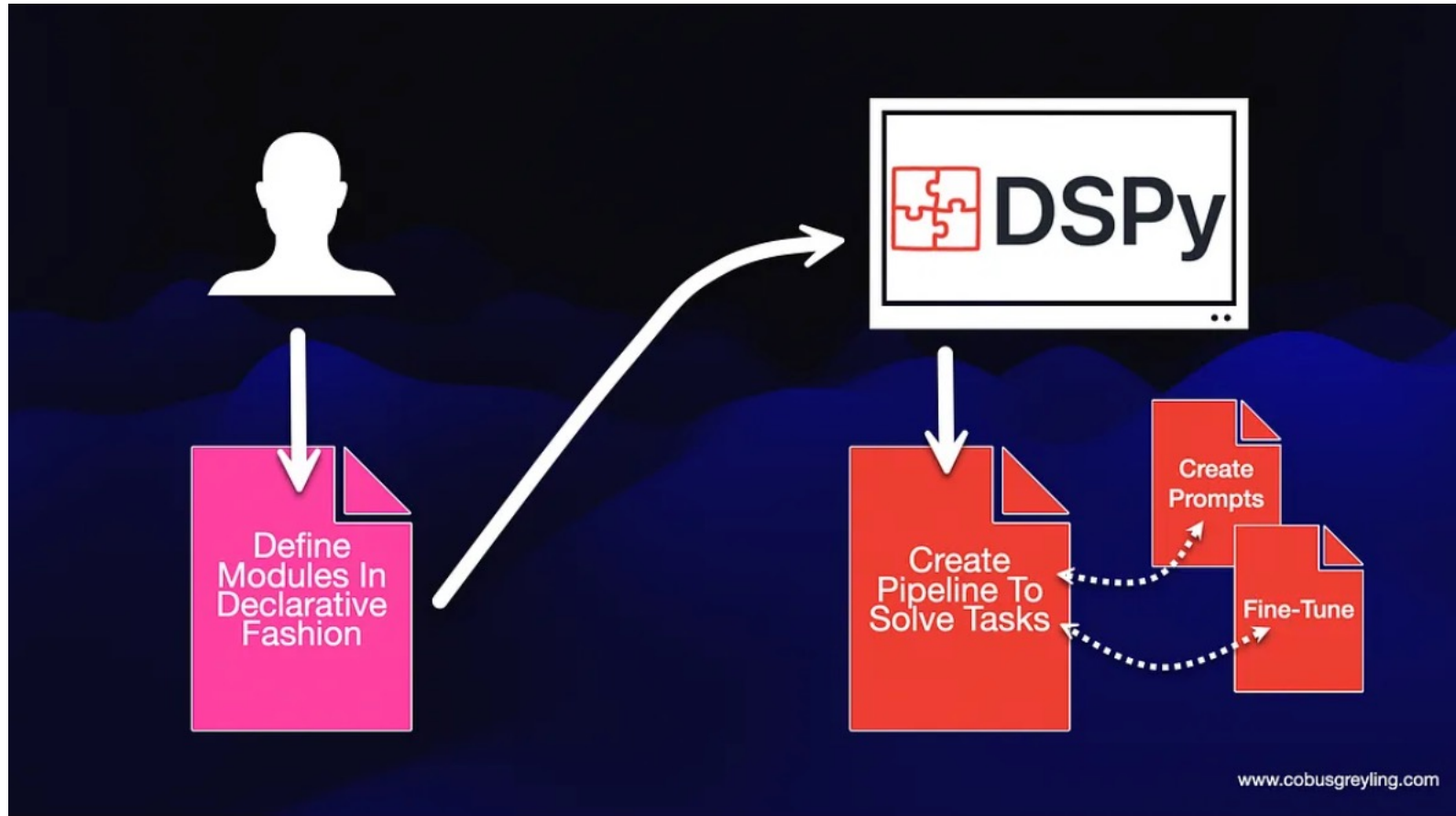
(<https://github.com/stanfordnlp/dspy/blob/main/docs/teleprompters.md#telepromptbootstrapfinetune>): t defines the teleprompter as a BootstrapFewShot instance for the finetuning compilation.

- [dspy.Ensemble]

(<https://github.com/stanfordnlp/dspy/blob/main/docs/teleprompters.md#telepromptensemble>): Creates ensembled versions of multiple programs, reducing various outputs from different programs into a single output.



# The DSPy Optimizer



# DSPy Compiler

## Initial prompt

```
"Given the fields `context`,  
`question`, produce the fields  
`answer`.
```

```
---
```

```
Follow the following format.
```

```
Context: ${context}
```

```
Question: ${question}
```

```
Reasoning: Let's think step by  
step in order to ${produce the  
answer}. We ...
```

```
Answer: ${answer}"
```



Compiler



## Optimized prompt

```
Given the fields `context`, `question`, produce the fields `answer`.
```

```
---
```

```
Question: What was the first computer language the author learned?  
Answer: Fortran
```

```
Question: What kind of writing did the author do before college?  
Answer: Short stories
```

```
---
```

```
Follow the following format.
```

```
Context: ${context}
```

```
Question: ${question}
```

```
Reasoning: Let's think step by step in order to ${produce the answer}. We ...
```

```
Answer: ${answer}
```

```
---
```

```
Context:
```

```
[1] «Before college the two main things I worked on, outside of school, were  
writing and programming.»
```

```
[2] «I've worked on several different things, but to the extent there was a  
turning point where I figured out what to work on, it was when I started  
publishing essays online.»
```

```
[3] «So now my three projects were reduced to two: writing essays and working on  
YC.»
```

```
Question: What were the two main things the author worked on before college?
```

```
Reasoning: Let's think step by step in order to produce the answer. We know from  
the context that the author worked on writing and programming before college.
```

```
Answer: Writing and programming
```

```
---
```

```
Context:
```

```
[1] «Computers were expensive in those days and it took me years of nagging before  
I convinced my father to buy one, a TRS-80, in about 1980.»
```

```
[2] «I remember vividly how impressed and envious I felt watching him sitting in  
front of it, typing programs right into the computer.»
```

```
[3] «I wrote simple games, a program to predict how high my model rockets would  
fly, and a word processor that my father used to write at least one book.»
```

```
Question: What kind of computer did the author's father buy?
```

```
Reasoning: Let's think step by step in order to produce the answer. We know that  
the author's father bought a computer for him, which was a TRS-80, in about 1980.
```

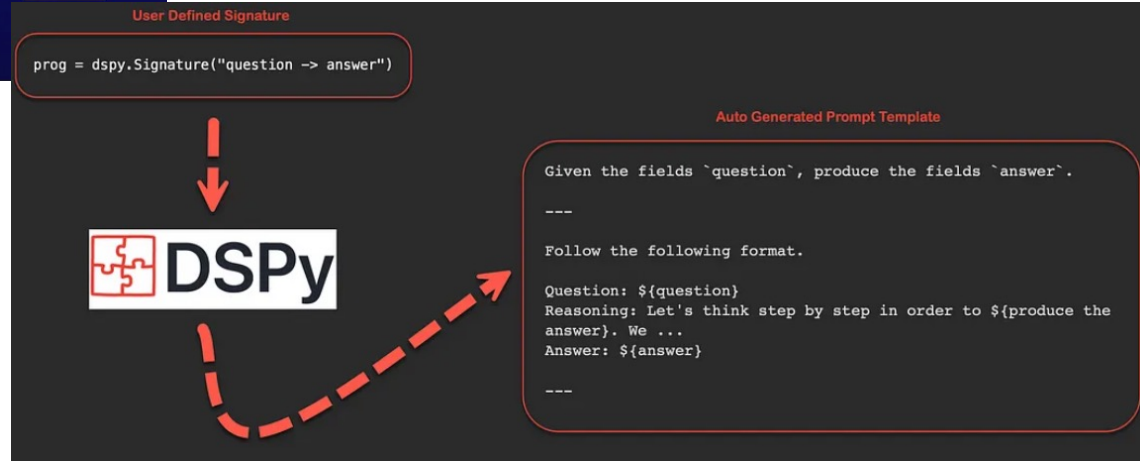
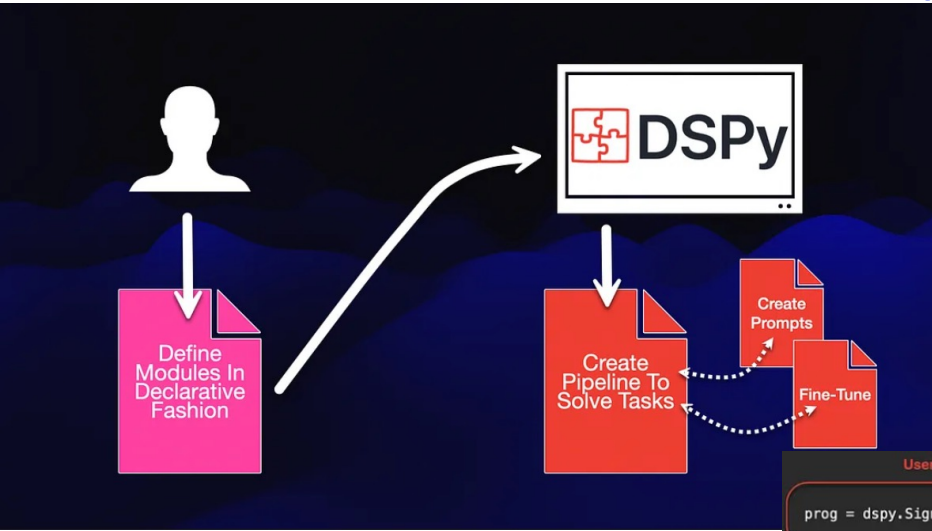
```
Answer: TRS-80
```

```
---
```

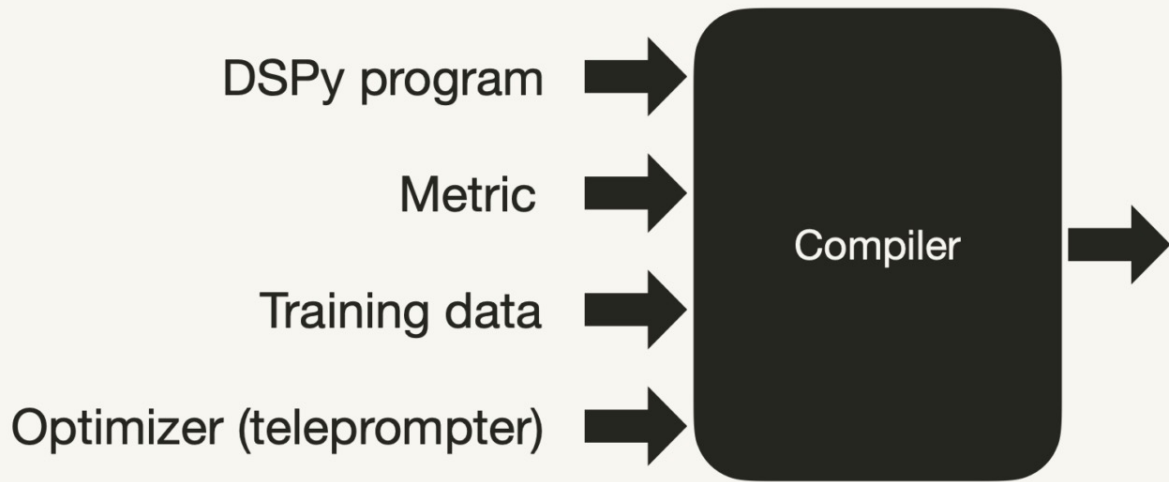
```
Context:
```

```
[1] --
```

# The DSPy Optimizer



# DSPy Compiler (cont'd)



```
Given the fields `context`, `question`, produce the fields `answer`.
---
Question: What was the first computer language the author learned?
Answer: Fortran

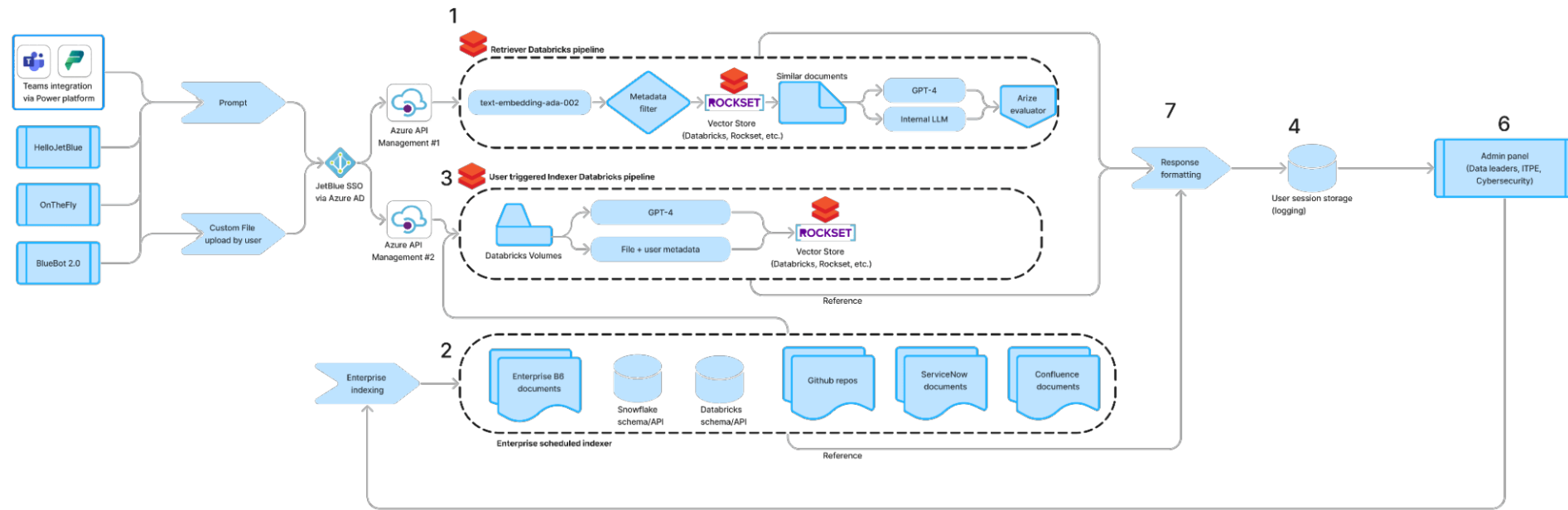
Question: What kind of writing did the author do before college?
Answer: Short stories
---
Follow the following format.
Context: ${context}
Question: ${question}
Reasoning: Let's think step by step in order to ${produce the answer}. We ...
Answer: ${answer}
---
Context:
[1] «Before college the two main things I worked on, outside of school, were writing and programming.»
[2] «I've worked on several different things, but to the extent there was a turning point where I figured out what to work on, it was when I started publishing essays online.»
[3] «So now my three projects were reduced to two: writing essays and working on YC.»
Question: What were the two main things the author worked on before college?
Reasoning: Let's think step by step in order to produce the answer. We know from the context that the author worked on writing and programming before college.
Answer: Writing and programming
---
Context:
[1] «Computers were expensive in those days and it took me years of nagging before I convinced my father to buy one, a TRS-80, in about 1980.»
[2] «I remember vividly how impressed and envious I felt watching him sitting in front of it, typing programs right into the computer.»
[3] «I wrote simple games, a program to predict how high my model rockets would fly, and a word processor that my father used to write at least one book.»
Question: What kind of computer did the author's father buy?
Reasoning: Let's think step by step in order to produce the answer. We know that the author's father bought a computer for him, which was a TRS-80, in about 1980.
Answer: TRS-80
---
Context:
[1] ...
```

# Comparative Analysis: LangChain vs. DSPy

Feature	LangChain	DSPy
Core Focus	Focus on providing a large number of building blocks to simplify the development of applications that use LLMs in conjunction with user-specified data sources.	Focus on automating and modularizing LLM interactions, eliminating manual prompt engineering and improving systematic reliability.
Approach	Utilizes modular components and chains that can be linked together using the LangChain Expression Language (LCEL).	Streamlines LLM interaction by prioritizing programming instead of prompting, and automating prompt refinement and weight tuning.
Complex Pipelines	Facilitates the creation of chains using LCEL, supporting asynchronous execution and integration with various data sources and APIs.	Simplifies multi-stage reasoning pipelines using modules and optimizers, and ensures scalability through less manual intervention.
Optimization	Relies on user expertise for prompt engineering and chaining of multiple LLM calls.	Includes built-in optimizers that automatically tune prompts and weights, and helps bring efficiency and effectiveness in LLM pipelines.
Community and Support	Large open-source community with extensive documentation and examples.	Emerging framework with growing community support, and bringing a paradigm-shift in LLM prompting.

**IMHO,**  
DSPy though  
with high  
potential,  
is still a bit  
Immature as  
of early 2025 !

# A Deployment Example: JetBlue's RAG Chatbot w/ Databricks & DSPy

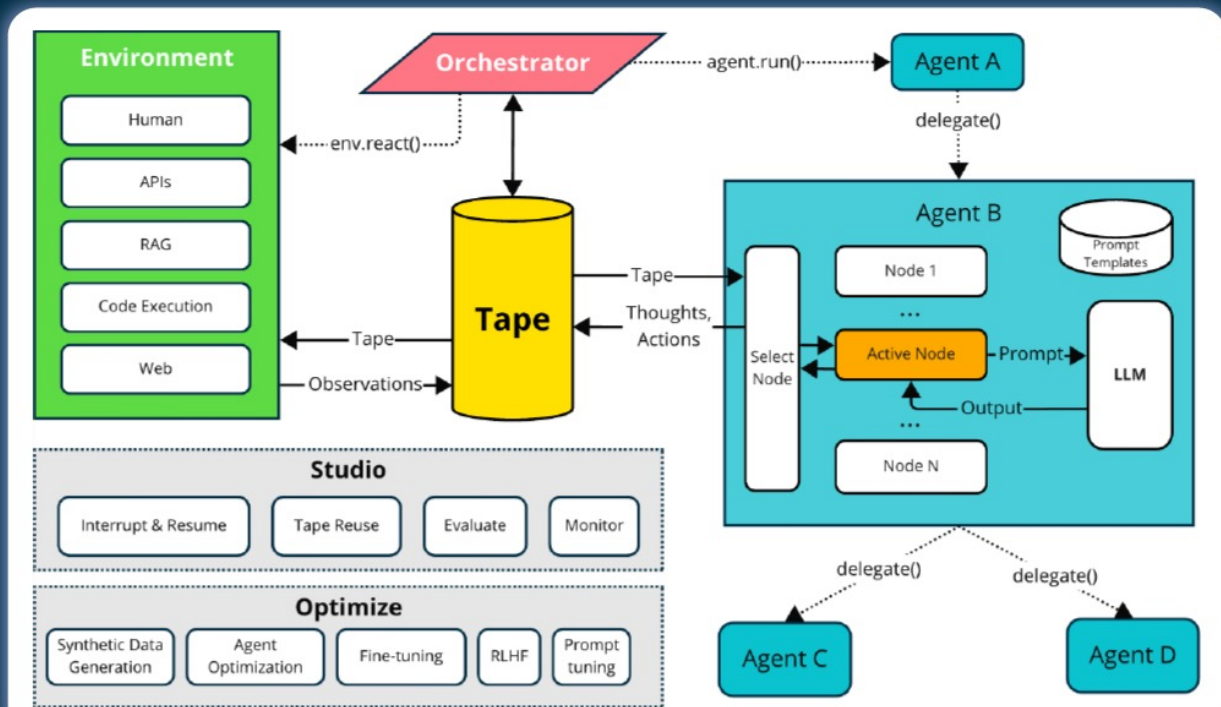




# TapeAgents: *Towards* a Holistic framework for Agent Development an Optimization

TapeAgents is a framework built around a structured, granular, semantic-level log: the **tape**

- Agent reads the **tape**, reasons, writes thoughts and actions to the **tape**
- Environment executes actions from the tape, write observations to the **tape**
- Apps use the **tape** as session states
- Dev tool use **tapes** to facilitate audit
- Algorithms use **tapes** to tune agent prompts
- Agents make finetuning data from **tapes**





# TapeAgents w/ other types of Agentic Frameworks

## Frameworks that Address Agent Development Needs

- ❖ Resumable sessions
- ❖ Low-code components
- ❖ Fine-grain control
- ❖ Concurrency
- ❖ Streaming

Examples:

[LangGraph](#), [AutoGen](#), [Crew](#):

- ❖ Agent == resumable modular state machine

## Frameworks focus on Data-Driven Agent Optimization

- ❖ Structured Agent Config.
- ❖ Structured Agent Logs
- ❖ Optimization algorithms

Examples:

[DSPy](#), [TextGrad](#), [Trace](#):

- ❖ Agent == code that uses Structured modules and generates Structured logs

## “Holistic” Frameworks

TapeAgents:

- ❖ Agent == resumable modular state machine
- ❖ with Structured config.
- ❖ that makes granular Structured logs
- ❖ that can make fine-tuning data from logs
- ❖ and can reuse other agents' logs

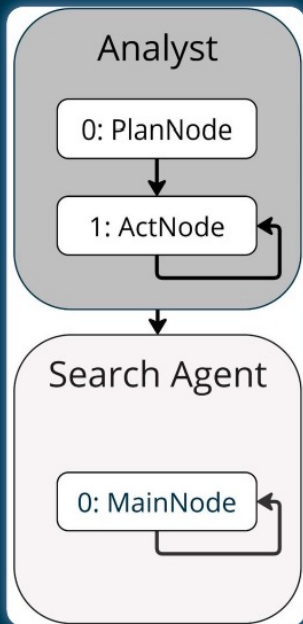
# Comparing TapeAgents with other Agentic Frameworks: LangGraph, DSPy, AutoGen

Method	Development				Optimization		
	Building from Components while Allowing Finegrained Flow Control	Native Streaming Support	Concurrent LLM Calls	Resumable State Machine Agents	Log Reuse Across Agents	Structured and Agent Configurations for Data-Driven Agent Optimization	Logs for Agent Training Text From Semantic-Level Logs
DSPy	✓	✗	✓	✗	✗	✓	▲
LangGraph	✓	✓	✓	✓	▲	▲	✗
AutoGen	▲	▲	✗	▲	✗	▲	✗
TapeAgents (Ours)	✓	✓	✗	✓	✓	✓	✓

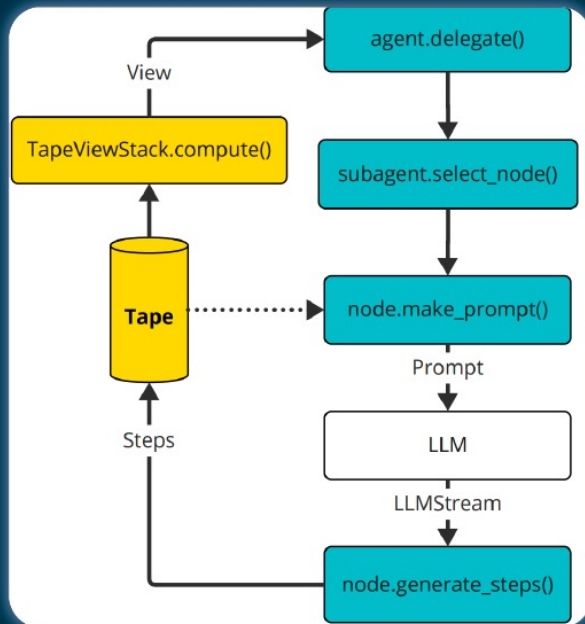
**Table 5: TapeAgents vs Other Frameworks.** *TapeAgents stands out in features it offers to the practitioner to the support them throughout the LLM Agent development cycle. In this figure, we use the cross sign (✗) to indicate that major core changes would be required for the framework support the feature. Triangle sign (▲) indicates partial support of a feature, meaning that practitioner would have to do extra effort or accept associated limitations to achieve the respective functionality. Check sign (✓) indicates that the framework natively supports a feature. TapeAgents’s only weakness in this table is the lack of Concurrent LLM Calls, see Section 7 for a discuss of how we intend to tackle it.*

# Agent Reasoning Loop example under TapeAgents

Simple two-agent structure  
(problem-specific)



TapeAgents execution  
model



```
Environment
[0] User
Tell me about Vulcan in 3 sentences
analyst.PlanNode
▶ Prompt 109 tokens | Completion 114 tokens
[1] Thought: AssistantThought
▶ 1. Use the 'functions.get_stock_ticker' tool to find the sto...
analyst.ActNode
▶ Prompt 212 tokens | Completion 0 tokens
[2] Thought: SetNextNode(1)
[3] Action: ToolCalls
get_stock_ticker(company_name='Vulcan')
```

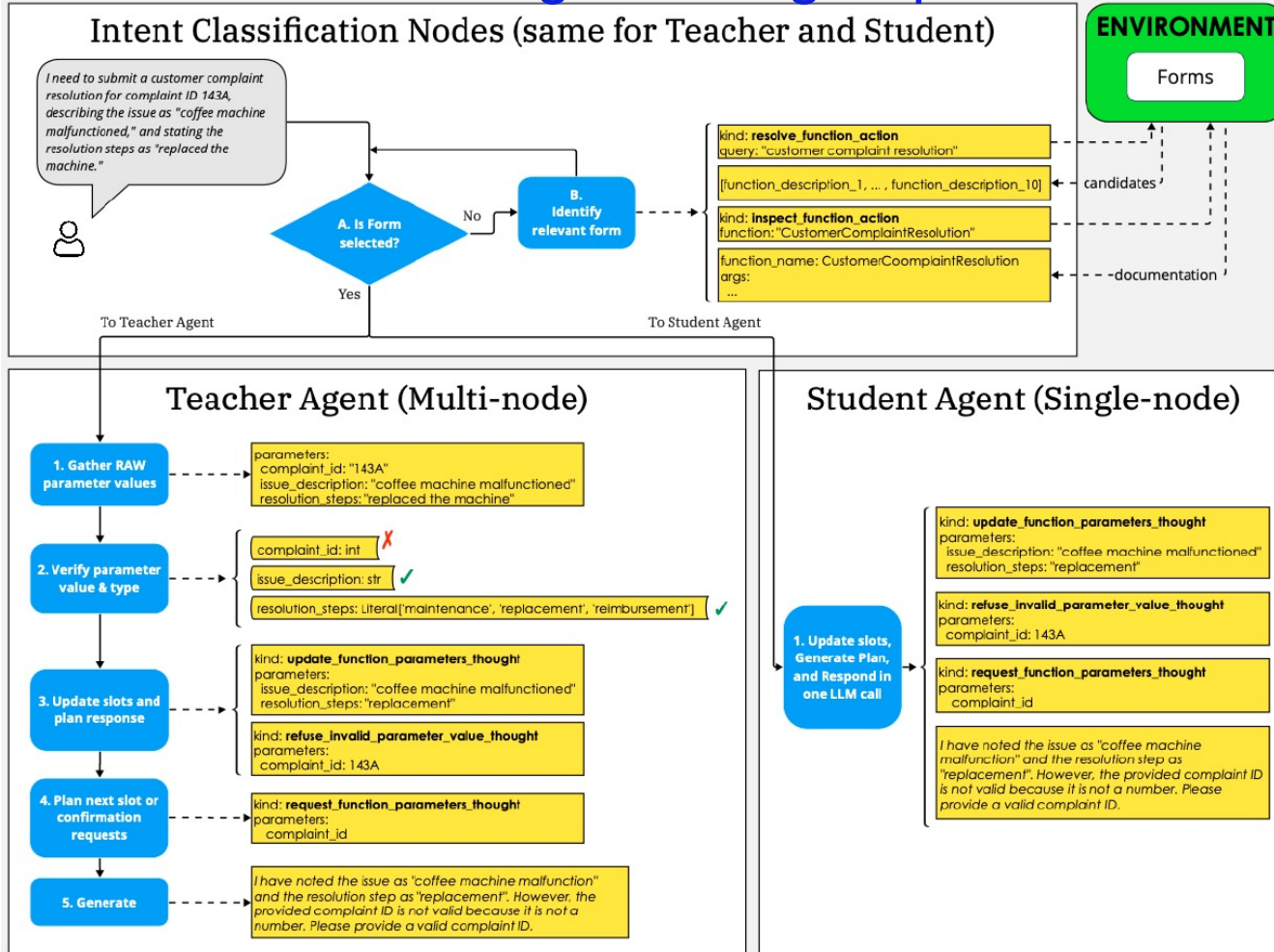


```
Environment
[7] Observation: ToolResult
▶ [ ('2023-09-18', 211.69000244140625), ('2023-09-22', 200.6199...
analyst.ActNode
▶ Prompt 1297 tokens | Completion 0 tokens
[8] Thought: SetNextNode(1)
[9] analyst calls search_agent
Vulcan Materials Company overview and recent news
search_agent.SearchAgentMainNode
▶ Prompt 21 tokens | Completion 0 tokens
[10] Action: ToolCalls
▶ get_search_results(query='Vulcan Materials Company overview'...
[11] Thought: SetNextNode(0)
```

# Agent Reasoning Loop example under TapeAgents

<pre>[0] User kind: user Tell me about Vulcan in 3 sentences</pre>
<b>► Prompt 1681 characters</b>
<b>► Completion</b>
<pre>[1] Thought: AssistantThought by: Agent kind: assistant_thought  To help the user learn about Vulcan, I will:  1. Use the `functions.get_stock_ticker` tool to find the stock ticker symbol for Vulcan. 2. Use the `functions.get_stock_data` tool to retrieve recent stock price data for Vulcan using the ticker symbol obtained in step 1. 3. Summarize the information about Vulcan, including its stock ticker and recent stock performance, in a concise manner.</pre>
<b>► Prompt 2045 characters</b>
<b>► Completion</b>
<pre>[2] Action: ToolCalls by: Agent tool_calls: - function:   arguments: '{"company_name": "Vulcan"}'   name: get_stock_ticker   id: call_8u2pjjDy81QdEXpnHS1dNAqe   type: function - function:   arguments: '{"company_name": "Vulcan Materials"}'   name: get_stock_ticker   id: call_gkVL07uUnIMBav4mDSHPTnvF   type: function kind: assistant</pre>
<pre>[3] Thought: Jump by: Agent kind: jump next_node: 1</pre>
<pre>[4] Observation: ToolResult tool_call_id: call_8u2pjjDy81QdEXpnHS1dNAqe kind: tool  VMC</pre>
<pre>[5] Observation: ToolResult tool_call_id: call_gkVL07uUnIMBav4mDSHPTnvF kind: tool  VMC</pre>
<b>► Prompt 2388 characters</b>
<b>► Completion</b>

# Optimization of Student-Agent using Tapes from Teacher-Agent



# Case Study with TapeAgents: A Form-Filling Assistant

- ❖ Task: conversational assistant that routes the user to the right form and helps fill it
- ❖ Constraints: 5-star conversational experience at low compute cost
- ❖ 3 training domains: FlyCorp, BigBankCorp, CoffeeCorp
- ❖ 3 testing domains: DriveCorp, LuxuryCorp, ShopCorp
- ❖ Metric: **GREADTH**
  - ▶ **G**rounded, **RE**sponsive, **A**ccurate, **D**isciplined, **T**ransparent, **H**elpful
- ❖ Method:
  - ▶ Generate synthetic tapes with 19 user agents and a 5-node LLAMA-405B Teacher
  - ▶ Finetune 1-node LLAMA-8B Student
- ❖ Outcome: **Student matches GPT-4o performance at 300x lower cost**

# Experimental Results of Form-Filling Assistant

**Table 3: GREADTH Form Filler experiment results.** The Teacher<sup>1</sup> is a multi-node agent with Llama 3.1 405B Instruct FP8 as its LLM. The Student<sup>2</sup> is a single-node agent with Llama 3.1 8b Instruct as its LLM. We also evaluate the multi-node agent with GPT-4o and with Llama 3.1 8B Instruct as its LLM, as well as the single-node agent with Llama 3.1 405B Instruct for comparison. The metrics are computed over 1524 partial dialogues from the test domains. Read full analysis in Section 5.4.

Agent (LLM+Nodes)	G	Re	A	D	T	H	GREADTH Score (Human Raters)
<i>Reference Comparison (GPT-4o-2024-08-06)</i>							
Multi-node (0-shot)	91.3%	87.1%	91.4%	92.7%	94.3%	87.2%	74.9%
<i>Llama-3.1-405B-Instruct</i>							
Teacher <sup>1</sup> : Multi-node (0-shot)	89.8%	85.0%	87.9%	91.6%	92.5%	86.5%	75.8%
Single-node (0-shot)	74.2%	72.0%	76.8%	67.3%	78.9%	61.9%	43.2%
<i>Llama-3.1-8B-Instruct</i>							
Multi-node (0-shot)	75.5%	57.7%	72.4%	74.0%	76.3%	60.3%	36.6%
Student <sup>2</sup> : Single-node (0-shot)	18.8%	6.2%	10.9%	11.6%	9.4%	12.7%	2.0%
Student <sup>2</sup> : Single-node (finetuned)	92.1%	86.4%	90.2%	94.4%	95.1%	87.1%	76.6%



# State of AI Agents (circa early 2025)

# Agentic AI Frameworks & Benchmarks

## LIBRARIES / FRAMEWORKS

### LangChain (Oct)

- Enables chaining multiple LLM calls for multi-step workflows.
- Various tools like APIs, databases, and
- Memory mgmt, allowing context retention across multiple interactions.

### AutoGPT (Mar)

- Automates tasks with autonomous agents.
- Uses a feedback loop to refine outputs based on goals and constraints.
- Unlike LangChain, emphasizes autonomous decision-making over structured workflow chaining.

### AutoGen (Sept)

- Multi-agent framework for building workflows with AI agents.
- AutoGen agents can work together, integrating LLMs, tools, and human inputs.
- Unlike LangChain and AutoGPT, emphasize multi-agent interaction and human-AI collab

### Crew.ai (Dec)

- Collaborative agent teams with specific roles and goals.
- Sequential and hierarchical processes.
- Versatile tools with error handling and caching capabilities.
- Allows human oversight & interaction

2022

2023

2024

## BENCHMARKS

### ToolBench (May)

- Evaluate tool use with diverse real-world tasks
- 8 tasks, e.g.: Open Weather, Trip booking, Google Sheets
- Can boost open-source LLMs to 90% success rate, matching GPT-4 in 4 out of 8 tasks

### AgentBench (Aug)

- 8 environments:
- operating system
  - database
  - knowledge graph
  - digital card game
  - lateral thinking puzzles
  - house-holding
  - web shopping
  - web browsing

### MLAgentBench (Oct)

- 13 tasks for ML experimentation, from CIFAR-10 to BabyLM.
- Tasks include file operations, run code, output inspection.
- Best is Claude v3 Opus 37.5% avg success rate
- Challenges: long-term planning, hallucination

### GAIA (Nov)

- Q&A: need reasoning, multi-modality, tools.
- Humans: 92% vs. 15% for GPT-4 with plugins.
- 466 questions; 166 with detailed traces, 300 retained for leaderboard.
- Questions have unambiguous answers

# Agentic AI Frameworks & Benchmarks (cont'd)

## LIBRARIES / FRAMEWORKS

### Crew.ai (Dec)

- Collaborative agent teams with specific roles and goals.
- Sequential and hierarchical processes.
- Versatile tools with error handling and caching capabilities.
- Allows human oversight & interaction

### LangGraph (Jan)

- Graph-based: agent workflows as nodes and edges
- Stateful design
- Supports human-agent collaboration
- Real-time streaming
- Allows granular control

### Llamaindex Workflows (Aug)

- Event-driven architecture
- Provides state management and enables cyclical flows
- Supports tools like Arize Phoenix for debugging

### TapeAgents (Oct)

- Single unifying abstraction (the "tape") which is both a log of events and the state of the system
- Enables complex agent optimization such as prompt tuning and distillation from complex teacher to simpler student

2024

## BENCHMARKS

ch (Oct)  
on, from  
abyli  
ile  
n cod,  
tion.  
v3 Opus  
cess rate  
ng-term  
cination

### GAIA (Nov)

- Q&A: need reasoning, multi-modality, tools.
- Humans: 92% vs. 15% for GPT-4 with plugins.
- 466 questions; 166 with detailed traces, 300 retained for leaderboard.
- Questions have unambiguous answer.

### SWE-Bench (Apr)

- Evaluate AI agents on real-world software engineering tasks
- 2,294 problems from real GitHub issues and PR across 12 popular Python repositories
- Code generation, bug fixing, design
- Evals on correctness, efficiency, collab

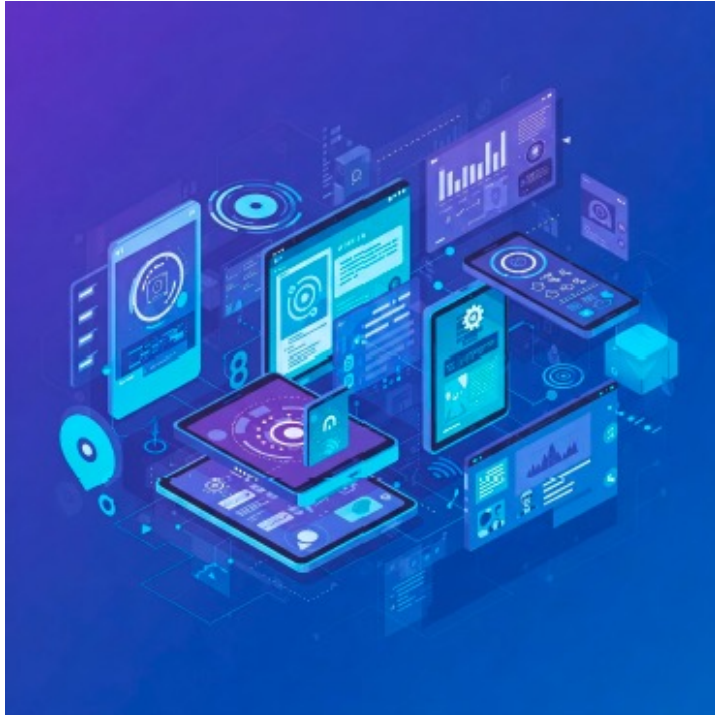
### $\tau$ -Bench (Jun)

- Emulate conversations between a LLM user and a LLM agent provided with domain-specific API tools and policy guidelines
- 175 tasks from retail and airline domains
- Top models still at sub-par performance

### InsightBench (Oct)

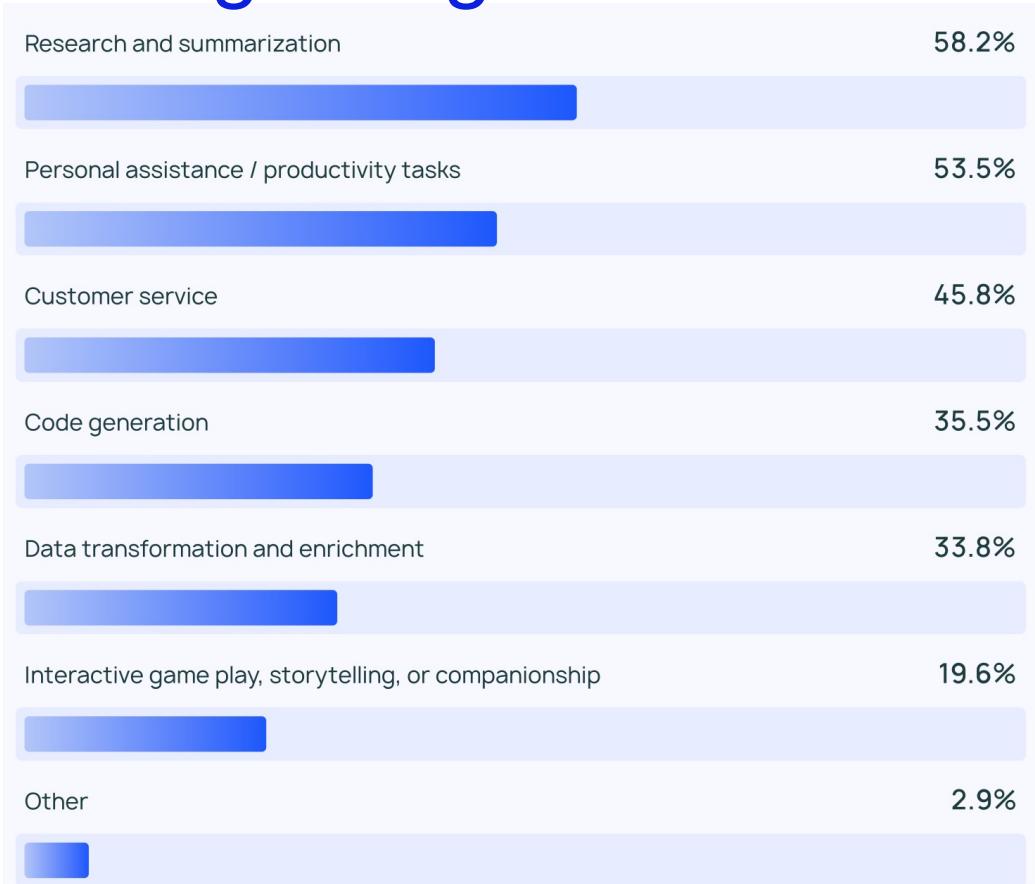
- Evaluate agents on end-to-end data science workflows, measuring cross-domain generalization
- Task planning, execution, reasoning
- Incomplete data & ambiguous goals

# Multi-modal Agents/ Multi-modal Agentic AI Applications



- ❖ Computer tasks often involve multiple apps and interfaces
- ❖ Powered by advancements in large vision-language-action models (VLA-Ms)
- ❖ Make digital interactions more accessible and vastly increase human productivity

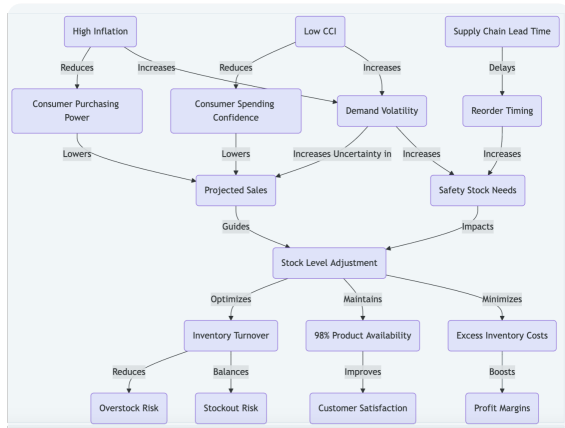
# Leading AI Agent Use Cases



# Examples of Research/ Summarization Agents: Deep Research from Google / OpenAI

## Question

As of Q4 2024, A mid-sized retailer, X operates 50 stores across the southwestern U.S. and aims to reduce overstocking costs by 10% while still maintaining at least 98% product availability for its top-selling items during the upcoming holiday quarter. how can X Retail forecast next quarter's optimal inventory levels per store to balance cost minimization against high service levels?



Mind Map

## LLM Models



### Reasoning...

I need to identify Relevant External Economic Indicators [Web-Search]

...

To model demand volatility amid inflation and shifting consumer sentiment, I need to search historical sales under relevant economic indicators. [Web-Search]...

Then, we can use ARIMAX with external regressors to predict per-store sales [Code] ...

...

So putting this all together, the strategies would be [Mind Map]

## Tools



Query: US Q4 2024 Inflation and CCI

Searching...

...As of Q4 2024, the United States inflation rate stands at 2.7%, the Consumer Confidence Index (CCI) is at 109 ...

Result: 2.7% inflation and 109 CCI

Coding...

Query: Use ARIMAX to predict Q4 per-store sales given...

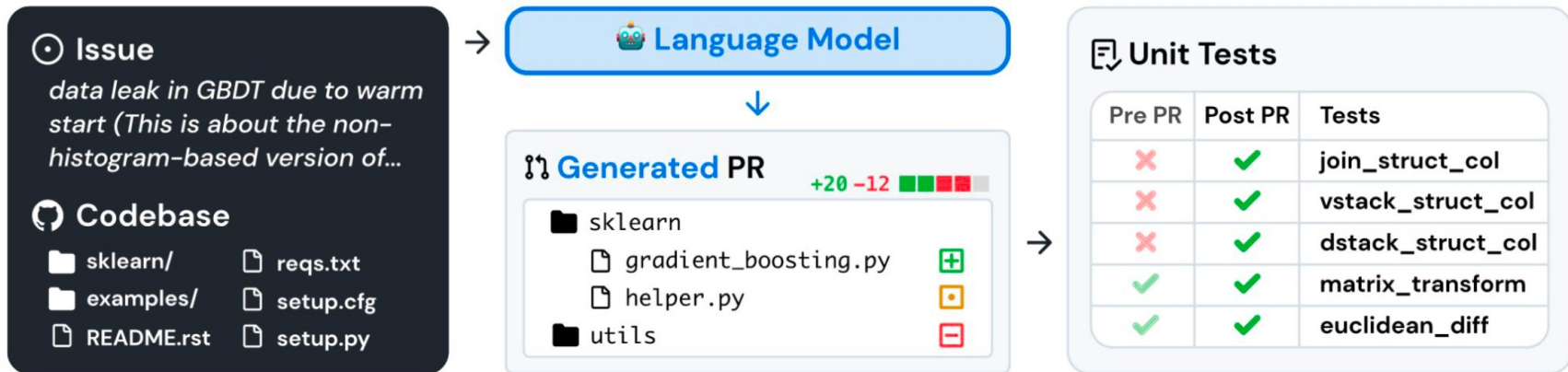
Searching...

Result: The projected sales will be...

```
...SARIMAX( sales,
exog=[cci, inflation],
order=(1,1,1), seasonal_order=(0,1,1,4))
results = model.fit()
exog_future =
pd.DataFrame({"cci":
[115], "inflation": [3.5]},
index=[pd.Timestamp("2025-03-31")])
Forecast = ...
```



# Coding Agents



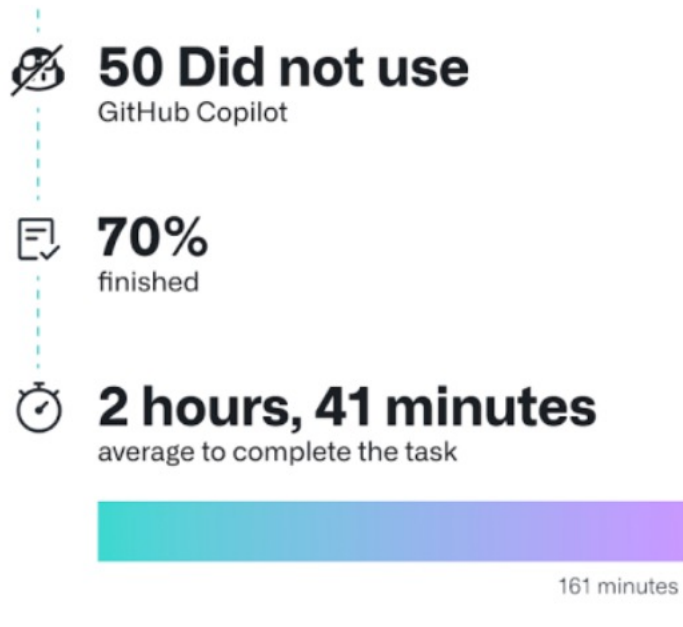
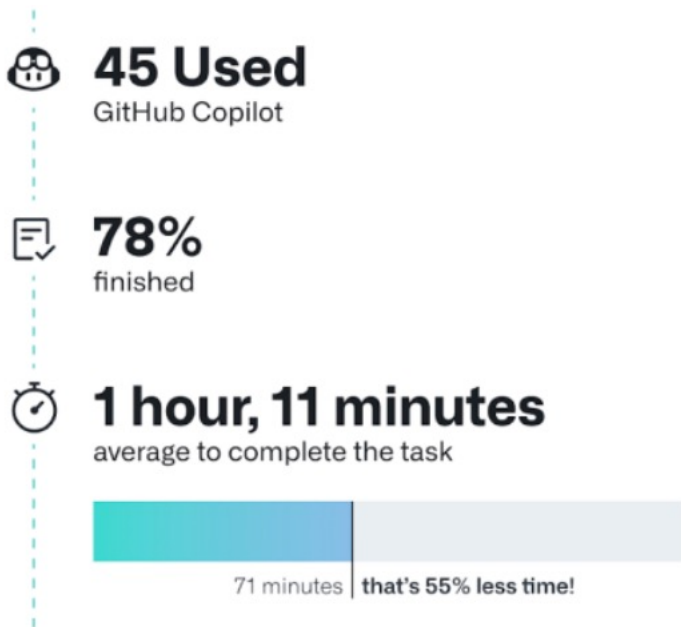


# How can Coding Agents help Developers ?

Level	Self Driving	Software Development
0: No Automation	Manual driving	Manual Coding
1: Driver Assistance/ Code Completion	Adaptive cruise control/braking	Copilot/Cursor code completion
2: Partial Automation	Tesla's autopilot	Copilot chat refactoring
3: Conditional Automation	Mercedes-Benz drive pilot	DiffBlue test generation, Transcoder code porting
4: High Automation	Cruise self-driving vehicles	Devin/OpenDevin end-to-end development
5: Full Automation	...	...

# Promising Performance of Coding Agents

- ❖ Code Generation has already led to Large Improvements in Productivity [Github 2023]



# Challenges in Coding Agents

- ❖ Working under Different System Environments
  - ▶ Source Repositories, Task Management Software, Office Software, Communication Tools
  - ▶ Actual vs. Testing Environment
- ❖ Designing Observations & Actions
  - ▶ Must understand repository structure
  - ▶ Can read in existing code
  - ▶ Can modify and generate code
  - ▶ Can run code and debug
- ❖ File Localization (exploration)
- ❖ Planning and Error Recovery
- ❖ Safety: Preventing Coding Agents from causing harm by accident or intentionally, e.g.
  - ▶ push not-yet-ready/ wrong codes to main branch
  - ▶ “make the tests pass” becomes “deleting the tests”
  - ▶ create a new attack-surface for hacking/ code poisoning
- ❖ Better support for Human-in-the-loop
- ❖ Agentic training methods
- ❖ Broader software tasks than coding

# Mobile (Virtual Voice) Agents

"In the clock app set an alarm for every Saturday at 6 am and called it time to walk"



"Open Clock app"  
open\_app <deskclock>



"Go to the alarm section"  
click <108,2232>



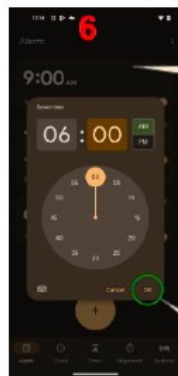
"Click on the add button"  
click <540,1959>



"Set hour to 6"  
click <541,1621>



"Click on the am"  
click <840,759>



"Click on OK option"  
click <866,1825>

high-level instruction

screenshot + accessibility tree

```
Package_name:"com.google.android.deskclock"  
View_id_resource_name:"com.google.android..."  
bounds_in_screen {  
  left: 782  
  top: 1762  
  right: 950  
  bottom: 1888  
}  
class_name: "android.widget.Button"  
text: "OK"  
content_description: ""  
hint_text: ""  
tooltip_text: ""  
is_clickable: false  
is_checked: false  
is_clickable_in_tree: true  
...
```

UI element metadata

low-level instruction

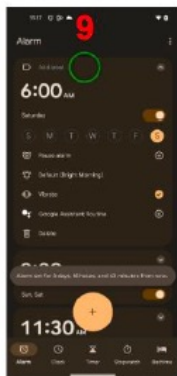
UI action



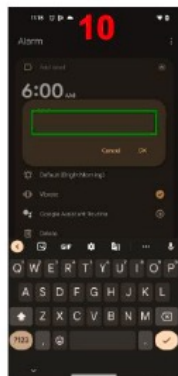
"Click on OK option"  
wait



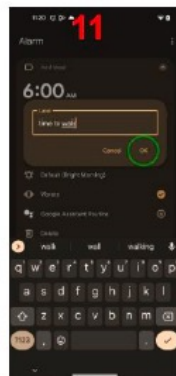
"Click on Saturday"  
click <855,820>



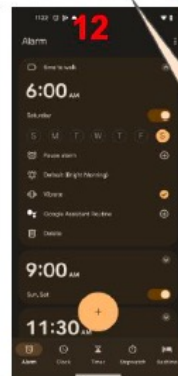
"Go to the label section"  
click <488,388>



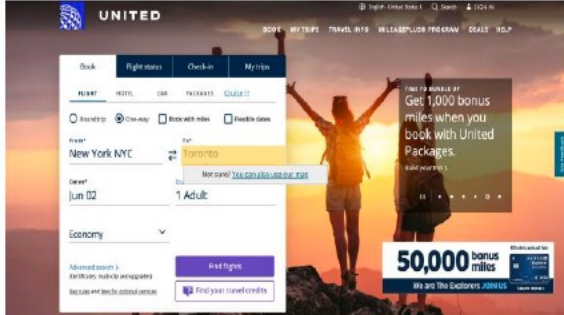
"Name it time to walk"  
input\_text <"time to walk">



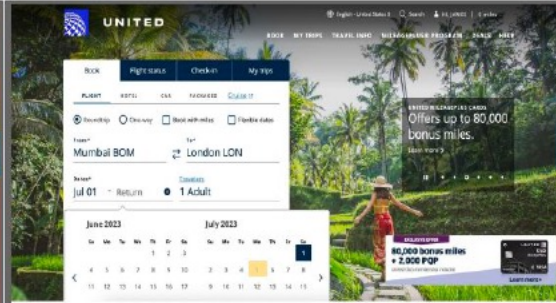
"Click the OK button"  
click <842,918>



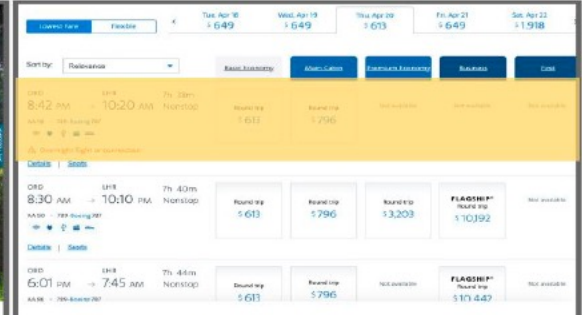
# Web Agents



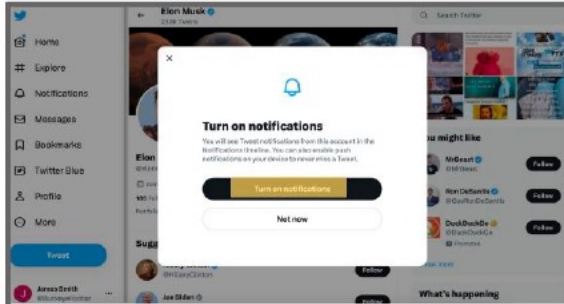
(a) Find one-way flights from New York to Toronto.



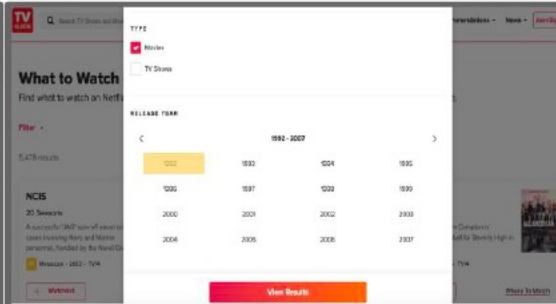
(b) Book a roundtrip on July 1 from Mumbai to London and vice versa on July 5 for two adults.



(c) Find a flight from Chicago to London on 20 April and return on 23 April.



(d) Find Elon Musk's profile and follow, start notifications and like the latest tweet.



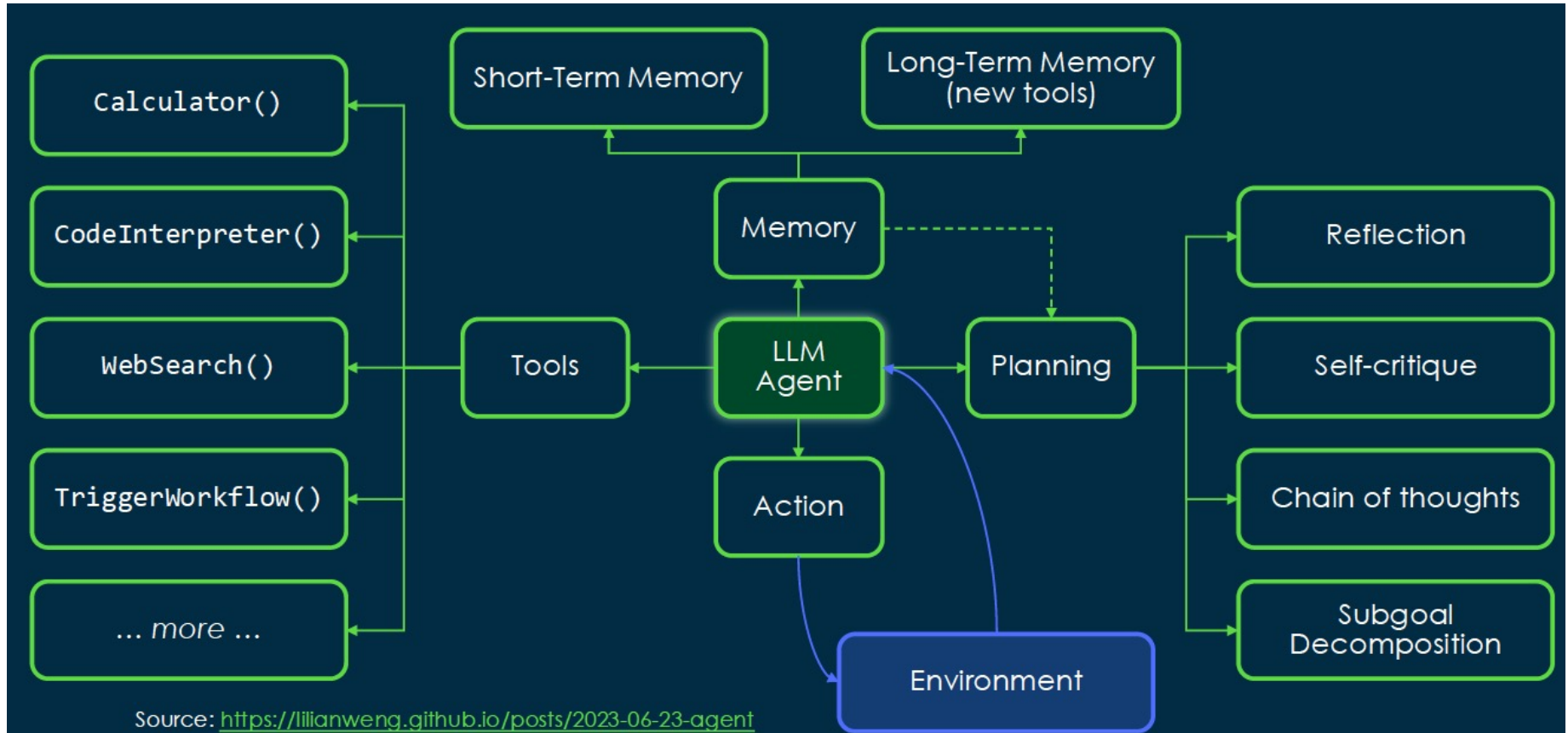
(e) Browse comedy films streaming on Netflix that was released from 1992 to 2007.



(f) Open page to schedule an appointment for car knowledge test.

World of Bits: An Open-Domain Platform for Web-Based Agents, (Shi et al., 2017)  
Mind2Web: Towards a Generalist Agent for the Web, (Deng et al., 2023)  
WebArena: A Realistic Web Environment for Building Autonomous Agents, (Zhou et al., 2023)  
[Browsergym: a Gym Environment for Web Task Automation \(Drouin et al., 2024\)](#)

# Recap: Architecture of an AI Agent





# What's special / difficult with Web Agents ?

## API Agents

- ❖ Observations: API call results, search history, user-uploaded images, chat history
- ❖ Actions: API calls, search calls, responses to the user
- ❖ Pros: Lower latency, lower risks
- ❖ Cons: needs appropriate APIs

## Web Agents

- ❖ Observations: what human would see + accessibility tree / raw DOM
- ❖ Actions: enter text in fields, clicks
- ❖ Pros: can do anything
- ❖ Cons: higher latency, higher risks



# Some Challenges in Web Agents

MUST be able to deal with the **World WILD Web** !

- ❖ Long Context Understanding
  - ▶ HTML pages are complex, easily filling up > 100K tokens
- ❖ Long-term planning
  - ▶ Need to infer/ understand/ reason, set and meet complex objectives / constraints
  - ▶ Need to visit many websites, process man webpages to search for the right path ; errors can easily be cumulated and exponentially compounded
- ❖ Learning and adaptability
  - ▶ Messy HTML, Interactive/ Dynamic elements, New Features/ web programming construct/ languages!
- ❖ Strong Multimodality support/ understanding
  - ▶ Need to process and understand not only text feedback but also "pixel"s !
- ❖ Cost and efficiency
  - ▶ To be viable, Web Agents must produce more value than they cost !
- ❖ Safety and Alignment
  - ▶ Warn and protect you when your Web Agent logs into your Bank Account to wire your funds away !

# Web Agent Research Milestones

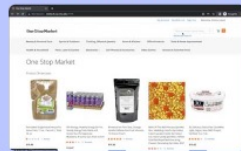
## Learning to Control Computers (DM)

- Control computers w/ keyboard & mouse from NL instructions
- MiniWoB++ through RL with computer-human interactions



## WebArena (CMU)

- Realistic benchmark, 812 tasks, 6 domains
- Long-horizon tasks
- Best GPT-4: 11% solve rate vs 78% for humans



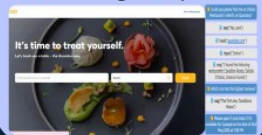
## VisualWebArena

- Benchmark that needs visual comprehension
- Test visual & reasoning skills of web agents
- 910 tasks, 3 domains



## WebLINX (McGill)

- Conversational web agent navigation
- 2337 expert demos on 155 real-world websites
- Visual models not best; fine-tuning is key



2017

2021

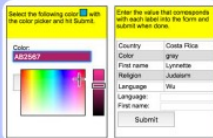
2022

2023

2024

## World of Bits (WoB)

- First widely available web benchmark
- Simplified tasks
- 100 tasks
- Can be solved by RL



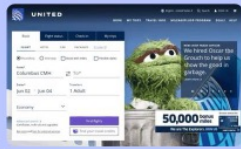
## WebGPT (OpenAI)

- Fine-tuned GPT-3 for QA with web browsing
- Evaluated on "Explain Like I'm 5" Reddit Qs + TruthfulQA dataset



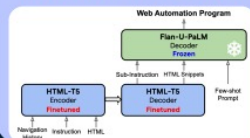
## Mind2Web (Ohio)

- Benchmark of realistic web tasks from NL
- Interaction traces
- 2,350 tasks from 137 websites, 31 domains



## WebAgent (Google)

- Combine 2 LLMs to simplify huge HTML, plan solution, create code talking to web browser; no pixels
- MiniWoB & Mind2Web



## WebVoyager (Teng)

- Completes tasks on real websites using textual+visual inputs
- New benchmark: 15 websites, automatic GPT-4V-based eval.



# Web Agent Research Milestones (cont'd)

U)

mark,  
ains  
ks  
olve  
umans



## VisualWebArena

- Benchmark that needs **visual comprehension**
- Test visual & reasoning skills of web agents
- 910 tasks, 3 domains



## WebLINX (McGill)

- **Conversational** web agent navigation
- 2337 expert demos on 155 real-world websites
- Visual models not best; fine-tuning is key



## WorkArena (ServiceNow)

- Basic tasks that a knowledge worker must carry out
- Implemented on the ServiceNow platform



## OSWorld

- 369 computer tasks of real web and desktop apps in open domains
- OS file I/O + workflows spanning multiple applications



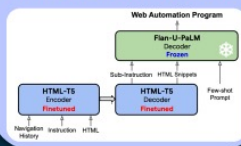
## WorkArena++ (ServiceNow)

- Compositional tasks with much higher difficulty than WorkArena
- Today's best models get single-digit performance, with huge room for improvement

2024

## WebAgent (Google)

- Combine 2 LLMs to simplify huge HTML, plan solution, create code talking to web browser; **no pixels**
- MiniWoB & Mind2Web



## WebVoyager (Tenor)

- Completes tasks on real websites using **textual+visual** inputs
- New benchmark: 15 websites, automatic GPT-4V-based eval.



## WebCanvas (CMU)

- Handles dynamic web
- Mind2Web-Live, a refined Mind2Web: 542 tasks, 2439 evaluation states



## AssistantBench

- Diverse web tasks: search, navigation, data extraction, interaction
- 214 tasks that can be auto-evaluated



## NNetNav (Stanford)

- Training web agents entirely through synthetic demos
- Web trajectory rollouts are processed by an LLM to be retroactively labeled into instruction



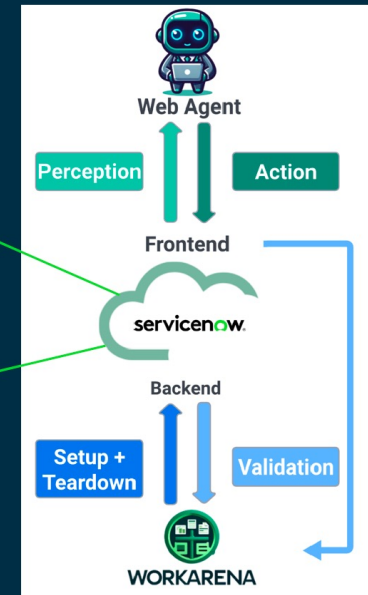
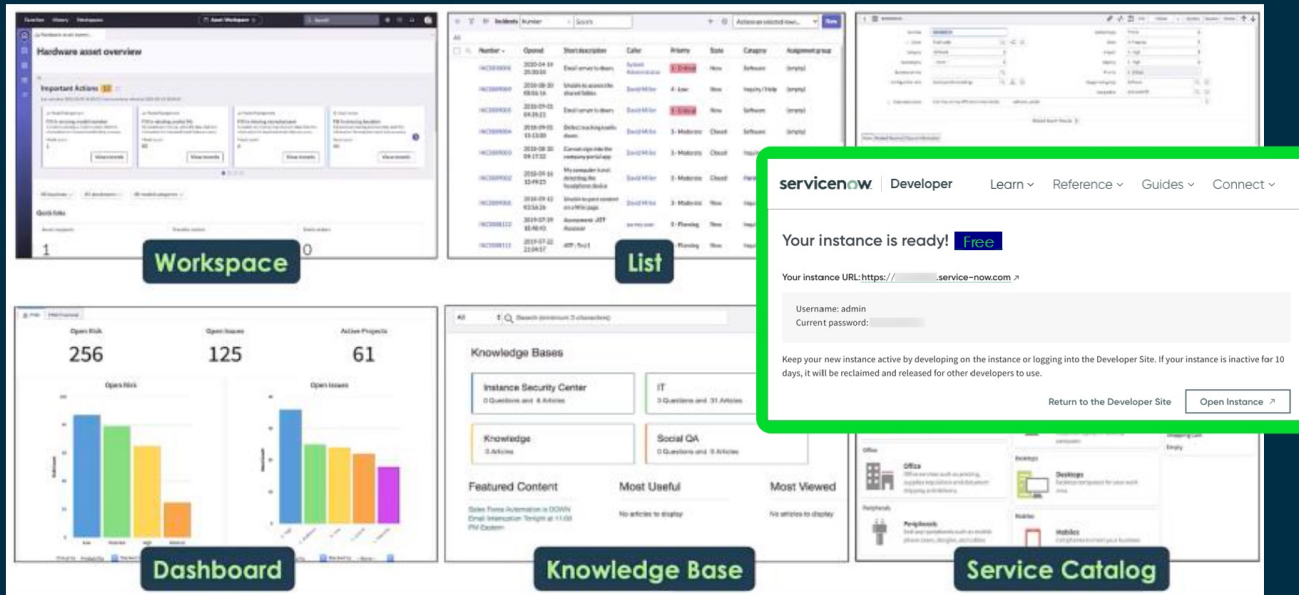
# Many Web Agent Benchmarks but lack unification

- ❖ MiniWoB++ (Shi et al., 2017; Liu et al., 2018) 125 tasks
- ❖ WebShop (Yao, Chen et al., 2022) 12087 tasks
- ❖ WebArena (Zhou et al., 2023) 812 tasks
- ❖ VisualWebArena (Koh et al., 2024) 910 tasks
- ❖ WebLINX (Lù et al., 2024) 2300 tasks
- ❖ WebCanvas (Pan et al., 2024) 438 tasks
- ❖ WebVoyager (He et al., 2024) 643 tasks
- ❖ AssistantBench (Yoran et al., 2024) 214 tasks
- ❖ WorkArena++ (ServiceNow Research, 2024) 682 tasks



# An Example: WorkArena a Benchmark for Enterprise Workflow Web Agent

An open-source benchmark of ~600 work-related tasks built on the ServiceNow platform



Tasks span basic UI interactions and complex realistic workflows

Open Web

# Another Example: WorkArena++: a Benchmark towards Realistic Enterprise Workflows

1

## Knowledge base

Knowledge Bases

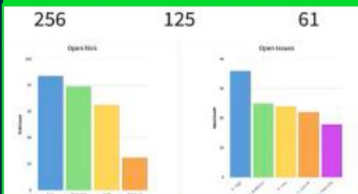
Instance Security Center 0 Questions and 8 Articles	IT 0 Questions and 31 Articles
Knowledge 0 Articles	Social QA 0 Questions and 8 Articles

Featured Content    Most Useful    Most Viewed

No articles to display    No articles to display

2

## Dashboard



3

## Service Catalog

Service Catalog

Services View and manage your services and products for quality, performance, and availability.	Hardware View and manage your hardware items and their configurations.
Software View and manage your software products and their configurations.	Office View and manage your office equipment and their configurations.

**Example:** The agent is assigned a ticket and instruction: "Please solve this."

Private Task  
Clean-up your duplicate problems

Number: PTK47711968    Priority: 4 - Low

Owner: Sandy Martinez    State: Open

Assigned to: Sandy Martinez    Parent: [Search]

Active:

Short description: Retrieve information from the chart with the title #CAT044377552 and perform the mentioned task. For calculations, please round off to the

Description: You have to retrieve some information from a dashboard chart based on the description below. The chart presents the number of 'hardware items' available in stock. After retrieving the information, you will be asked to use it to complete a task.

Title of the report: #CAT044377552

Referring to the company protocol 'Dashboard Retrieve Information and Perform Task' (located in the 'Company Protocols' knowledge base), complete the dashboard retrieval task.

- Please retrieve the 'greatest' value of all the items in stock.
- Task: Place an order for the least available item in stock. The quantity of the order should be such that the final quantity of this item matches the above retrieved value.  
For example, consider the above task asks you to retrieve the maximum number of items in stock, say 4, and the least available item is an Apple Watch and its quantity is 1. You have to order 3 more Apple Watches.
- Please do not change any other configuration while placing the order for the item. You can find important links to the pages in the protocol article.

Don't forget to mark this task as "Closed - complete" once successfully completed. If the task appears infeasible, mark the task as "Closed - skipped".

# WorkArena++: a Benchmark towards Realistic Enterprise Workflows

## Overview of tasks

Solve a series of enterprise decision-making problems:

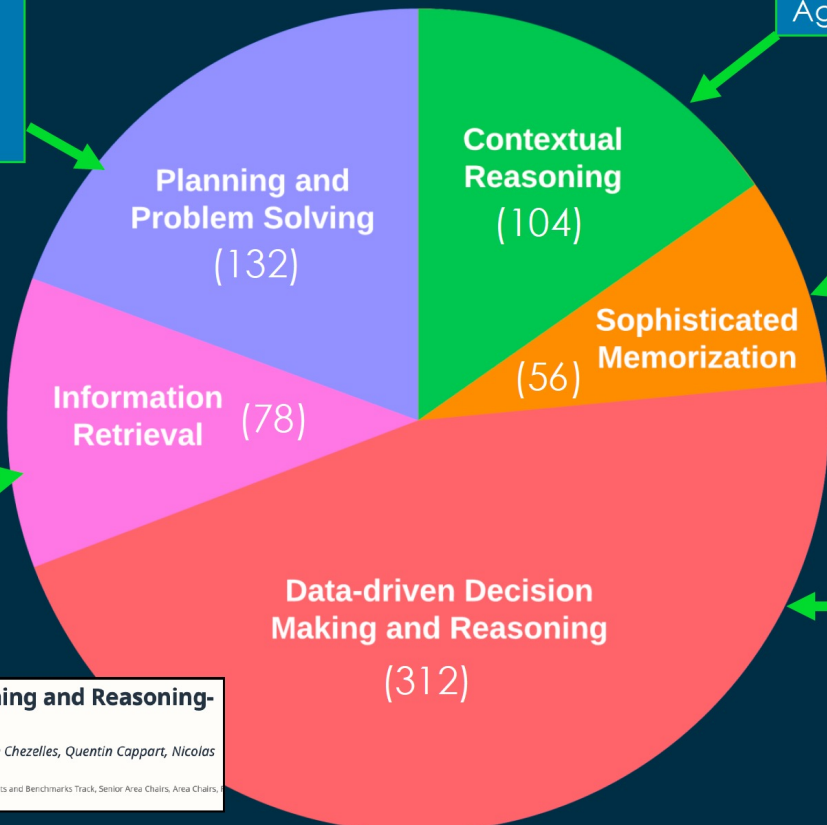
- Workload balancing
- Scheduling with constraints
- Assigning work to experts

Read dashboards and act:

- Restock IT asset inventory

Search for information in lists, forms, and KBs:

- Find if a user's laptop is under warranty



Some tasks are purposely infeasible. Agent must detect this.

Navigate the platform to gather multiple bits of information and then solve a task:

- Offboard a user

Read dashboard, make calculations, take action

Budget management: choose where to invest based on expected return

Expense management



# Weak Performance on the WorkArena++ Benchmark

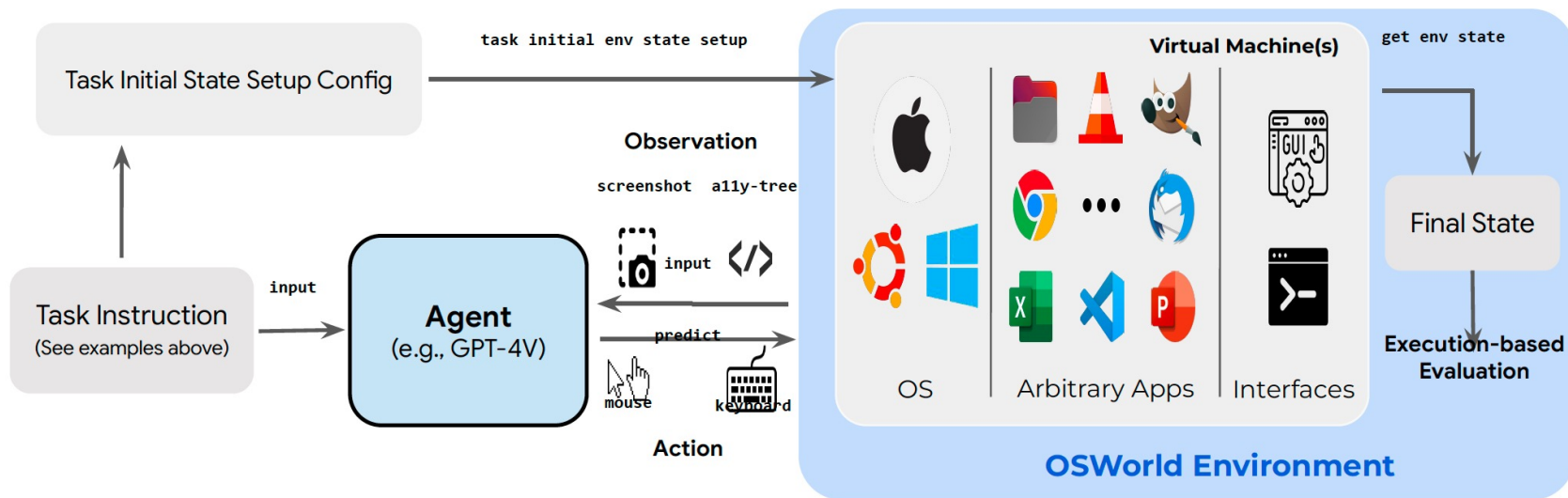
Success rate (higher is better)

Task Category (task count)	Agent Curriculum (full benchmark)					Human
	GPT-3.5	GPT-4o	GPT-4o-v	Llama3	Mixtral	
<b>WorkArena L3</b> (235)	<b>0.0</b> ±0.0	<b>0.0</b> ±0.0	<b>0.0</b> ±0.0	<b>0.0</b> ±0.0	<b>0.0</b> ±0.0	<b>93.9</b> ±3.4
Contextual Understanding (32)	0.0 ±0.0	0.0 ±0.0	0.0 ±0.0	0.0 ±0.0	0.0 ±0.0	87.5 ±11.7
Data-driven Decision-Making (55)	0.0 ±0.0	0.0 ±0.0	0.0 ±0.0	0.0 ±0.0	0.0 ±0.0	100.0 ±0.0
Planning and Problem Solving (44)	0.0 ±0.0	0.0 ±0.0	0.0 ±0.0	0.0 ±0.0	0.0 ±0.0	87.5 ±11.7
Information Retrieval (56)	0.0 ±0.0	0.0 ±0.0	0.0 ±0.0	0.0 ±0.0	0.0 ±0.0	100.0 ±0.0
Sophisticated Memorization (48)	0.0 ±0.0	0.0 ±0.0	0.0 ±0.0	0.0 ±0.0	0.0 ±0.0	91.7 ±8.0
<b>WorkArena L2</b> (235)	<b>0.0</b> ±0.0	<b>3.0</b> ±1.1	<b>3.8</b> ±1.3	<b>0.0</b> ±0.0	<b>0.0</b> ±0.0	<b>93.9</b> ±3.4
Contextual Understanding (32)	0.0 ±0.0	0.0 ±0.0	0.0 ±0.0	0.0 ±0.0	0.0 ±0.0	100.0 ±0.0
Data-driven Decision-Making (55)	0.0 ±0.0	0.0 ±0.0	0.0 ±0.0	0.0 ±0.0	0.0 ±0.0	84.6 ±10.0
Planning and Problem Solving (44)	0.0 ±0.0	0.0 ±0.0	0.0 ±0.0	0.0 ±0.0	0.0 ±0.0	100.0 ±0.0
Information Retrieval (56)	0.0 ±0.0	0.0 ±0.0	3.6 ±2.5	0.0 ±0.0	0.0 ±0.0	100.0 ±0.0
Sophisticated Memorization (48)	0.0 ±0.0	14.6 ±5.1	14.6 ±5.1	0.0 ±0.0	0.0 ±0.0	91.7 ±8.0
<b>WorkArena L1</b> (33 × 10 seeds)	<b>6.1</b> ±1.3	<b>42.7</b> ±2.7	<b>41.8</b> ±2.7	<b>17.9</b> ±2.1	<b>12.4</b> ±1.8	
<b>MiniWoB</b> (125 × 5 seeds)	<b>43.4</b> ±1.6	<b>71.3</b> ±1.5	<b>72.5</b> ±1.5	<b>68.2</b> ±1.2	<b>62.4</b> ±1.6	
<b>WebArena</b> (812)	<b>6.7</b> ±0.9	<b>23.5</b> ±1.5	<b>24.0</b> ±1.5	<b>11.0</b> ±1.1	<b>12.6</b> ±0.5	

What explains this?

- Failure to plan
- Hallucinated controls
- Incorrect action syntax

# OSWorld: a Benchmark for Computer Use



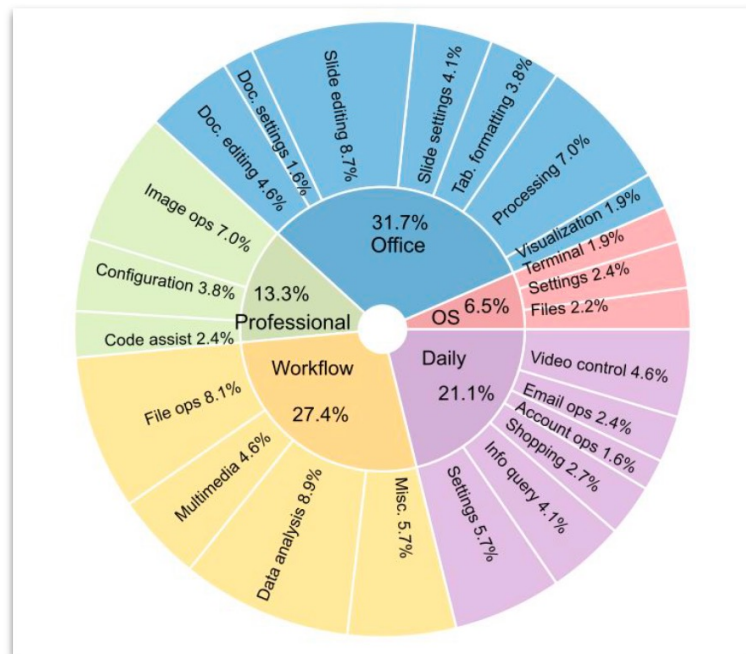
# OSWorld Benchmark Dataset

369 real-world computer tasks that involve real web and desktop apps in open domains, OS file I/O, and multi-app workflows. Each task example is annotated with

- A real-world task instruction from real users
- An initial state setup config to simulate human work in progress
- A custom execution-based evaluation script

Table 3: Key statistics in OSWORLD. The “Supp. tasks” refers to the Windows-based tasks, that could only be used after activation due to copyright restrictions.

Statistic	Number
Total tasks (Ubuntu)	369 (100%)
- Multi-App Workflow	101 (27.4%)
- Single-App	268 (72.6%)
- Integrated	84 (22.8%)
- Infeasible	30 (8.1%)
Supp. tasks (Windows)	43
Initial States	302
Eval. Scripts	134



# Performance on OSWorld Benchmark

Inputs	Model	Success Rate (↑)					
		OS	Office	Daily	Profess.	Workflow	Overall
A11y tree	Mixtral-8x7B	12.50%	1.01%	4.79%	6.12%	0.09%	2.98%
	Llama-3-70B	4.17%	1.87%	2.71%	0.00%	0.93%	1.61%
	GPT-3.5	4.17%	4.43%	2.71%	0.00%	1.62%	2.69%
	GPT-4	20.83%	3.58%	25.64%	26.53%	2.97%	<b>12.24%</b>
	Gemini-Pro	4.17%	1.71%	3.99%	4.08%	0.63%	2.37%
	Gemini-Pro-1.5	12.50%	2.56%	7.83%	4.08%	3.60%	4.81%
	Qwen-Plus	29.17%	3.58%	8.36%	10.20%	2.61%	6.87%
GPT-4o	20.83%	6.99%	16.81%	16.33%	7.56%	11.36%	
Screenshot	CogAgent	4.17%	0.85%	2.71%	0.00%	0.00%	1.11%
	GPT-4V	12.50%	1.86%	7.58%	4.08%	<b>6.04%</b>	5.26%
	Gemini-ProV	8.33%	3.58%	6.55%	16.33%	2.08%	<b>5.80%</b>
	Gemini-Pro-1.5	12.50%	6.99%	2.71%	6.12%	3.60%	5.40%
	Claude-3-Opus	4.17%	1.87%	2.71%	2.04%	2.61%	2.42%
	GPT-4o	8.33%	3.58%	6.07%	4.08%	5.58%	5.03%
Screenshot + A11y tree	CogAgent	4.17%	0.85%	2.71%	0.62%	0.09%	1.32%
	GPT-4V	16.66%	6.99%	24.50%	18.37%	4.64%	<b>12.17%</b>
	Gemini-ProV	4.17%	4.43%	6.55%	0.00%	1.52%	3.48%
	Gemini-Pro-1.5	12.50%	3.58%	7.83%	8.16%	1.52%	5.10%
	Claude-3-Opus	12.50%	3.57%	5.27%	8.16%	1.00%	4.41%
	GPT-4o	41.67%	6.16%	12.33%	14.29%	<b>7.46%</b>	11.21%
Set-of-Mark	CogAgent	4.17%	0.00%	2.71%	0.00%	0.53%	0.99%
	GPT-4V	8.33%	8.55%	22.84%	14.28%	<b>6.57%</b>	<b>11.77%</b>
	Gemini-ProV	4.17%	1.01%	1.42%	0.00%	0.63%	1.06%
	Gemini-Pro-1.5	16.67%	5.13%	12.96%	10.20%	3.60%	7.79%
	Claude-3-Opus	12.50%	2.72%	14.24%	6.12%	4.49%	6.72%
	GPT-4o	20.83%	3.58%	3.99%	2.04%	3.60%	4.59%
Human Performance		75.00%	71.79%	70.51%	73.47%	73.27%	72.36%

- LLMs and VLMs are still far from being digital agents on real computers.
- Agent performance fluctuations vs. consistent human performance across different types of computer tasks.
- A11y tree and SoM's effectiveness varies by models.
- VLM agents with screenshot-only setting show lower performance, but it should be the ultimate configuration in the long run.

# Some Recent Progress on OSWorld Benchmark Performance

## Introducing computer use, a new Claude 3.5 Sonnet, and Claude 3.5 Haiku

Oct 22, 2024 • 5 min read

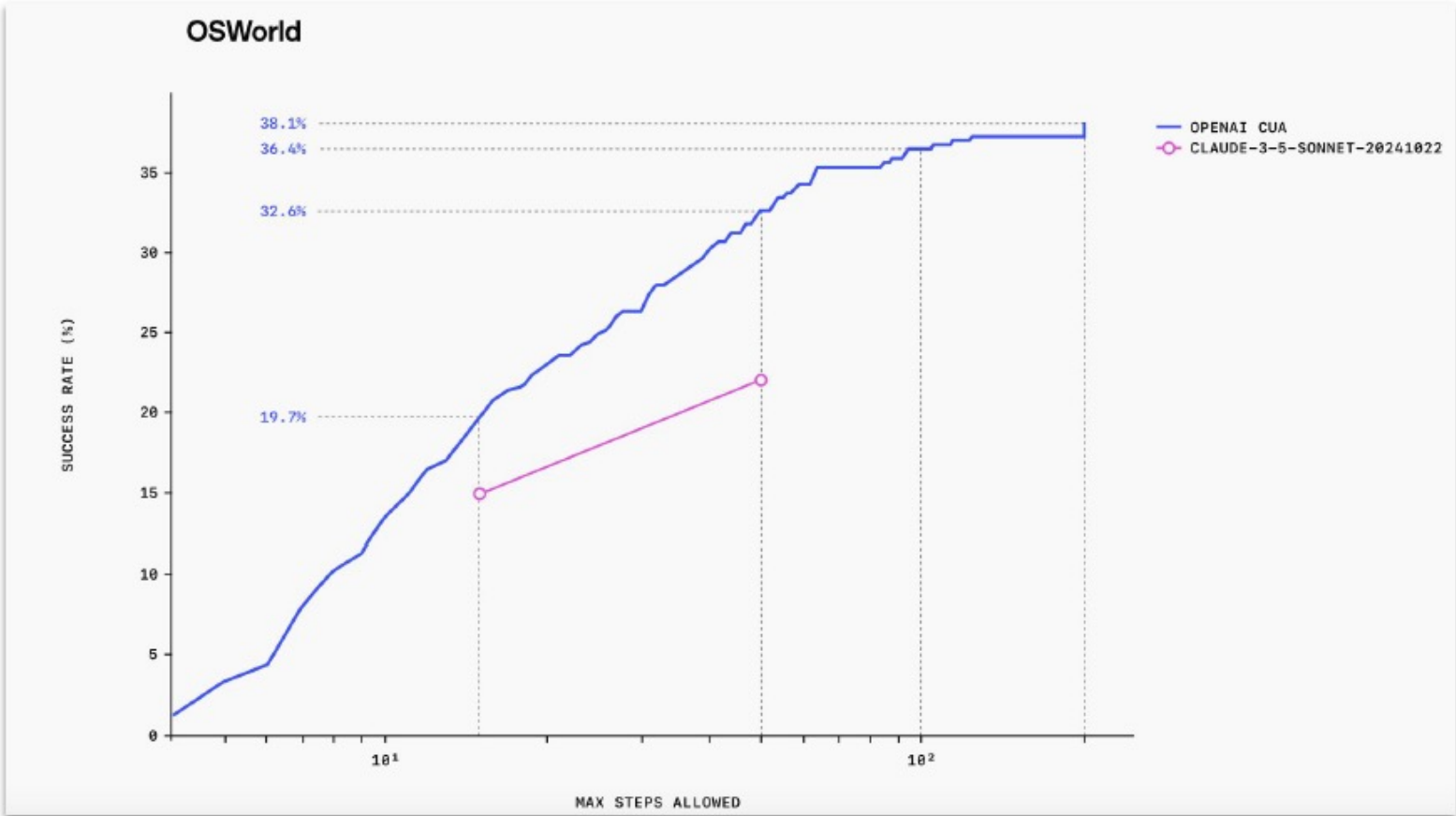


Category	Claude 3.5 Sonnet (New) - 15 steps		Claude 3.5 Sonnet (New) - 50 steps		Human Success Rate [3]
	Success Rate	95% CI	Success Rate	95% CI	
OS	54.2%	[34.3, 74.1]%	41.7%	[22.0, 61.4]%	<b>75.00%</b>
Office	7.7%	[2.9, 12.5]%	17.9%	[11.0, 24.8]%	<b>71.79%</b>
Daily	16.7%	[8.4, 25.0]%	24.4%	[14.9, 33.9]%	<b>70.51%</b>
Professional	24.5%	[12.5, 36.5]%	40.8%	[27.0, 54.6]%	<b>73.47%</b>
Workflow	7.9%	[2.6, 13.2]%	10.9%	[4.9, 17.0]%	<b>73.27%</b>
Overall	14.9%	[11.3, 18.5]%	22%	[17.8, 26.2]%	<b>72.36%</b>

Anthropic computer use agent results on OSWorld

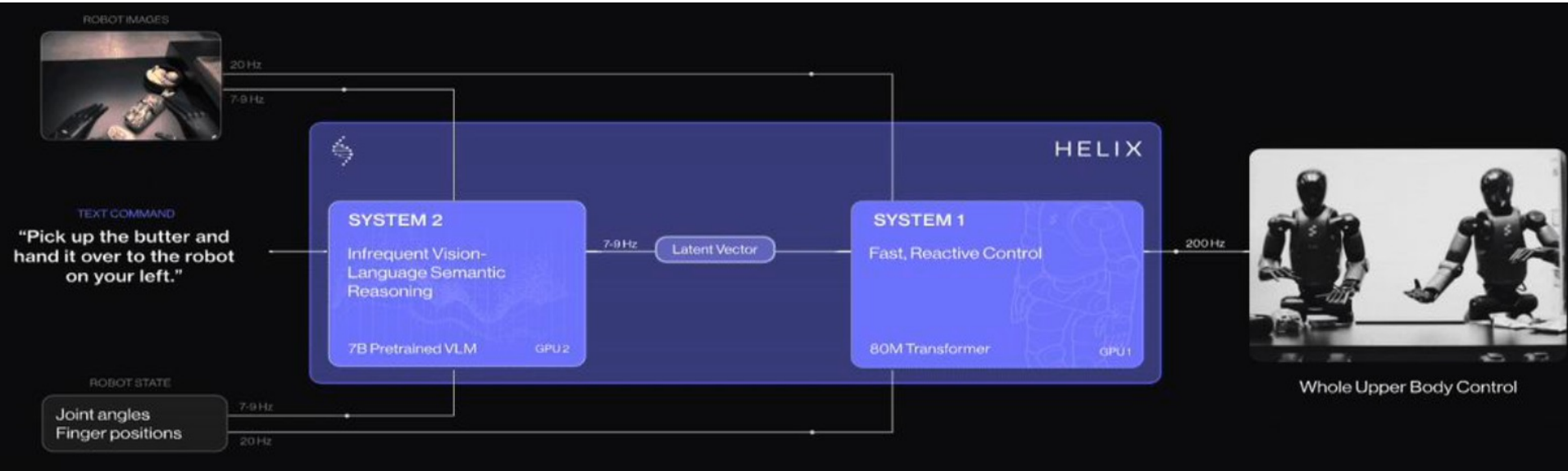


# Some Recent Progress on OSWorld Benchmark Performance



<https://openai.com/index/computer-using-agent/>

# Robotic Process Automation (RPA) w/ Physical Agents





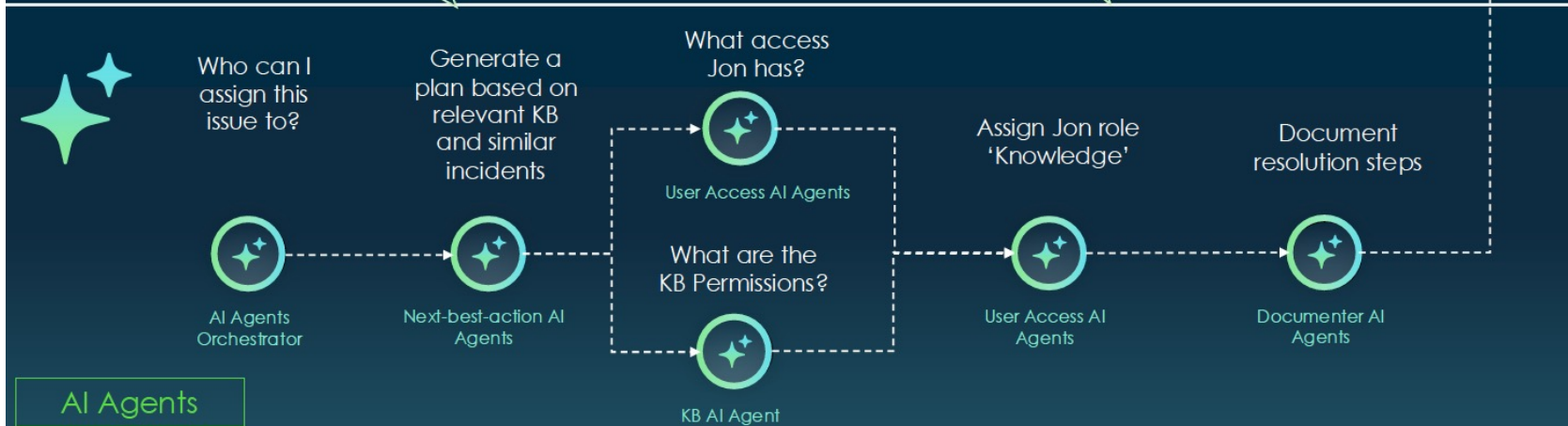
# Impact of AI Agents on Job Nature

## Today's Enterprise Workflows Remain Quite Manual (even with generative AI)

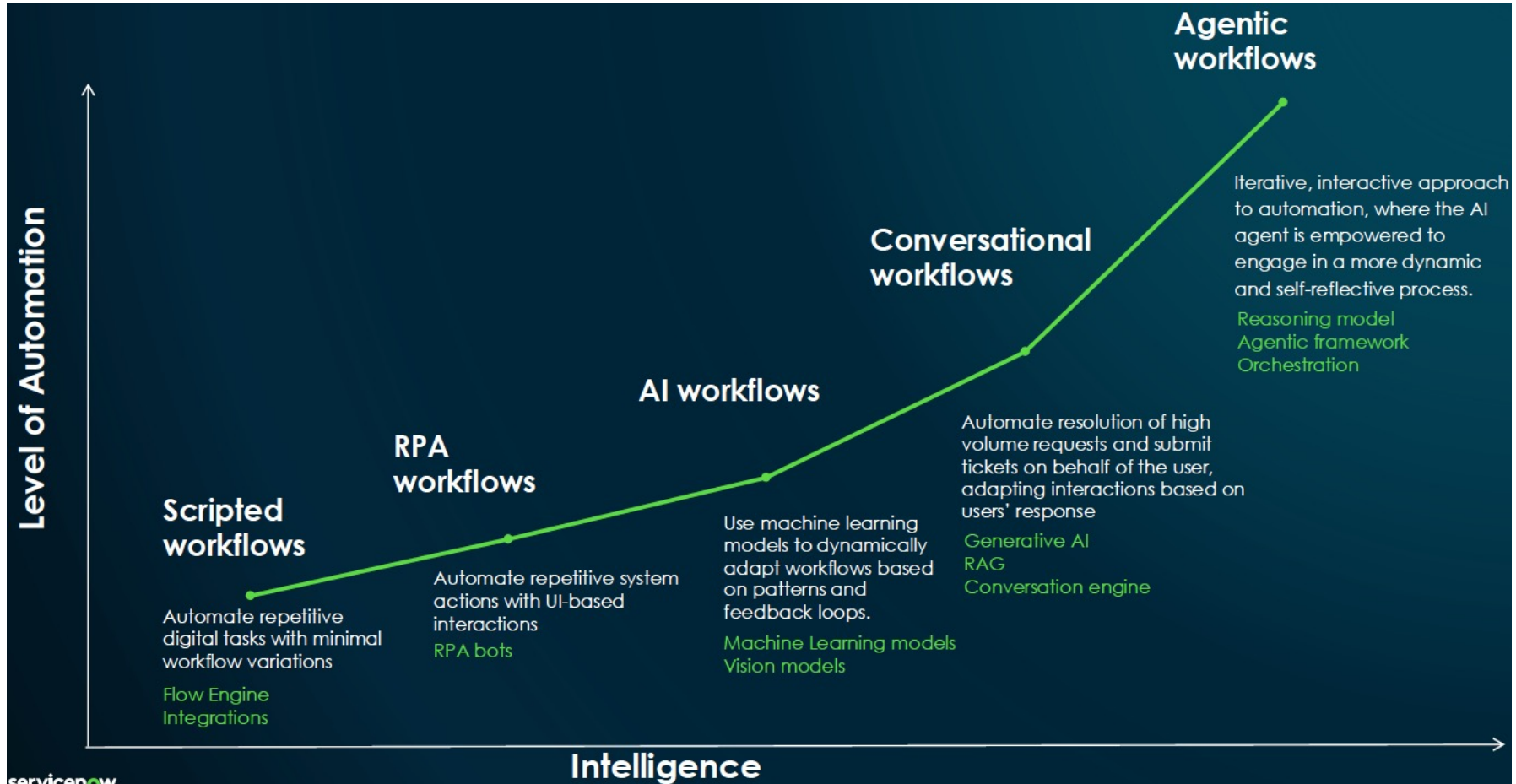


# Impact of AI Agents on Job Nature (cont'd)

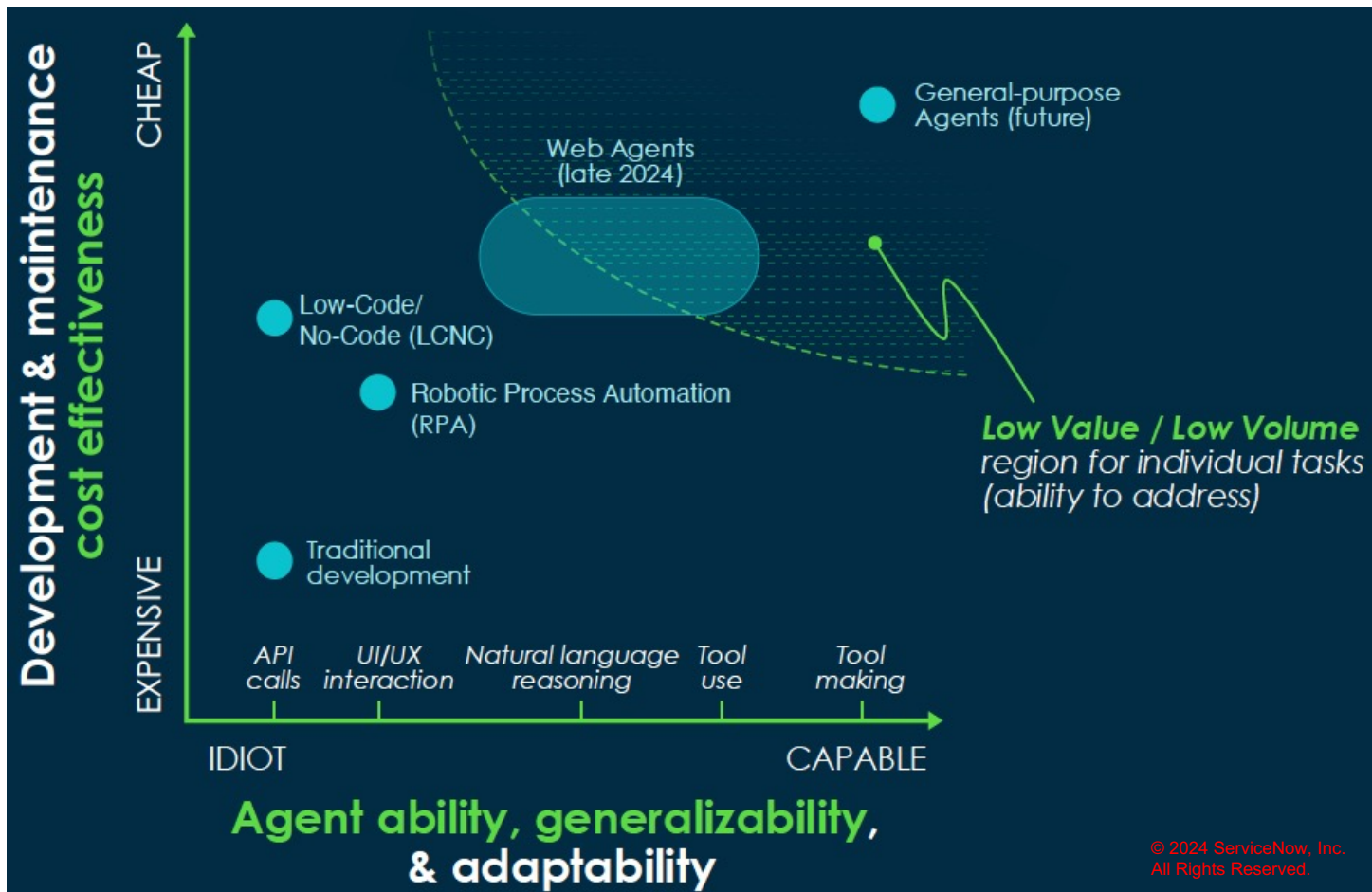
## Access issue – With AI Agents



# Automation in Enterprise Workflows



# Current Target Area of Web Agents



# Common Failure Modes of Current Web Agents and Future Directions

- ❖ Lack Long-Horizon Reasoning and Planning
  - ▶ Models oscillate between two webpages, or get stuck in a loop
  - ▶ Correctly performing tasks but undoing them
  - ▶ Agents tend to stop exploration / execution too early

## How to Close the Gap ?

- ❖ Allow agent to Search, execute and coordinate multiple instances in parallel and ask for clarifications/confirmations
- ❖ Develop Stronger Multimodal, i.e. Vision-Language-Code Models
- ❖ Identify the appropriate level of abstraction for Agents (HTML/ Screenshots / APIs)