

# AI Agents in LangGraph

Last updated: July 4, 2024

## Contents

<b>1</b>	<b>Introduction</b>	<b>1</b>
1.1	Agent Workflow . . . . .	1
<b>2</b>	<b>Build an Agent from Scratch</b>	<b>2</b>
2.1	Generic Agent . . . . .	2
2.2	ReAct Agent . . . . .	2
2.3	Add Loop . . . . .	4
<b>3</b>	<b>LangGraph Components</b>	<b>5</b>
3.1	Prompt Templates . . . . .	5
3.2	Tools . . . . .	5
3.3	Introduction to LangGraph . . . . .	5
3.3.1	Graphs . . . . .	5
3.3.2	State . . . . .	5
3.4	Building an Agent . . . . .	6
<b>4</b>	<b>Agentic Search Tools</b>	<b>9</b>
4.1	How an Agent Uses Search . . . . .	9
4.1.1	Regular Search . . . . .	9
4.1.2	Agentic Search . . . . .	10
<b>5</b>	<b>Persistence and Streaming</b>	<b>12</b>
5.1	Agent with Persistence and Streaming . . . . .	12
5.1.1	Streaming Tokens . . . . .	14
<b>6</b>	<b>Human in the Loop</b>	<b>15</b>
6.1	Setting Up the Agent . . . . .	15
6.2	Manual Human Approval . . . . .	15
6.2.1	Continue After Interrupt . . . . .	16
6.3	Snapshot and State Management . . . . .	17
6.3.1	Modify State . . . . .	17
6.3.2	Time Travel . . . . .	18
6.3.3	Go Back in Time and Edit . . . . .	18
6.3.4	Add Message to a State . . . . .	19
<b>7</b>	<b>Essay Writer</b>	<b>20</b>
7.1	Overview . . . . .	20
7.2	Implementation . . . . .	20
<b>8</b>	<b>LangChain Resources</b>	<b>24</b>
<b>9</b>	<b>Conclusion</b>	<b>25</b>
9.1	Multi-Agent Architecture . . . . .	25
9.2	Supervisor Architecture . . . . .	25
9.3	Plan and Execute Architecture . . . . .	25
9.4	Language Agent Tree Search Architecture . . . . .	25

# 1 Introduction

AI agents are systems designed to *perform tasks autonomously* by interacting with each other and their environment. These agents follow a structured workflow, known as agent workflow, which allows them to achieve better results compared to one-shot requests by *acting iteratively*.

## 1.1 Agent Workflow

The key design patterns of an agentic workflow include:

- **Planning:** Identifying the steps needed to build and execute the workflow.
- **Tool Use:** Knowing which tools are available and how to utilize them effectively.
- **Reflection:** Iteratively improving results, often using multiple LLMs to provide useful suggestions for enhancing the editing cycle.
- **Multi-agent Communication:** Each LLM can be assigned a specific role, with a unique prompt to perform that role.
- **Memory:** Tracking progress and results over multiple steps to maintain continuity.

Some of these capabilities are inherent to the LLMs themselves, while others are implemented outside of the LLMs by the framework they operate within.

There are two additional capabilities useful in constructing an agent:

- *Human Input Reception:* Allowing agents to receive guidance at critical points.
- *Persistence:* The ability to store the current state of information so that it can be referenced later.

## 2 Build an Agent from Scratch

We will build an agent based on the ReAct pattern (Reason + Act). First, an LLM thinks about what to do and then decides on the actions to take. This action is then executed in an environment, and an observation is returned. With this observation, the LLM repeats (thinks again about what to do, decides on another action, and so on) until it decides that it is done.

### 2.1 Generic Agent

To start, let's import all the necessary libraries:

```
1 import openai
2 import re
3 import httpx
4 import os
5 from dotenv import load_dotenv
6
7 _ = load_dotenv()
8 from openai import OpenAI
```

Next, we initialize the LLM (we use OpenAI in this case):

```
1 client = OpenAI()
```

Now, let's create the actual agent by defining a class for the agent:

```
1 class Agent:
2     def __init__(self, system=""):
3         self.system = system
4         self.messages = []
5         if self.system:
6             self.messages.append({"role": "system", "content": system})
7
8     def __call__(self, message):
9         self.messages.append({"role": "user", "content": message})
10        result = self.execute()
11        self.messages.append({"role": "assistant", "content": result})
12        return result
13
14    def execute(self):
15        completion = client.chat.completions.create(
16            model="gpt-4o",
17            temperature=0,
18            messages=self.messages)
19        return completion.choices[0].message.content
```

The agent needs to be parameterized by a system agent, so we allow the user to pass that in and save it as an attribute. We also keep track of the list of messages over time (the first being the system message).

The `__call__` method *takes the user's message, appends it to the list*, executes the *execute* function which *returns a response, and appends the response to the list*. The *execute* method is what actually calls the LLM with all the accumulated messages so far and returns the model's response. This is our generic agent.

### 2.2 ReAct Agent

The ReAct agent requires a specific system prompt:

```
1 prompt = """
2 You run in a loop of Thought, Action, PAUSE, Observation.
3 At the end of the loop you output an Answer
4 Use Thought to describe your thoughts about the question you have been asked.
5 Use Action to run one of the actions available to you - then return PAUSE.
6 Observation will be the result of running those actions.
7
8 Your available actions are:
9
10 calculate:
11 e.g. calculate: 4 * 7 / 3
12 Runs a calculation and returns the number - uses Python so be sure to /
```

```

13 use floating point syntax if necessary
14
15 average_dog_weight:
16 e.g. average_dog_weight: Collie
17 returns average weight of a dog when given the breed
18
19 Example session:
20
21 Question: How much does a Bulldog weigh?
22 Thought: I should look the dogs weight using average_dog_weight
23 Action: average_dog_weight: Bulldog
24 PAUSE
25
26 You will be called again with this:
27
28 Observation: A Bulldog weights 51 lbs
29
30 You then output:
31
32 Answer: A bulldog weights 51 lbs
33 """.strip()

```

Besides explaining the cycle of thinking, acting, and observing the result of the action, it's important to explain in the prompt what actions the agent can take and provide an example of how to use them. In this case, these are two functions we define here:

```

1 def calculate(what):
2     return eval(what)
3
4 def average_dog_weight(name):
5     if name in "Scottish Terrier":
6         return("Scottish Terriers average 20 lbs")
7     elif name in "Border Collie":
8         return("a Border Collies average weight is 37 lbs")
9     elif name in "Toy Poodle":
10        return("a toy poodles average weight is 7 lbs")
11    else:
12        return("An average dog weights 50 lbs")
13
14 known_actions = {
15     "calculate": calculate,
16     "average_dog_weight": average_dog_weight
17 }

```

If we try it now, we first need to call the agent and then ask a question:

```

1 abot = Agent(prompt)
2 result = abot("How much does a toy poodle weigh?")
3 print(result)

```

The response will be something like: *Thought: I should look up the dog's weight using average\_dog\_weight for a Toy Poodle. Action: average\_dog\_weight: Toy Poodle. PAUSE*

This means that the function with that parameter needs to be called:

```

1 result = average_dog_weight("Toy Poodle")
2 print(result)

```

We will get the result *a toy poodles average weight is 7 lbs*; we can pass this as the next prompt to the LLM:

```

1 next_prompt = "Observation: {}".format(result)
2 abot(next_prompt)

```

If we want more details on what happened throughout the entire process, we can use the *messages* attribute to print all the messages in the conversation:

```

1 abot.messages

```

The final output will be *Answer: A Toy Poodle weighs 7 lbs.*

Now, let's see a more complex use case:

```

1 abot = Agent(prompt)
2 question = """I have 2 dogs, a border collie and a scottish terrier. \
3 What is their combined weight"""
4 abot(question)

```

The response will be *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. Action: average\_dog\_weight: Border Collie PAUSE*

Execute the action, insert it into a new prompt, and call the agent again:

```
1 next_prompt = "Observation: {}".format(average_dog_weight("Border Collie"))
2 abot(next_prompt)
```

The response will be *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. Action: average\_dog\_weight: Scottish Terrier PAUSE*

Again, execute the action, insert it into a new prompt, and call the agent:

```
1 next_prompt = "Observation: {}".format(average_dog_weight("Scottish Terrier"))
2 abot(next_prompt)
```

The response will be *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. Action: calculate: 37 + 20 PAUSE*

Execute this last action, which now calls a different function:

```
1 next_prompt = "Observation: {}".format(calculate("37 + 20"))
2 abot(next_prompt)
```

The final response will be *Answer: The combined weight of a Border Collie and a Scottish Terrier is 57 lbs.*

## 2.3 Add Loop

What we've seen so far is the process that needs to be executed. Now we want to put it into a *loop to automate it*.

The first step is to create a *regex* (regular expression) to match the action string. A regular expression is a sequence of characters that define a search pattern. It can be used for string matching, extraction, and replacement. In this case, the regex function will look for specific patterns in the text to determine the action to be taken:

```
1 action_re = re.compile('^Action: (\w+): (.*?)$')
2 # python regular expression to selection action
```

This splits the LLM response and determines whether we want to execute an action or if it is the final answer.

Now, let's create a *query* function. It takes a question and follows the same process we did manually:

```
1 def query(question, max_turns=5):
2     i = 0
3     bot = Agent(prompt) # initialize the agent with the default system prompt
4     next_prompt = question
5     while i < max_turns:
6         i += 1
7         result = bot(next_prompt) # call the agent with the last prompt
8         print(result)
9         actions = [
10             action_re.match(a)
11             for a in result.split('\n')
12             if action_re.match(a)
13         ] # parse the response and get a list of actions
14         if actions:
15             # There is an action to run
16             action, action_input = actions[0].groups() # function to call and the input
17             if action not in known_actions:
18                 raise Exception("Unknown action: {}: {}".format(action, action_input))
19             print("-- running {} {}".format(action, action_input))
20             # First look up in the action dictionary the action to take and
21             # we get back a function that we call with the inputs
22             observation = known_actions[action](action_input)
23             print("Observation:", observation)
24             next_prompt = "Observation: {}".format(observation)
25         else:
26             # If there aren't any action we have the final answer
27             return
```

Calling this function, for example, with the following code, will activate all the processes we previously did manually and get the final answer:

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

## 3 LangGraph Components

In the previous lesson, we built an agent from scratch; now we will create a similar one using *LangGraph*, introducing some of its components and features.

### 3.1 Prompt Templates

*PromptTemplate* allows you to create reusable prompts. This means you can create a prompt like this and use it multiple times with arbitrary variables.

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

In the *LangChain* Hub, there are various examples of such prompts, uploaded by users.

### 3.2 Tools

There are also many tools created by the community. An example is *TavilySearch*, which will be the search tool we use from now on.

```
1 from langchain_community.tools.tavily_search import TavilySearchResults
2 tool = TavilySearchResults(max_results=2)
3
4 self.tools = {t.name: t for t in tools}
5 self.model = model.bind_tools([tool])
```

### 3.3 Introduction to LangGraph

*LangGraph* is an extension of *LangChain* that supports graphs. The fundamental role of *LangGraph* is to describe and *orchestrate the control flow*, i.e., all the relationships between the components we previously created manually.

Specifically, it allows creating cyclic graphs, which support looping and relationships between elements. Single and multi-agent flows are described and represented as graphs, allowing for extremely controlled “flows.”

Additionally, it comes with built-in persistence, which is excellent for maintaining multiple conversations simultaneously or remembering previous interactions and actions. Persistence also allows for introducing human-in-the-loop features.

#### 3.3.1 Graphs

**Graphs** are composed of three basic elements:

- *Nodes*: Agents or functions.
- *Edges*: Connect nodes.
- *Conditional edges*: Used when you need to make decisions about which node to go to next.

#### 3.3.2 State

One of the most important aspects to understand in *LangGraph* is the state that is tracked over time, often called the **agent state**. This is *accessible to all parts* of the graph and can be stored in the persistence layer.

A simple state is:

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

where we have just a list of messages. The variable *messages* is a sequence of *BaseMessage*, which is a *LangChain* type.

A more complex state is:

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

Each of these four variables is of a different type. The first three are not annotated in any way, meaning when a new update is pushed to that variable, it overrides the existing value. Instead, *intermediate\_steps*, through the *add* operator, adds the new update without removing the previous ones.

## 3.4 Building an Agent

Now we will actually build the agent, composed as follows: the first node is *call\_openai* which calls the LLM; the second is a conditional edge which checks if there is an action to be taken and we will call it *exists\_action*. From here it splits into two paths: either there is no action and it goes to the final node, or there is an action and it goes to the respective action node, then loops back to *call\_openai*.

First, load an environment variable, our OpenAI key:

```
1 from dotenv import load_dotenv
2 _ = load_dotenv()
```

Next, import all necessary libraries:

```
1 from langgraph.graph import StateGraph, END
2 from typing import TypedDict, Annotated
3 import operator
4 from langchain_core.messages import AnyMessage, SystemMessage, HumanMessage, ToolMessage
5 from langchain_openai import ChatOpenAI
6 from langchain_community.tools.tavily_search import TavilySearchResults
```

Although *langchain\_openai* may seem specific to OpenAI, it can be easily swapped with any other supported model without changing anything else in the code.

The search tool we will use is *Tavily*:

```
1 tool = TavilySearchResults(max_results=4)
2 print(type(tool))
3 print(tool.name)
```

First, initialize it with *max\_results=4*, which means we will get a maximum of 4 responses from the search API. The name of the tool is *tavily\_search\_results\_json*, and this is the name the LLM will use to call this tool.

Now, create the agent state, simply an annotated list of messages:

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

Next, create the actual agent. As we have seen, we will need three functions: the first to call OpenAI, the second to check if there is an action to take, and the last to take that action:

```
1 class Agent:
2
3     def __init__(self, model, tools, system=""):
4         self.system = system
5         graph = StateGraph(AgentState)
6         graph.add_node("llm", self.call_openai)
7         graph.add_node("action", self.take_action)
8         graph.add_conditional_edges(
9             "llm",
10            self.exists_action,
11            {True: "action", False: END}
12        )
13        graph.add_edge("action", "llm")
14        graph.set_entry_point("llm")
15        self.graph = graph.compile()
16        self.tools = {t.name: t for t in tools}
17        self.model = model.bind_tools(tools)
18
19    def exists_action(self, state: AgentState):
20        result = state['messages'][-1]
21        return len(result.tool_calls) > 0
22
23    def call_openai(self, state: AgentState):
24        messages = state['messages']
25        if self.system:
26            messages = [SystemMessage(content=self.system)] + messages
27        message = self.model.invoke(messages)
28        return {'messages': [message]}
29
30    def take_action(self, state: AgentState):
31        tool_calls = state['messages'][-1].tool_calls
32        results = []
33        for t in tool_calls:
34            print(f"Calling: {t}")
```

```

35         if not t['name'] in self.tools:      # check for bad tool name from LLM
36             print("\n ...bad tool name...")
37             result = "bad tool name, retry" # instruct LLM to retry if bad
38         else:
39             result = self.tools[t['name']].invoke(t['args'])
40             results.append(ToolMessage(tool_call_id=t['id'], name=t['name'], content=str(result)))
41         print("Back to the model!")
42         return {'messages': results}

```

The *agent* is parameterized by three things: the *model* to use, the *tools* to call, and the *system message* (saved as an attribute).

First, initialize the graph with the *AgentState*.

Then, add the node that calls the LLM and the action node (the order in which they are added is not important since there are no links yet). Finally, add the conditional edge. This edge follows the LLM call node and checks if there is an action to perform. If an action is needed, it routes to the action node; if not, it routes to the end node. The *first argument* of the *conditional edge* is the *starting node*, the *second* argument is the *function* that *determines the next step*, and the *third* argument is a *dictionary* that maps the function's response to the next node. Next, add a standard edge between the action node and the LLM node, and establish the entry point of the graph as the LLM node. After completing the setup, compile the graph to transform it into a *LangChain* runnable. Save the compiled graph as a class attribute, along with the tools and model. For the tools, create a dictionary that maps each tool's name to the tool itself. For the model, call *bind\_tools*, passing in the list of tools; this informs the model of the available tools.

For the LLM node, create the *call\_openai* function: retrieve the list of messages from the agent state, add the system message, and call the model. Finally, return a dictionary with the message obtained from the LLM.

For the action node, create the *take\_action* function: if we are in this node, we know that one or more tools need to be called. The *tool\_calls* attribute will be present in the last message, which we retrieve. Loop through these tool calls, find the correct tool (using the name in the dictionary), and invoke it with the arguments from the tool call. Add a *ToolMessage* to the state with the results of the invocation and return the new messages obtained from this call. If the *tool\_name* is incorrect, print an error message, indicating to the model that the provided tool name is invalid and to retry.

For the conditional edge, create the *exists\_action* function: take the last message and return whether there are tools to call (*True* if there are, *False* if not).

This defines our agent; now we can use it:

```

1 prompt = """You are a smart research assistant. Use the search engine to look up information. \
2 You are allowed to make multiple calls (either together or in sequence). \
3 Only look up information when you are sure of what you want. \
4 If you need to look up some information before asking a follow up question, you are allowed to do that!
5 """
6
7 model = ChatOpenAI(model="gpt-3.5-turbo") #reduce inference cost
8 abot = Agent(model, [tool], system=prompt)

```

We have defined the system prompt and the model to use; for the list of tools, we pass only one, our *TavilySearch* tool.

This graph we have created can also be visualized graphically:

```

1 from IPython.display import Image
2
3 Image(abot.graph.get_graph().draw_png())

```

Now, let's call the agent:

```

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

```

The *HumanMessage* represents the user message; we add it as the sole element of the message list because the agent state takes a list of messages.

The result we obtain is a list of calls that have been made (first the tool, then back to the model); the final response is the last message, obtained by looking at the *content* attribute:

```

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

```

If we ask a more complex question:

```

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

```

The tool is called twice (once for each search, through *parallel tool calling*) and then returns to the model.

Let's try one last, even more complex example:



```
1 query = "Who won the super bowl in 2024? In what state is the winning team headquarters located? \  
2 What is the GDP of that state? Answer each question."  
3 messages = [HumanMessage(content=query)]  
4  
5 model = ChatOpenAI(model="gpt-4o") # requires more advanced model  
6 abot = Agent(model, [tool], system=prompt)  
7 result = abot.graph.invoke({"messages": messages})
```

First, it calls the tool to find out who won. Then it returns to the model. Next, it calls the tool again to find out the location of the headquarters. It returns to the model and finally makes one last call to find the state's GDP and returns to the model for the final answer.

In this case, the tool calls are *not parallelizable* because the *response* from the *first call* is *needed* to formulate the *second*.

## 4 Agentic Search Tools

This section explores the differences between agentic search and standard search and how to use the agentic search tool.

### 4.1 How an Agent Uses Search

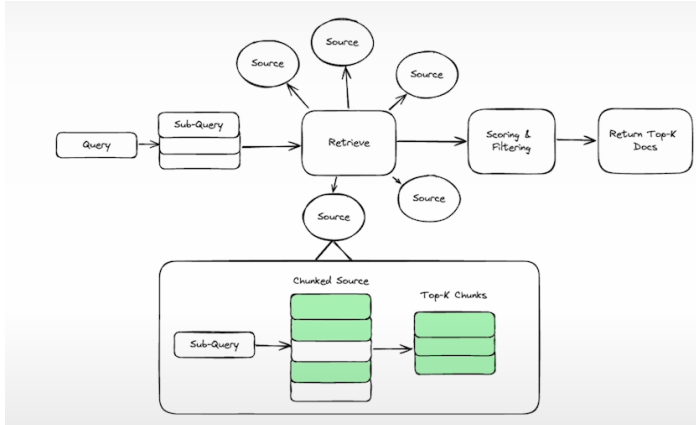


Figure 1: Example of a basic search tool implementation.

*search* to extract the top K chunks. After extracting data from the sources, the search tool would then *score the results* and filter out less relevant information.

All these steps are done automatically when, for example, we execute the following search:

```

1 from dotenv import load_dotenv
2 import os
3 from tavily import TavilyClient
4
5 # load environment variables from .env file
6 _ = load_dotenv()
7
8 # connect
9 client = TavilyClient(api_key=os.environ.get("TAVILY_API_KEY"))
10
11 # run search
12 result = client.search("What is in Nvidia's new Blackwell GPU?",
13                        include_answer=True)
14
15 # print the answer
16 result["answer"]

```

Now let's look in more detail at the difference between regular search and agentic search.

#### 4.1.1 Regular Search

Let's create a simple query about the weather in a specific location:

```

1 city = "San Francisco"
2
3 query = f"""
4     what is the current weather in {city}?
5     Should I travel there today?
6 """

```

We will test a regular search using *DuckDuckGo* and we want to get the links that contain the answer:

```

1 import requests
2 from bs4 import BeautifulSoup
3 from duckduckgo_search import DDGS
4 import re
5
6 ddg = DDGS()

```

In zero-shot learning, a *LLM* receives a prompt and produces a response based on the weights of its model.

However, this has two *limitations*:

- *Time*: If the model is not trained on recent events, such as yesterday's game results, it cannot provide an answer.
- *Source Verification*: We often want to know the sources of the obtained information, for example, to reduce hallucinations.

Let's analyze the example in Fig. 1 step by step. If the agent wants to, it sends the query to the search tool; the first step is to understand the question and *break it into sub-questions* if necessary. This way, it can handle complex queries. Then, for each sub-query, the search tool looks for the best source, choosing among many. From the best source, only the most relevant information for the sub-query is extracted. A basic implementation of this can be achieved by *chunking* the sources and then using *vector*

```

7
8 def search(query, max_results=6):
9     try:
10         results = ddg.text(query, max_results=max_results)
11         return [i["href"] for i in results]
12     except Exception as e:
13         print(f"returning previous results due to exception reaching ddg.")
14         results = [ # cover case where DDG rate limits due to high deeplearning.ai volume
15                     "https://weather.com/weather/today/1/USCA0987:1:US",
16                     "https://weather.com/weather/hourbyhour/1/54f9d8baac32496f6b5497b4bf7a277c3e2e6cc5625de69680e6169e7e38e9a8"
17                 ]
18         return results
19
20 for i in search(query):
21     print(i)

```

From the links obtained with this response, we want to extract the actual information to answer the question:

```

1 def scrape_weather_info(url):
2     """Scrape content from the given URL"""
3     if not url:
4         return "Weather information could not be found."
5
6     # fetch data
7     headers = {'User-Agent': 'Mozilla/5.0'}
8     response = requests.get(url, headers=headers)
9     if response.status_code != 200:
10         return "Failed to retrieve the webpage."
11
12     # parse result
13     soup = BeautifulSoup(response.text, 'html.parser')
14     return soup
15
16 # use DuckDuckGo to find websites and take the first result
17 url = search(query)[0]
18
19 # scrape first website
20 soup = scrape_weather_info(url)
21
22 print(f"Website: {url}\n\n")
23 print(str(soup.body)[:50000]) # limit long outputs

```

*BeautifulSoup* extracts the HTML, then we print the output; this is very long and contains a lot of unnecessary information. We want to extract the headers and some content:

```

1 # extract text
2 weather_data = []
3 for tag in soup.find_all(['h1', 'h2', 'h3', 'p']):
4     text = tag.get_text(" ", strip=True)
5     weather_data.append(text)
6
7 # combine all elements into a single string
8 weather_data = "\n".join(weather_data)
9
10 # remove all spaces from the combined text
11 weather_data = re.sub(r'\s+', ' ', weather_data)
12
13 print(f"Website: {url}\n\n")
14 print(weather_data)

```

Now the output is much better; we get readable and understandable text that answers the question, but it is still much longer than necessary.

Now let's try with agentic search.

### 4.1.2 Agentic Search

Let's perform the same query and call *Tavily* to get an answer:

```

1 # run search
2 result = client.search(query, max_results=1)
3

```

```
4  # print first result
5  data = result["results"][0]["content"]
6  print(data)
```

The first response (from the search tool) is a simple JSON with a lot of information about the weather in San Francisco.

We can further analyze it and highlight the information more clearly:

```
1  import json
2  from pygments import highlight, lexers, formatters
3
4  # parse JSON
5  parsed_json = json.loads(data.replace("'", ''))
6
7  # pretty print JSON with syntax highlighting
8  formatted_json = json.dumps(parsed_json, indent=4)
9  colorful_json = highlight(formatted_json,
10                           lexers.JsonLexer(),
11                           formatters.TerminalFormatter())
12
13 print(colorful_json)
```

We still get a JSON, but it is well-formatted. It might not be what a human wants to see, but this is the exact answer an agent would want: *structured data*.

## 5 Persistence and Streaming

When working with agents, we often need to handle longer running tasks. For these types of tasks, two important concepts are persistence and streaming.

**Persistence** allows you to maintain the state of an agent at a particular point in time. This means you can return to that state and resume from it in future interactions. This is crucial for long-running applications.

**Streaming**, on the other hand, enables you to emit a series of signals about what is happening at any given moment. This provides real-time updates on the actions the agent is taking.

### 5.1 Agent with Persistence and Streaming

First, we create the agent, as done in previous sections:

```

1  from dotenv import load_dotenv
2
3  _ = load_dotenv()
4
5  from langgraph.graph import StateGraph, END
6  from typing import TypedDict, Annotated
7  import operator
8  from langchain_core.messages import AnyMessage, SystemMessage, HumanMessage, ToolMessage
9  from langchain_openai import ChatOpenAI
10 from langchain_community.tools.tavily_search import TavilySearchResults
11
12 tool = TavilySearchResults(max_results=2)
13
14 class AgentState(TypedDict):
15     messages: Annotated[list[AnyMessage], operator.add]
16
17     from langgraph.checkpoint.sqlite import SqliteSaver
18
19 memory = SqliteSaver.from_conn_string(":memory:")
20
21 class Agent:
22     def __init__(self, model, tools, checkpointer, system=""):
23         self.system = system
24         graph = StateGraph(AgentState)
25         graph.add_node("llm", self.call_openai)
26         graph.add_node("action", self.take_action)
27         graph.add_conditional_edges("llm", self.exists_action, {True: "action", False: END})
28         graph.add_edge("action", "llm")
29         graph.set_entry_point("llm")
30         self.graph = graph.compile(checkpointer=checkpointer)
31         self.tools = {t.name: t for t in tools}
32         self.model = model.bind_tools(tools)
33
34     def call_openai(self, state: AgentState):
35         messages = state['messages']
36         if self.system:
37             messages = [SystemMessage(content=self.system)] + messages
38         message = self.model.invoke(messages)
39         return {'messages': [message]}
40
41     def exists_action(self, state: AgentState):
42         result = state['messages'][-1]
43         return len(result.tool_calls) > 0
44
45     def take_action(self, state: AgentState):
46         tool_calls = state['messages'][-1].tool_calls
47         results = []
48         for t in tool_calls:
49             print(f"Calling: {t}")
50             result = self.tools[t['name']].invoke(t['args'])
51             results.append(ToolMessage(tool_call_id=t['id'], name=t['name'], content=str(result)))
52         print("Back to the model!")
53         return {'messages': results}

```

Compared to previous implementations, we have added a new line of code. To add persistence, we rely on the concept

of a *checkpoint*, which is responsible for checkpointing the state after and between every node. To add persistence to this agent, we used a *SqLiteSaver* which uses SQLite, a built-in database under the hood. Here, we use an in-memory database; this database is very simple, and if the notebook is refreshed, it disappears. However, it can easily be connected to an external database. Alternatively, other checkpointers can use *Redis*, *Postgres*, and other more persistent databases.

Once this checkpoint is initialized, it can be used by passing it to *graph.compile*. To facilitate this, we added a new parameter to the agent: the checkpoint.

Now we can create our agent:

```
1 prompt = """You are a smart research assistant. Use the search engine to look up information. \
2 You are allowed to make multiple calls (either together or in sequence). \
3 Only look up information when you are sure of what you want. \
4 If you need to look up some information before asking a follow up question, you are allowed to do that!
5 """
6 model = ChatOpenAI(model="gpt-4o")
7 abot = Agent(model, [tool], system=prompt, checkpoint=memory)
```

To use our agent more effectively, we need to introduce the concept of streaming. There are two important aspects of streaming:

1. *Streaming individual messages*: This involves streaming both the AI message that determines which action to take and the observation message that represents the result of that action. This ensures real-time updates on the actions being taken by the AI and their outcomes.
2. *Streaming tokens*: For each token generated by the LLM call, we might want to stream the output. This provides real-time updates of the generated text, token by token.

For now, let's start by focusing only on streaming individual messages. This will allow us to observe the decision-making process of the AI in real-time, as it receives new observations and decides on subsequent actions.

Let's create the human message:

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

Add the concept of *thread config* which will be used to track the different threads inside the persistence checkpoint. This allows us to have *multiple conversations* going on at the same time, necessary for production applications where there are typically many users, and also to *resume* a *conversation* we had *previously*.

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

This thread config is simply a dictionary with a key *configurable* containing the thread ID.

Now let's call the graph; instead of using *invoke*, we use *stream*:

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

We pass it the dictionary of messages and the thread. What we get is a stream of events, which represent updates to the state over time. Since our state has only one key, the messages key, we loop through and print it.

Executing this code, we get a stream of results. First, we get an *AIMessage*, which tells us to call *Tavily*. Next, we have the *ToolMessage*, the result of the *Tavily* call. Finally, we get another *AIMessage* with the final answer from the LLM.

With this stream method, we get all intermediate results and have a clear understanding of what is happening.

Now let's call the agent with another message:

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

Using the same thread ID as before, we expect to continue the previous conversation about the weather in SF. Indeed, the first message is an *AIMessage* indicating a need to look up the weather in LA. It knows we are still talking about the weather because we have introduced the concept of persistence with the checkpoint.

Again, if we use the same thread ID, the conversation continues, and we can ask questions without having to re-establish the context:

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

However, if we change the thread ID, it is as if we change the conversation, and the model no longer knows what we are talking about:

```

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

```

### 5.1.1 Streaming Tokens

To track tokens, we use the `astream_events` method that comes on all *LangChain* and *LangGraph* objects; this is an asynchronous method, meaning we need to use an asynchronous checkpoint.

To do this, we can import *AsyncSqliteSaver* and pass it to the agent (it is very similar to the one used before but asynchronous):

```

1 from langgraph.checkpoint.aiosqlite import AsyncSqliteSaver
2
3 memory = AsyncSqliteSaver.from_conn_string(":memory:")
4 abot = Agent(model, [tool], system=prompt, checkpointer=memory)

```

We use a new thread ID to start a new conversation:

```

1 messages = [HumanMessage(content="What is the weather in SF?")]
2 thread = {"configurable": {"thread_id": "4"}}
3 async for event in abot.graph.astream_events({"messages": messages}, thread, version="v1"):
4     kind = event["event"]
5     if kind == "on_chat_model_stream":
6         content = event["data"]["chunk"].content
7         if content:
8             # Empty content in the context of OpenAI means
9             # that the model is asking for a tool to be invoked.
10            # So we only print non-empty content
11            print(content, end="|")

```

We iterate over a different type of event that represents updates from the underlying stream. We want to look for events corresponding to new tokens; this type of event is called *on\_chat\_model\_stream*. When we see this type of event, we want to view it and print its content (separated with "|").

The output of this code first shows a call to the search tool; in this case, it didn't stream out anything because there is effectively no content to stream, it is just a function call. Next, we have the final response, and in this case, we see it token by token, separated by "|", exactly as we asked to print them.

## 6 Human in the Loop

There are many situations where it is beneficial to have a human in the loop to monitor and approve an agent's actions. This can be easily implemented with *LangGraph*.

**Note:** The code in this section is running in a later version of *LangGraph*. The later version has a couple of key additions:

- Additional state information is stored to memory and displayed when using *get\_state()* or *get\_state\_history()*.
- State is additionally stored every state transition while previously it was stored at an interrupt or at the end. These change the command output slightly, but are a useful addition to the information available.

### 6.1 Setting Up the Agent

Let's start with what we already know:

```

1 from dotenv import load_dotenv
2
3 _ = load_dotenv()
4
5 from langgraph.graph import StateGraph, END
6 from typing import TypedDict, Annotated
7 import operator
8 from langchain_core.messages import AnyMessage, SystemMessage, HumanMessage, ToolMessage
9 from langchain_openai import ChatOpenAI
10 from langchain_community.tools.tavily_search import TavilySearchResults
11 from langgraph.checkpoint.sqlite import SqliteSaver
12
13 memory = SqliteSaver.from_conn_string(":memory:")

```

When creating the *AgentState*, let's make a small modification. In previous examples, we annotated the messages list with the operator *add()*, which added new messages to the existing list. For human-in-the-loop interactions, however, *we might want to replace existing messages*. To achieve this, we create a function that checks for messages with the same ID and replaces the old one with the new one; otherwise, it adds the new messages as usual.

```

1 from uuid import uuid4
2 from langchain_core.messages import AnyMessage, SystemMessage, HumanMessage, AIMessage
3
4 def reduce_messages(left: list[AnyMessage], right: list[AnyMessage]) -> list[AnyMessage]:
5     # assign ids to messages that don't have them
6     for message in right:
7         if not message.id:
8             message.id = str(uuid4())
9     # merge the new messages with the existing messages
10    merged = left.copy()
11    for message in right:
12        for i, existing in enumerate(merged):
13            # replace any existing messages with the same id
14            if existing.id == message.id:
15                merged[i] = message
16                break
17        else:
18            # append any new messages to the end
19            merged.append(message)
20    return merged
21
22 class AgentState(TypedDict):
23     messages: Annotated[list[AnyMessage], reduce_messages]

```

Next, we use the same tool we have used before:

```

1 tool = TavilySearchResults(max_results=2)

```

### 6.2 Manual Human Approval

For the agent creation, we make a slight modification:

```

1 class Agent:
2     def __init__(self, model, tools, system="", checkpointer=None):

```



```

3     self.system = system
4     graph = StateGraph(AgentState)
5     graph.add_node("llm", self.call_openai)
6     graph.add_node("action", self.take_action)
7     graph.add_conditional_edges("llm", self.exists_action, {True: "action", False: END})
8     graph.add_edge("action", "llm")
9     graph.set_entry_point("llm")
10    self.graph = graph.compile(
11        checkpointer=checkpointer,
12        interrupt_before=["action"]
13    )
14    self.tools = {t.name: t for t in tools}
15    self.model = model.bind_tools(tools)
16
17    def call_openai(self, state: AgentState):
18        messages = state['messages']
19        if self.system:
20            messages = [SystemMessage(content=self.system)] + messages
21        message = self.model.invoke(messages)
22        return {'messages': [message]}
23
24    def exists_action(self, state: AgentState):
25        print(state)
26        result = state['messages'][-1]
27        return len(result.tool_calls) > 0
28
29    def take_action(self, state: AgentState):
30        tool_calls = state['messages'][-1].tool_calls
31        results = []
32        for t in tool_calls:
33            print(f"Calling: {t}")
34            result = self.tools[t['name']].invoke(t['args'])
35            results.append(ToolMessage(tool_call_id=t['id'], name=t['name'], content=str(result)))
36        print("Back to the model!")
37        return {'messages': results}

```

When we compile the graph, besides passing the checkpointer, we also pass the *interrupt\_before* parameter with the value `["action"]`. This causes the process to *interrupt before reaching the action node*. This is done because we will add a step requiring human approval before executing the tools, ensuring they are executed correctly.

This interruption happens before executing any tool; it is also possible to specify interruptions only before certain tools and not others.

Initialize the agent and call it, also passing the thread configuration:

```

1  prompt = """You are a smart research assistant. Use the search engine to look up information. \
2  You are allowed to make multiple calls (either together or in sequence). \
3  Only look up information when you are sure of what you want. \
4  If you need to look up some information before asking a follow up question, you are allowed to do that!
5  """
6  model = ChatOpenAI(model="gpt-3.5-turbo")
7  abot = Agent(model, [tool], system=prompt, checkpointer=memory)
8
9  messages = [HumanMessage(content="Whats the weather in SF?")]
10 thread = {"configurable": {"thread_id": "1"}}
11 for event in abot.graph.stream({"messages": messages}, thread):
12     for v in event.values():
13         print(v)

```

We get an interruption after the *AIMessage* because the next step would be to call the tool.

We can view the current state of the graph for this thread:

```
1 abot.graph.get_state(thread)
```

We can also see what the next step would be, i.e., the node to be called next:

```
1 abot.graph.get_state(thread).next
```

### 6.2.1 Continue After Interrupt

To continue, simply call the stream again with the same thread configuration and pass *None* as the input:

```

1 for event in abot.graph.stream(None, thread):
2     for v in event.values():
3         print(v)

```

This way, we get the *ToolMessage* and the final response.

If we now execute:

```

1 abot.graph.get_state(thread).next

```

We get an empty response, meaning there are no more steps to take.

We can combine everything in a loop, asking the user whether to proceed:

```

1 messages = [HumanMessage("Whats the weather in LA?")]
2 thread = {"configurable": {"thread_id": "2"}}
3 for event in abot.graph.stream({"messages": messages}, thread):
4     for v in event.values():
5         print(v)
6 while abot.graph.get_state(thread).next:
7     print("\n", abot.graph.get_state(thread), "\n")
8     _input = input("proceed?")
9     if _input != "y":
10        print("aborting")
11        break
12    for event in abot.graph.stream(None, thread):
13        for v in event.values():
14            print(v)

```

## 6.3 Snapshot and State Management

As a graph is executing, a snapshot of each state is stored in memory. This snapshot includes the agent state and other useful information, such as the thread and a unique identifier for each snapshot. This allows access to specific snapshots.

There are commands for accessing memory:

- `.get_state(thread)` returns the current state if only the thread ID is provided.
- `.get_state_history(thread)` returns an iterator over all the state snapshots. This iterator can be used to access all the unique identifiers for each state.

An example usage is, given the thread ID and the unique identifier `thread_ts`, one can access a specific state and use it with the command `.invoke(None, thread, thread_ts)`; this will use that state as the starting point for the rest of the graph. This is akin to time travel. If `thread_ts` is not provided, the current state will be used as the starting point.

The `thread_ts` can be used to access a specific state, modify it, and then use `.update_state(thread, state.values)` to save it to memory. Then, using stream or invoke, the modified state will be used as the starting point.

Let's see how to implement these solutions.

### 6.3.1 Modify State

Let's see an example of modifying a state.

```

1 messages = [HumanMessage("Whats the weather in LA?")]
2 thread = {"configurable": {"thread_id": "3"}}
3 for event in abot.graph.stream({"messages": messages}, thread):
4     for v in event.values():
5         print(v)

```

If we check the state, we will see two messages: *HumanMessage* and *AIMessage*.

```

1 abot.graph.get_state(thread)

```

Suppose we want to modify the last message to check the weather in Louisiana instead of LA. First, we save the current state of the graph.

```

1 current_values = abot.graph.get_state(thread)

```

The last message will be the *AIMessage*.

```

1 current_values.values['messages'][-1]

```

We can also see the list of tool calls associated with this message.

```

1 current_values.values['messages'][-1].tool_calls

```

Let's update this tool call.

```

1 _id = current_values.values['messages'][-1].tool_calls[0]['id']
2 current_values.values['messages'][-1].tool_calls = [
3     {'name': 'tavily_search_results_json',
4       'args': {'query': 'current weather in Louisiana'},
5       'id': _id}
6 ]

```

First, we get the ID associated with this tool call. Then we update the tool calls property to be a list, which will have only one element because there is only one call. The keys are the same as before, but the arguments are different.

This does not take effect until we call `update_state`.

```

1 abot.graph.update_state(thread, current_values.values)

```

It takes the thread to know which conversation we are referring to and the new values we want to overwrite.

If we then continue, we will see that *Tavily* is called with the weather query for Louisiana, returns the information, and the model responds correctly.

```

1 for event in abot.graph.stream(None, thread):
2     for v in event.values():
3         print(v)

```

### 6.3.2 Time Travel

One important thing to note is that we're actually keeping a *running list of all these states*. So when we modify a state, we are actually creating a new state, which becomes the current state.

This allows us to go back and visit previous states. To do this, we call `get_state_history` and save all the states in a variable.

```

1 states = []
2 for state in abot.graph.get_state_history(thread):
3     print(state)
4     print('--')
5     states.append(state)

```

Suppose we want to go back to the weather search for LA; we need to find the correct state and restart from that point with the stream.

```

1 to_replay = states[-3]
2
3 for event in abot.graph.stream(None, to_replay.config):
4     for k, v in event.items():
5         print(v)

```

### 6.3.3 Go Back in Time and Edit

Another thing we can do is go back in time and then edit from there.

We can do the same as before, where we modify the state, but this time we modify a previous state instead of the last one.

```

1 _id = to_replay.values['messages'][-1].tool_calls[0]['id']
2 to_replay.values['messages'][-1].tool_calls = [{'name': 'tavily_search_results_json',
3         'args': {'query': 'current weather in LA, accuweather'},
4         'id': _id}]

```

We modified it to search for the weather information on *accuweather* instead of other sites.

We can update the state of this *to\_replay* and get back a new branch state.

```

1 branch_state = abot.graph.update_state(to_replay.config, to_replay.values)

```

Then, we call the stream from this *branch\_state*.

```

1 for event in abot.graph.stream(None, branch_state):
2     for k, v in event.items():
3         if k != "__end__":
4             print(v)

```

### 6.3.4 Add Message to a State

Another thing we can do is add messages to a state at a given time.

Suppose that, starting from *to\_replay*, instead of calling *Tavily*, we wanted to mock a response. We can do this by adding a new message to the state to mock out a response.

First, we get the ID of the tool call.

```
1 _id = to_replay.values['messages'][-1].tool_calls[0]['id']
```

Then, we create the *state\_update*, which is a list of messages. The structure is that of a *ToolMessage* to which we add the content, as if it were the actual response obtained from the tool.

```
1 state_update = {"messages": [ToolMessage(
2     tool_call_id=_id,
3     name="tavily_search_results_json",
4     content="54 degree celcius",
5 )]}
```

This is a new message, so when we add it, it will not replace any existing messages but will be added.

Now, we update the graph's state.

```
1 branch_and_add = abot.graph.update_state(
2     to_replay.config,
3     state_update,
4     as_node="action")
```

Because we're adding and pretending that an action has taken place, as opposed to modifying the existing state, we need to add this *as\_node="action"* parameter.

What we're doing is indicating that the *state\_update* is not just a modification but that we are performing this update as if we were the action node.

If we now call the stream in this new configuration, it does not perform any action (because the role of the action is taken by the message we added); instead, it calls the model and responds to the question with an *AIMessage*.

```
1 for event in abot.graph.stream(None, branch_and_add):
2     for k, v in event.items():
3         print(v)
```

## 7 Essay Writer

In this section, we will create a compact version of an AI researcher, specifically an essay writer.

### 7.1 Overview

Here is the general workflow of this agent:

- First, we create a *plan* for the essay; this will happen once upfront.
- Then, based on that plan, we will conduct some *research*; this involves calling *Tavily* and obtaining relevant documents.
- From there, we move into *generation*; following the plan and using the researched documents, we will write the essay.
- After generating the draft, we will either finish or continue.
- If we continue, we will enter a *reflection* phase where we generate a critique of the current essay.
- Based on the critique, we will conduct *additional research* with *Tavily* to gather more information.
- We will append the new information to the existing set of documents and return to the *generation* phase, continuing this cycle until certain criteria are met and we conclude.

### 7.2 Implementation

First, we import the environment variables and standard functions, and set up our in-memory SQLite checkpoints:

```
1 from dotenv import load_dotenv
2
3 _ = load_dotenv()
4
5 from langgraph.graph import StateGraph, END
6 from typing import TypedDict, Annotated, List
7 import operator
8 from langgraph.checkpoint.sqlite import SqliteSaver
9 from langchain_core.messages import AnyMessage, SystemMessage, HumanMessage, AIMessage, ChatMessage
10
11 memory = SqliteSaver.from_conn_string(":memory:")
```

Next, we define the agent state:

```
1 class AgentState(TypedDict):
2     task: str
3     plan: str
4     draft: str
5     critique: str
6     content: List[str]
7     revision_number: int
8     max_revisions: int
```

We want to keep track of several things:

- *task*: The user-provided topic for the essay.
- *plan*: The outline of the essay generated by the planning agent.
- *draft*: The current draft of the essay.
- *critique*: Feedback and recommendations from the critique agent.
- *content*: A list of documents obtained through *Tavily* research.
- *revision\_number*: The number of revisions made so far.
- *max\_revisions*: The maximum number of revisions allowed.

We initialize the model:

```
1 from langchain_openai import ChatOpenAI
2 model = ChatOpenAI(model="gpt-3.5-turbo", temperature=0)
```

Now, we create the prompts for each stage of the agent. First, the planning prompt:

```
1 PLAN_PROMPT = """You are an expert writer tasked with writing a high level outline of an essay. \
2 Write such an outline for the user provided topic. Give an outline of the essay along with any relevant notes \
3 or instructions for the sections."""
```

Next, the writer prompt, which will write the essay given the content researched:

```

1 WRITER_PROMPT = """You are an essay assistant tasked with writing excellent 5-paragraph essays.\
2 Generate the best essay possible for the user's request and the initial outline. \
3 If the user provides critique, respond with a revised version of your previous attempts. \
4 Utilize all the information below as needed:
5
6 -----
7
8 {content}"""

```

We also need a reflection prompt for critiquing the draft:

```

1 REFLECTION_PROMPT = """You are a teacher grading an essay submission. \
2 Generate critique and recommendations for the user's submission. \
3 Provide detailed recommendations, including requests for length, depth, style, etc."""

```

Next, the research plan prompt for generating queries to pass to *Tavily* after the planning step:

```

1 RESEARCH_PLAN_PROMPT = """You are a researcher charged with providing information that can \
2 be used when writing the following essay. Generate a list of search queries that will gather \
3 any relevant information. Only generate 3 queries max."""

```

Lastly, the research critique prompt for generating queries after receiving critique:

```

1 RESEARCH_CRITIQUE_PROMPT = """You are a researcher charged with providing information that can \
2 be used when making any requested revisions (as outlined below). \
3 Generate a list of search queries that will gather any relevant information. Only generate 3 queries max."""

```

To ensure we get a list of strings from the LLM, we use function calling with a *pydantic* model, that just represents the result that we want to get back:

```

1 from langchain_core.pydantic_v1 import BaseModel
2
3 class Queries(BaseModel):
4     queries: List[str]

```

We import the *Tavily client* instead of *Tavily tool*, because we're working with it in a slightly unconventional way:

```

1 from tavily import TavilyClient
2 import os
3 tavily = TavilyClient(api_key=os.environ["TAVILY_API_KEY"])

```

Now we can create the agents for each node, starting with the planning agent:

```

1 def plan_node(state: AgentState):
2     messages = [
3         SystemMessage(content=PLAN_PROMPT),
4         HumanMessage(content=state['task'])
5     ]
6     response = model.invoke(messages)
7     return {"plan": response.content}

```

This node takes the state, creates a list of messages, one of which is the planning prompt (as the system prompt), and the other is a *HumanMessage* representing the task. We pass these messages to the model, which returns a message. We take the content of that message and set it as the plan. This updates the plan key to be this response.

The next node is the agent that takes the plan and performs some research:

```

1 def research_plan_node(state: AgentState):
2     queries = model.with_structured_output(Queries).invoke([
3         SystemMessage(content=RESEARCH_PLAN_PROMPT),
4         HumanMessage(content=state['task'])
5     ])
6     content = state['content'] or []
7     for q in queries.queries:
8         response = tavily.search(query=q, max_results=2)
9         for r in response['results']:
10             content.append(r['content'])
11     return {"content": content}

```

This node takes the agent state and first generates some queries. By using *model.with\_structured\_output(Queries)*, we ensure the response follows the pydantic object that we defined, which is a list of queries. We then obtain a list of documents (and if it's empty, we create an empty list). We loop over the generated queries and search them in *Tavily*. We append the results to the content and update the content key.

Next, we create the node that writes the draft:

```

1 def generation_node(state: AgentState):
2     content = "\n\n".join(state['content'] or [])
3     user_message = HumanMessage(
4         content=f"{state['task']}\n\nHere is my plan:\n\n{state['plan']}")
5     messages = [
6         SystemMessage(
7             content=WRITER_PROMPT.format(content=content)
8         ),
9         user_message
10    ]
11    response = model.invoke(messages)
12    return {
13        "draft": response.content,
14        "revision_number": state.get("revision_number", 1) + 1
15    }

```

First, we prepare the content by joining the list of strings into one. Then, we create the user message by combining the task and the plan. We create a list of messages that includes the *SystemMessage* and the combined task and plan. We pass these messages to the model and obtain a response. We then update two things in the state: the draft (which will be the model's response) and the revision number (which will be incremented by 1).

Next, we need the reflection node:

```

1 def reflection_node(state: AgentState):
2     messages = [
3         SystemMessage(content=REFLECTION_PROMPT),
4         HumanMessage(content=state['draft'])
5     ]
6     response = model.invoke(messages)
7     return {"critique": response.content}

```

This node sets the system prompt to the reflection prompt and the newly created draft. We pass these messages to the model and take the response. Finally, we update the critique with the value of the response.

Lastly, we create the node for research after the critique:

```

1 def research_critique_node(state: AgentState):
2     queries = model.with_structured_output(Queries).invoke([
3         SystemMessage(content=RESEARCH_CRITIQUE_PROMPT),
4         HumanMessage(content=state['critique'])
5     ])
6     content = state['content'] or []
7     for q in queries.queries:
8         response = tavily.search(query=q, max_results=2)
9         for r in response['results']:
10             content.append(r['content'])
11     return {"content": content}

```

This node is very similar to the previous research node, but it takes the critique as input.

The last thing we need is the *should\_continue* condition, which checks the revision number. If it is greater than the maximum number, we end; otherwise, we continue to reflect:

```

1 def should_continue(state):
2     if state["revision_number"] > state["max_revisions"]:
3         return END
4     return "reflect"

```

Now that we have all the nodes and conditional edges, we need to assemble them into a graph. First, we initialize the graph with the agent state:

```

1 builder = StateGraph(AgentState)

```

We then add all the created nodes and set the entry point:

```

1 builder.add_node("planner", plan_node)
2 builder.add_node("generate", generation_node)
3 builder.add_node("reflect", reflection_node)
4 builder.add_node("research_plan", research_plan_node)
5 builder.add_node("research_critique", research_critique_node)
6
7 builder.set_entry_point("planner")

```

We add the conditional edge:

```

1 builder.add_conditional_edges(
2     "generate",
3     should_continue,
4     {END: END, "reflect": "reflect"}
5 )

```

And finally, we add all the standard edges:

```

1 builder.add_edge("planner", "research_plan")
2 builder.add_edge("research_plan", "generate")
3
4 builder.add_edge("reflect", "research_critique")
5 builder.add_edge("research_critique", "generate")

```

Now we can compile the graph and pass in the checkpointer we created earlier:

```

1 graph = builder.compile(checkpointer=memory)

```

For better understanding and to verify everything is set up correctly, we can visualize the created graph:

```

1 from IPython.display import Image
2
3 Image(graph.get_graph().draw_png())

```

Finally, we can test our writing agent:

```

1 thread = {"configurable": {"thread_id": "1"}}
2 for s in graph.stream({
3     'task': "what is the difference between langchain and langsmith",
4     "max_revisions": 2,
5     "revision_number": 1,
6 }, thread):
7     print(s)

```

We call *graph.stream* to observe everything that happens.



## 8 LangChain Resources

In this course, we have covered only the basics of implementing an agent. This section aims to guide you on where to find further and more in-depth information.

The first useful resource is definitely the *LangChain* documentation. It provides a great high-level overview of all the packages and services in the *LangChain* ecosystem. There are also many tutorials useful for building various applications with agents. Additionally, it offers an introduction to LangServe, which is an easy way to turn your *LangChain* application into a web server. Helping out with all of this is LangSmith: it can assist with debugging, prompt management, and also with monitoring in production, and it also has a playground.

Another useful resource is the *LangChain GitHub* repository, which has many good resources and templates to get started. The *LangGraph GitHub* repository has in-depth documentation: reference docs, tutorials, and how-to guides.

The prompt hub is a good place for inspiration and to see what other expert prompters are doing.

## 9 Conclusion

In this course, we have explored how to build agents, including some rather complex ones. In this final section, we will look at some additional agent flows that we haven't covered but are great to know about.

### 9.1 Multi-Agent Architecture

We touched on this concept when creating the essay writer agent. This architecture involves multiple agents working on the *same shared state*. These agents could simply be different prompts and LLMs, or they could each have different tools that they can call. The crucial aspect is that they all work on the same shared state, passing it from one agent to the next. This allows for a highly collaborative and integrated approach to task management.

### 9.2 Supervisor Architecture

In this architecture, there is a supervisor agent that oversees and coordinates several *sub-agents*. The supervisor determines the input for each sub-agent, which can have their own *distinct states*. These sub-agents might operate as independent graphs rather than sharing a single state. The key difference from the multi-agent framework is the emphasis on the supervisor's role in routing and coordinating the sub-agents. This setup is useful for complex workflows where different tasks require specialized handling.

### 9.3 Plan and Execute Architecture

This architecture begins with an explicit planning step. The agent formulates a plan outlining the steps needed to complete a task. It then starts *executing the plan step-by-step*. The agent may perform a step, return to update the plan if necessary, and then proceed to the next step. This cycle continues until the plan is fully executed. Afterward, the agent checks if the task was successfully completed. If not, it replans and repeats the process until the objective is met. This method ensures a structured approach to task execution and allows for adjustments as new information becomes available.

### 9.4 Language Agent Tree Search Architecture

This architecture involves a tree search over possible actions. The agent generates an action and then reflects on its consequences. Based on this reflection, it generates subsequent sub-actions and reflects further. Through these reflections, the agent can determine the *optimal path by backpropagating and updating parent nodes* with new information. This iterative process helps inform future decisions from previous states, allowing for a dynamic and flexible approach to problem-solving.