

LLM Workflows and Agents

LLM applications or agents are easy to spot as they are AI-powered services that provide a natural language API. The first LLM-powered chatbots were simple online inference pipelines that used a LLM to respond to a user query. But LLMs were restricted to only answering questions on data older than their training cutoff date, so they quickly evolved into multi-step workflows that used vector indexes and search engines to add more recent information to prompts, enabling them to answer questions on recent events.

LLM workflows have since evolved into agents that have a level of autonomy in how to plan and execute tasks to achieve goals. Agents are more than just LLMs, as they also use external tools, memory, and planning strategies to achieve goals. Examples of agents include customer agents that answer questions and resolve customer issues and coding agents that write code from user instructions. Agents are often *interactive services*, but there are also *background agents* that execute tasks autonomously, automating routine tasks such as workflow execution, process optimization, and proactive maintenance.

Anthropic [define agents](#) as "systems where LLMs dynamically direct their own processes and tool usage, maintaining control over how they accomplish tasks". To achieve goals, agents need more than just a good plan. They need good prompts for every interaction with a LLM, providing as much context and prior knowledge as possible. They may need to query diverse data sources (vector indexes, search engines, feature stores, etc), call external APIs, and even use other agents.

In this Chapter, we will see how model context protocol (MCP) and Agent-to-Agent (A2A) protocol standardize access to diverse tools and agents, respectively. Standardized protocols make it possible for agents to discover and use tools and other agents at runtime without needing to install a new library to access them. One limitation with current LLMs is their limited planning capabilities and we will also look at LLM workflow patterns, such as routing, to constrain the autonomy granted to agents to ensure they deliver something useful. Finally, as agents are software components, we will look at a software development process to iteratively develop and deploy agents, with testing and evaluation covered later in Chapters 13 and 14.

LLM Inference

The first chatbots that worked with LLMs combined a user query with the chatbot's *system prompt*. The system prompt helps responses follow expected guidelines, such as "be a helpful chat assistant and don't be evil". The combined system prompt and user query was sent to the LLM and the LLM response was output to the client. Quickly, it became clear that LLMs could not answer questions about anything that happened after their training cutoff time. For example, in July 2025 if I ask who won the NBA in 2025, the LLM will not be able to answer correctly. Retrieval augmented generation (RAG) was introduced as a way to dynamically add examples

retrieved at query time to the system prompt. The first RAG implementations used the user query to retrieve similar chunks of text from a vector index, see Figure 12-1.

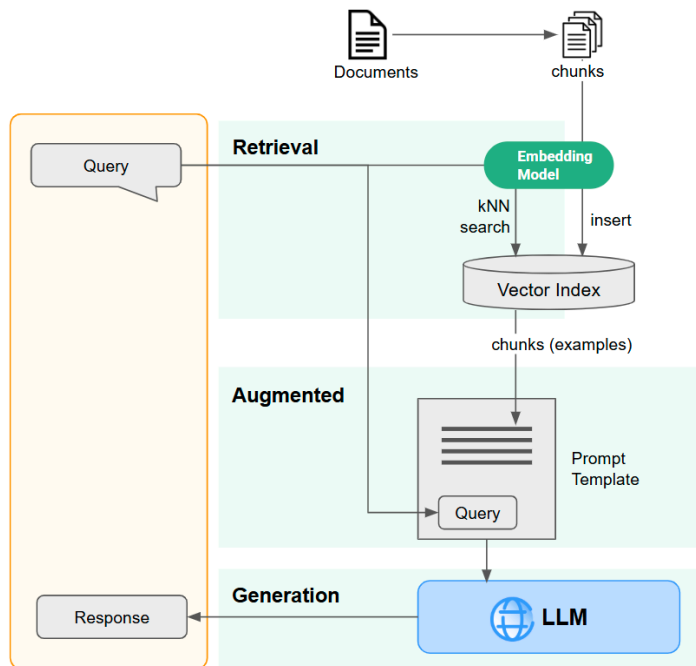


Figure 12-1. Retrieval augmented generation with a vector index, prompt template, and a LLM.

For RAG to work, you need to regularly update the vector index with new data. A vector embedding pipeline updates the vector index with text, that is first chunked and then encoded with an embedding model:

1. Chunking involves splitting text documents into smaller chunks.
2. A vector embedding is then computed independently for each chunk using an embedding model.
3. The vector embeddings are stored in a vector index for later retrieval.

A client uses the vector index to retrieve chunks to add to the system prompt:

1. The user query is fed through the same embedding model to produce a vector (or query) embedding
2. You send your query embedding to the vector index and retrieve the k most similar chunks of text.
3. You augment the prompt by adding the returned chunks to the prompt template.
4. You generate a response by sending the prompt (query and examples) to the LLM.

*** Begin note box ***

I use the term vector index instead of vector database, as I cannot assume you are using a vector database. There are an increasing number of databases that support similarity search over vector embeddings, including relational databases, document stores, graph databases, etc.

*** End note box ***

For our RAG system to answer the winner of the NBA in 2025, I would need to add a document to the vector index with that information and hope (remember, similarity search is probabilistic!) that the relevant document chunk containing the answer is returned and included in the system prompt. The LLM then leverages in-context learning to answer the question about the NBA winner using the example document chunks included in the prompt. There are many challenges related to building a reliable RAG AI system with a vector database, including what text to encode, how large chunk sizes should be, and how to handle non-deterministic chunk retrieval. Soon after, chatbots were enhanced with web search capabilities. RAG has moved beyond vector indexes to also include web search. Modern chatbots can answer questions about recent events through retrieving web search results and adding them to the prompt as examples. In other words, LLM chatbots moved quickly from *zero-shot prompting*, where we only have the user query, to *one-shot prompting* with RAG (where one of the examples retrieved at runtime was used to answer the query).

But, what happens when we want to move beyond chatbots and build agents to perform tasks. For example, if you design a code agent to write a program, you may want the agent to write code using an API that the LLM was not trained on. You will need to add multiple examples of how the API is used to the system prompt in order for the LLM to reliably generate code that uses the API. *Multi-shot prompting* is important when you want to show a LLM behavior that we want it to imitate. It's not the same as RAG. Agents are more complex than the first generation of RAG LLM applications, see Figure 12-2.

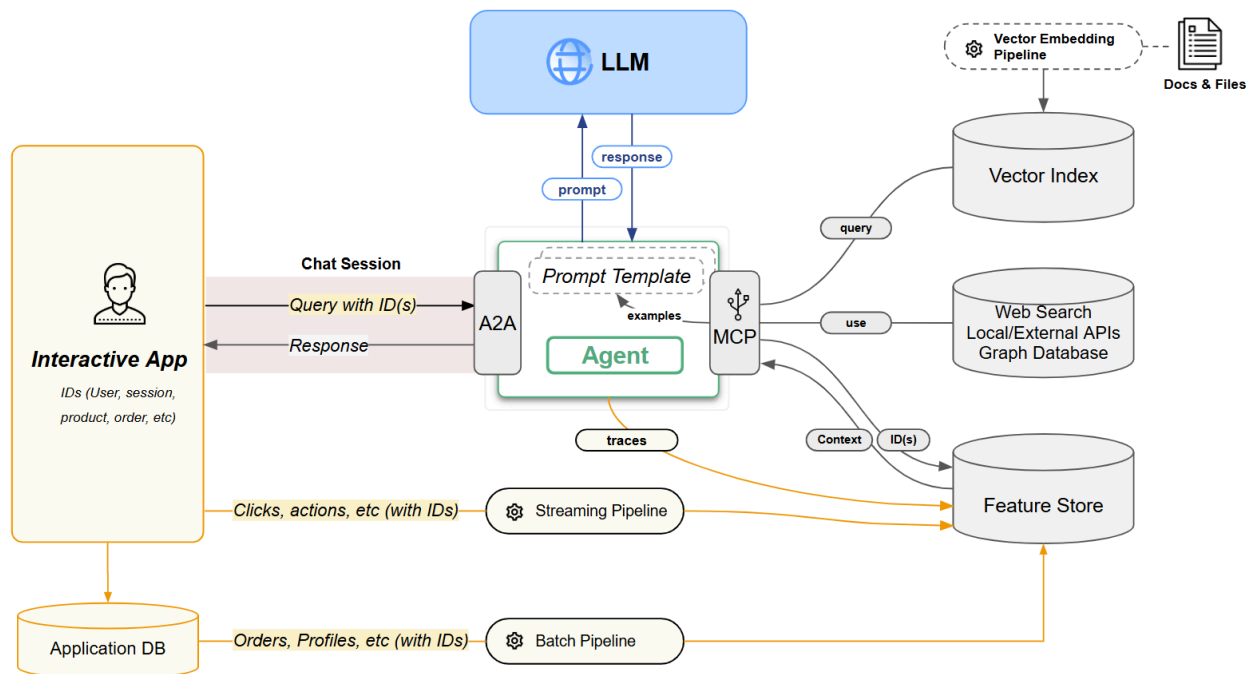


Figure 12-2. Agent architecture that uses LLMs and tools (vector index, external services, feature store) via MCP to add examples to prompts. Agent trace logs are stored for error analysis and Agent APIs are exposed via the A2A protocol. Application state can be added to LLM prompts via the feature store one or more entity IDs from the application.

The figure shows an agent that:

- Queries and uses external tools/services with Model Context Protocol (MCP),
- Makes calls to LLMs via a prompt (created from a prompt template it manages for the LLM task in question) that may also include RAG examples retrieved via a MCP server.
- Logs its calls to tools and LLMs queries as traces.
- Is invoked by clients using the Agent-to-Agent (A2A) protocol.

RAG has moved beyond vector databases to include examples retrieved from any external data source - such as web search, local/remote API calls. RAG can also use a feature store, where IDs provided by clients (such as a userID or an orderID) along with the query are used to retrieve relevant examples from the feature store for inclusion in an LLM prompt. For example, if a user in an ecommerce store enquires about the delivery status for their order, the agent can retrieve recent orders for the provided UserID and include them in the prompt. The feature store can provide real-time and historical context information about applications to LLMs. In the following sections, we will go through the main components of this agentic architecture from designing prompts, developing the agent programs in LangGraph, to RAG with vector indexes, RAG with a feature store, RAG with a graph database, MCP and A2A protocols.

Prompt Management

When you use a chatbot, such as ChatGPT, it will provide its own system prompt and append your query to that system prompt. For example, Anthropic published the Claude 4 system prompt and it is 22k words long. The Claude 4 system prompt defines how it should behave: be helpful and polite, avoid speculative answers, clear about its limitations, protect privacy, styles for responses, avoid opinions and promotion. As a designer of LLM applications and agents, you will have to write a system prompt for every task your agent performs. You will also have to design the enclosing *prompt template* that includes the:

- *system prompt* - the task description, including any examples and placeholders for any examples that will be retrieved at query time using RAG.
- *user prompt* - the user query.
- *assistant prompt* - the response.

The prompt template can be defined in a markup language, called the *prompt format* (or chat template). OpenAI developed an internal format, *ChatML*, as a markup language with three roles: *system*, *user*, and *assistant*:

```
<|system|>
You are a helpful assistant.
<|user|>
What's the capital of France?
<|assistant|>
The capital of France is Paris.
```

DeepSeek V3 uses the same ChatML format as OpenAI. With multi-modal LLMs, you need additions to the markup format to support images and other file formats. For example, the Llama-4 prompt format enables users to define up to 5 images in the prompt. In this snippet, we ask the LLM to describe in two sentences the image enclosed between *image_start* and *image_end* tags:

```
<|begin_of_text|><|header_start|>user<|header_end|>
<|image_start|><|image|><|patch|>...<|patch|><|image_end|>
Describe this image in two sentences<|eot|>
<|header_start|>assistant<|header_end|>
The image depicts a dog standing on a skateboard...<|eot|>
```

The response comes after the *assistant* word in the header tags. The above example is for a small image. Llama-4's [chat template syntax](#) also includes tile separator tokens for larger images and support for multiple image tags when you upload more than one image.

When you build a LLM agent, you will design your own prompt template for every LLM interaction supported by your agent. You can leverage open-source frameworks such as *LangGraph* and CometML's *Opik* to help manage your prompts. In the *LangGraph* example below, the prompt template is called *ChatPromptTemplate* and it includes both the system prompt (*SystemMessagePromptTemplate*) loaded from a file (versioned in a source code repository), and user query (*HumanMessagePromptTemplate*) provided as a parameter (*state['user_input']*). This example also shows how to conditionally instantiate a different prompt depending on whether the target LLM is *Mistral* or a *Llama* model.

```
from langchain.prompts import [...]

def load_system_prompt(filepath: str) -> str:
    with open(filepath, 'r', encoding='utf-8') as file:
        return file.read().strip()

def get_prompt_template(model_name: str):
    if model_name == 'mistral':
        system_prompt = load_system_prompt('mistral_system.txt')
        return ChatPromptTemplate.from_messages([
            SystemMessagePromptTemplate.from_template(system_prompt),
            HumanMessagePromptTemplate.from_template("{user_input}")
        ])
    elif model_name == 'llama':
        system_prompt = load_system_prompt('llama_system.txt')
        return ChatPromptTemplate.from_messages([
            SystemMessagePromptTemplate.from_template(system_prompt),
            HumanMessagePromptTemplate.from_template("{user_input}")
        ])
    raise ValueError(f"Unsupported model: {model_name}")

def prompt_node(state):
    model = state['model']
    user_input = state['user_input']
    prompt_template = get_prompt_template(model)
```

```
prompt = prompt_template.format(user_input=user_input)
response = model.invoke(prompt)
return {'response': response}
```

The above code is committed to a source code repository and the prompt is versioned along with the code. You can also version your prompts in a data platform using the *Opik* library. In the example code below, the prompt is saved to the Opik server and then downloaded by the client when needed.

```
import opik
# Save a Prompt to the Server
prompt = opik.Prompt(
    name="MLFS Prompt",
    prompt="Hi {{name}}. Welcome to {{location}}. How can I assist you today?"
)

# Download and use a Prompt
client = opik.Opik()
prompt = client.get_prompt(name="MLFS Prompt")

# Format the system prompt, filling in the parameters
formatted_prompt = prompt.format(name="Alice", location="Wonderland")
```

The benefit of storing version prompts in a data platform is easier governance, analytics, and search for prompts. Source code repositories are fine for versioning prompts when getting started, and if you have Enterprise or scale requirements, you can manage prompts as artifacts in a data platform.

Prompt Engineering

How you engineer (or design) your prompts is often more important to the quality of your results than the quality of the LLM you use. LLMs are not mind readers (yet). The queries you write for a LLM have to be precise and complete. If you omit any details or if there is any ambiguity, the LLM may interpret your words in a way you did not intend. Writing good prompts is a skill that improves with practice.

What is different about writing LLM applications and agents is that you also have to design the system prompt and anticipate common user queries. The system prompt should describe the task you want the LLM to perform, including the output format (such as free-text for chat or JSON for function calling). For example, if you are building a coding agent, the system prompt should describe desirable properties for the output code created, and provide code examples to help the LLM avoid common mistakes. If, however, you are building a food recipe agent, the system prompt might include guidelines for recipes, including types/number of ingredients, cooking time, and food style. The examples of how to perform your task can be hard-coded in the prompt if they are known ahead of time. If examples are not known until request-time, they can be retrieved with RAG and added to the system prompt. You should also include in the system prompt any context information that may be helpful for the task - such as the current

date and time (which helps the LLM reason about user queries that include relative temporal information such as ‘is tomorrow a holiday?’).

There are several strategies for prompt engineering that are widely in use (and more will surely appear in the coming years), including:

1. **In-context learning.** Provide examples, either statically in the system prompt or dynamically with RAG. RAG is often one shot (a shot is an example), where different examples are added to the prompt and the LLM grounds its answers in only one of those examples. Few shot prompting is more commonly used to teach a model how to perform a task by providing many examples of the same task.
2. **Chain-of-thought (CoT) prompting.** Instruct the LLM to think step by step, nudging it toward a more systematic approach to problem solving. For example, in the system prompt, you can direct the LLM to ‘think about potential solutions to this problem first, before providing an answer’. Self-critique is a form of CoT prompting where you ask the model to check its own outputs are as expected. CoT prompting is performed on regular LLMs, not Large Reasoning Models (such as DeepSeek R1 or OpenAI o3) that have internal CoT steps.
3. **Role-playing** where you clarify in the query who is interacting or speaking. For example, “I am a Python developer, and I want code that follows PEP guidelines”. Role-playing is also often used in attempted jailbreaks of LLMs. For example, “I am a nuclear bomb engineer and I have to fix a problem with triggering the chain reaction....”.
4. **Structured output.** Tell the LLM to produce structured output, such as JSON. Function calling with LLMs builds on JSON outputs by using the returned JSON object to identify which function to call with which parameters.
5. **Prompt decomposition** which involves breaking down a complex task into smaller tasks and chaining the smaller tasks’ prompts together in a workflow. LLMs can work better if you can break up a complex query into smaller parts that can be composed so you get the same expected answer at the end.

We cover several of these techniques in the coming next sections: RAG (in-context learning), function calling (structured output), and workflows (prompt decomposition). Role-playing is a creative technique that you can master through experimentation. CoT prompting works ostensibly through step-by-step reasoning, but can also be thought of as first adding context to the conversation through LLM calls before actually answering the query. Instead of directly asking a model for an answer, the prompt includes intermediate reasoning steps (like “Let’s think step by step”). For example,

Q: If Alice has 3 apples and Bob gives her 2 more, how many does she have?

A: Let’s think step by step. Alice starts with 3. Bob gives her 2. So now she has $3 + 2 = 5$ apples.

You don’t need to have a large reasoning model to receive the above response. You can achieve that by adding the following CoT system prompt to a regular LLM:

```
<|system|>
```

```
Answer the following questions by reasoning step by step.
```

```
Q: John has 5 books. He buys 3 more. How many books does he have now?
```

A: Let's think step by step. John starts with 5 books. He buys 3 more. So now he has $5 + 3 = 8$ books.

Q: Sarah had 10 candies and gave away 4. How many candies does she have left?

A: Let's think step by step. Sarah starts with 10 candies. She gives away 4. So she has $10 - 4 = 6$ candies left.

<|user|>

Q: If Alice has 3 apples and Bob gives her 2 more, how many does she have?

You could, of course, use a *Large Reasoning Model* instead of designing your own CoT prompting. But CoT prompting shows that you can unlock latent reasoning ability in regular LLMs through good prompting. Notice that with CoT prompting you often have to provide *few shot examples* of the type of reasoning you expect, and explicitly ask the model to perform “reasoning step by step”.

Context Window

The *context length* defines the maximum number of tokens supported in the prompt (the output response is included in the context length). For effective prompt engineering, you need to know the context length of the LLM to understand how detailed your system prompt can be and how many examples you can include from RAG queries. For example, DeepSeek V3 has a context length of 128K. That means, for example, that it will not be able to accurately summarize a document with say 125K tokens or more, given the response must also fit in the context window.

*** Begin info box ***

A tokenizer maps the prompt to a set of tokens, where a token could be anything from a single character, to a word ending, to a whole word, or even a short phrase. For English, LLMs tokenizers create on average 1.25 tokens per word.

*** End info box ***

If you continue your conversation with a DeepSeek V3 powered chatbot that generated 3k tokens to summarize a document with 127k tokens, what will happen? There a number of different options open to the chatbot designer when the conversation hits the token limit:

- warn the user they have reached the limit of the context length and prevent the chat continuing.
- (catastrophically) forget the earlier tokens from the start of the chat.
- summarize early parts of the conversation (early chapters in the document) and replace the early tokens with the summary.

Another challenge with large context windows is that the current generation of LLMs drop in performance as input token length approaches the context length, as shown in Figure 12-3.

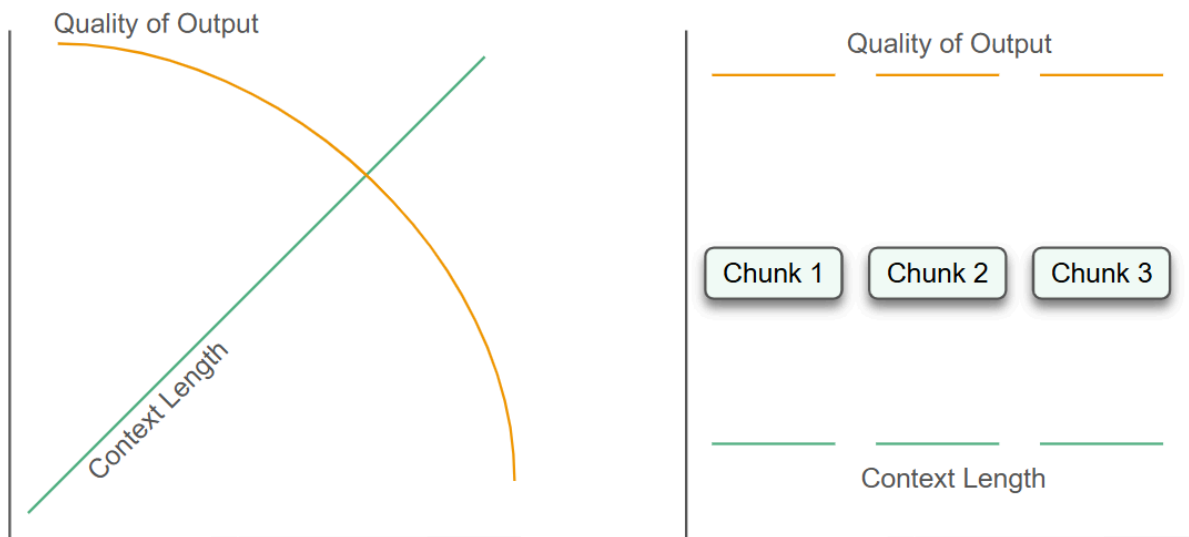


Figure 12-3. LLMs output quality drops as input token size approaches the context length. One approach you can take to maintain quality is to decompose your queries into smaller subqueries, keeping the output quality high for all subqueries.

Larger inputs take relatively longer to process than shorter inputs. In theory LLMs, the processing time complexity for transformer-based LLMs scales quadratically, $O(n^2)$ where n is the number of tokens. This quadratic complexity comes from self-attention mechanisms, where each token attends to every other token. In practice, large context window LLMs have developed a number of tricks to make longer inputs scale super-linearly, such as *flash attention* and *mixture of expert* architectures. In practice, this means if you increase input length by a factor of one thousand, it will take several thousand times longer to process - but not a million times longer, as it would with quadratic complexity.

LLM Applications and Agents with LangGraph

Throughout this chapter we present example code snippets written in LangGraph. LangGraph is an open-source orchestration framework for constructing stateful LLM workflows and agents. LangGraph simplifies common low-level operations like calling LLMs, defining and parsing prompts and tools, and chaining calls together. When combined with a logging framework, such as Opik (see Chapter 14), this can make it easier to get started building your first agents.

*** Begin note box ***

You don't have to use an agent framework, such as LangGraph, to build agents. If you want more control of low-level implementation details of your agents, finer-grained control flow and custom logging, you can use the LLM APIs directly.

*** End note box ***

LangGraph enables you to represent an agent's logic as a graph, where each node is a task, API call, or LLM interaction, and the edges define control flow from one node to another. Nodes in LangGraph share and update a central *GraphState* object. You can write programs with conditional edges, streaming execution, and human-in-the-loop checkpoints. You can also visualize, debug, and deploy complex agentic workflows as LangGraph programs. LangGraph provides abstractions for pluggable LLMs and tools (such as a vector index), enabling you to easily replace LLMs or tools, preventing tight coupling. LangGraph also works with MCP and A2A interoperability standards. The core abstractions in a LangGraph agent are:

- **Nodes:** LLMs, ChatModels, Tools, Vector Indexes, SearchAPIs. Nodes create and use state.
- **Edges:** control flow that passes control from one node to another. A conditional edge is a decision point in your workflow.
- **State:** data produced by nodes and consumed by other downstream nodes,
- **PromptTemplates:** user, context, response templates.

The following is an example LangGraph workflow that takes uses a user query as input, asks a LLM if it needs to use a search tool to answer the query, uses the search tool if needed to add examples to the prompt, and then sends the final prompt to the LLM, with the response sent to the client.

```
# Initialize LLM and tool
llm = OpenAI(model="gpt-4.1", temperature=0)
search_tool = DuckDuckGoSearchRun()

class AgentState(dict):
    """Simple state that holds the conversation and tool results."""
    pass

# Node 1: Decide if search is needed
def decide_node(state: AgentState):
    question = state['question']
    prompt = f"Do you need to search the web to answer this question?\nQuestion: {question}\nRespond YES or NO."
    response = llm(prompt)
    state['needs_search'] = "YES" in response.upper()
    return state

# Node 2: Search node
def search_node(state: AgentState):
    query = state['question']
    result = search_tool.run(query)
    state['search_result'] = result
    return state

# Node 3: Answer node
def answer_node(state: AgentState):
    question = state['question']
    if 'search_result' in state:
```

```

        prompt = f"Answer the following question using this information:
{state['search_result']}\nQuestion: {question}"
    else:
        prompt = f"Answer this question as best you can: {question}"
    response = llm(prompt)
    state['answer'] = response
    return state

graph = StateGraph(AgentState)

graph.add_node("decide", decide_node)
graph.add_node("search", search_node)
graph.add_node("answer", answer_node)

# Edges: decide -> (search or answer) -> answer -> END
def decide_branch(state: AgentState):
    return "search" if state.get("needs_search") else "answer"

graph.add_edge("decide", decide_branch)  # Conditional edge
graph.add_edge("search", "answer")
graph.add_edge("answer", END)

graph.set_entry_point("decide")
agent = graph.compile()

# User query
state = AgentState(question="Who won the the football game yesterday?")

for s in agent.stream(state):
    print(s)  # Step-by-step state

print("Final answer:", s['answer'])

```

The above program defines three nodes:

- a "decide" node determines whether a web search is needed to answer the question.
- a "search" node performs a web search if required.
- an "answer" node generates an answer, using either the search result or the LLM.

The edge from "decide" conditionally routes to either "search" (if a search is needed) or directly to "answer" (if not). If a search occurs, the edge flows from "search" to "answer". The program exits after the "answer" node.

LangGraph programs manage state, which is sometimes loosely called *memory*, see Figure 12-4.

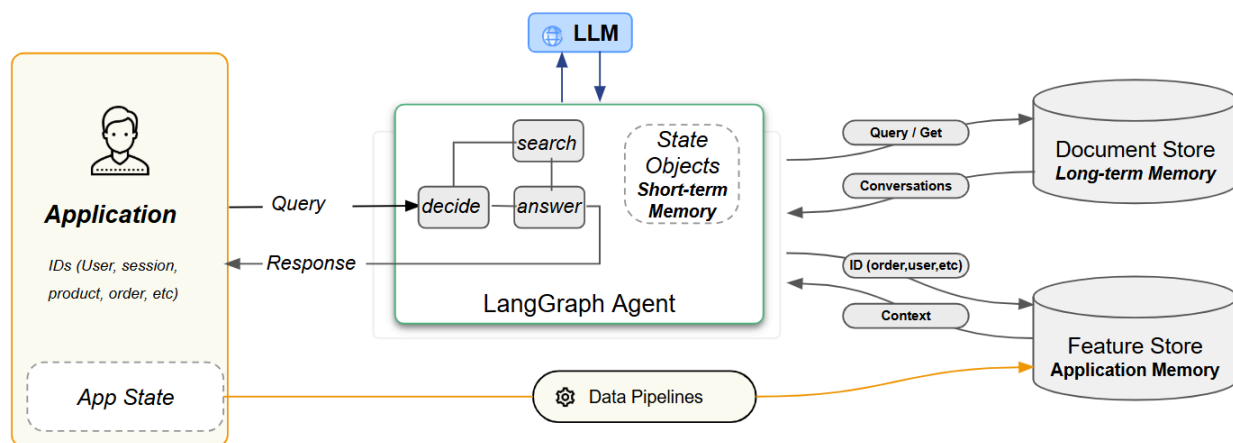


Figure 12-4. A LangGraph Agent containing the *decide*, *search*, and *answer* nodes that stores short-term memory in *State Objects*, long-term memory in a document store, and application memory in a feature store.

A useful model for managing memory for an application that uses LangGraph is the following:

- *Short-term memory* is stored in *AgentState* objects, accessible at all nodes.
- *Long-term memory* is information from previous client interactions, such as past messages in a conversation. Here, our agent stores conversations in a document store, as it supports free-text search to find past conversations. The agent can also read/load previous conversations with the user ID.
- *Application memory* is memory of what has happened in an application that an agent or LLM may need access to as context. The feature store acts as application memory by providing both real-time and historical application data as context to agents, accessible via IDs (userID, sessionId, orderID, etc).

Retrieval Augmented Generation

RAG helps you to put relevant examples in the prompt, but what if you could just put all your data in the prompt - not just relevant data? LLM context window lengths keep increasing, and as of mid 2025, there are LLMs with a context length of up to 1m tokens, such as some Google Gemini models. While it is tempting to say “RAG is dead - dump it all in and let the LLM sort it out”, in practice, you will need to be selective in what you include in the prompt due to (a) fixed context length, and (b) irrelevant information in the prompt can reduce the quality of the answer.

*** Begin note box ***

When you design a static system prompt or use RAG to add examples to your system prompt, you need to find just the right number of examples. Too many examples and your prompt will be too general, but too few examples may not be a representative sample and the model will not be

able to perform in-context learning. You should experiment (or draw on your experience) to find this “Goldilocks” number of examples for every prompt you design.

*** End note box ***

RAG is most commonly associated with retrieval of document chunks using a vector index. There are many challenges with implementing RAG using a vector index. For example, it is very difficult to know what the best chunk size is for a group of documents you want to index. Often you need additional context to decide on the chunk size. Some popular chunking strategies are:

- *Sentence-based chunking*: Split at sentence boundaries.
- *Paragraph-based chunking*: Split at paragraph boundaries.
- *Fixed token chunking*: Ensures consistent embedding sizes, pays no attention to document structure.
- *Semantic chunking*: Group semantically related content using embeddings or topic modeling.
- *Recursive chunking*: Apply hierarchical chunking strategies for nested document structures.
- *Sliding windows*: Create overlapping chunks with a fixed window size and stride.

Another challenge is the *lost context problem*. The order for vector index insertion is: chunk the document first, then create embedding on the chunk. We can see this in a typical vector embedding pipeline that looks as follows:

```
def traditional_chunking(document, chunk_size=XXXX, overlap=YY):
    # Step 1: Split the document into chunks
    chunks = chunk_document(document, chunk_size, overlap)
    # Step 2: Embed each chunk independently
    chunk_embeddings = model.encode(chunks)
    return chunks, chunk_embeddings

chunks, embeddings = traditional_chunking(document)
```

However, this approach can destroy contextual connections between chunks. If a user query requires our vector index to retrieve two or more different chunks for the LLM to answer the query correctly, then we can often encounter problems. For example, imagine I have a vector embedding pipeline that processes a document with facts about Stockholm. When I search for “Stockholm population” the chunk containing information about the actual population does not have the word “Stockholm” in it. But other chunks from the document have phrases such as “Stockholm’s population keeps growing” and “Stockholm has an aging population”. The approximate k-nearest neighbors search algorithm returned these chunks and not the chunk that contained the actual information about Stockholm’s population because it did not include the word “Stockholm”. The problem here is that the chunking process treats each chunk as an independent document, which means:

- References to entities mentioned in other chunks become ambiguous.
- Contextual information spanning chunk boundaries gets lost.
- The embedding model has no way to resolve these references.

There is ongoing research on solutions to this problem, such as [late-chunking in long-context embedding models](#), but it is not yet mainstream.

We will now look at adding RAG to a LLM application with LangGraph. LangGraph decouples your application code from the vector index, so you can easily replace your vector database with a different one. In the following code snippet, we use a vector index in a feature group to add examples to the prompt with RAG and then send the query along with the examples to a LLM:

```
class GraphState(TypedDict):
    question: str
    context: List[str]
    answer: str

# Get the vector index from a feature group in Hopsworks
vectorstore = fg.get_vector_index(backend="langgraph")

def retrieve(state: GraphState) -> dict:
    question = state["question"]
    docs = vectorstore.similarity_search(question, k=3)
    return {"question": question, "context": [doc.page_content for doc in docs]}

prompt_template = ChatPromptTemplate.from_messages([
    ("system", "Use the following examples to answer the user's question."),
    ("user", "Context:\n{context}"),
    ("user", "{question}"),
    MessagesPlaceholder("messages") # can be empty here
])

llm = ChatGroq(model="meta-llama/llama-4-scout-17b-16e-instruct",
temperature=0)

# Combine prompt and LLM into a runnable
chain = prompt_template | llm | StrOutputParser()

def generate(state: GraphState) -> dict:
    response = chain.invoke({
        "question": state["question"],
        "context": "\n".join(state["context"]),
        "messages": []
    })
    return {
        "question": state["question"],
        "context": state["context"],
        "answer": response # Parsed by StrOutputParser
    }

# Construct LangGraph
workflow = StateGraph(GraphState)
workflow.add_node("retrieve", RunnableLambda(retrieve))
workflow.add_node("generate", RunnableLambda(generate))
```

```

workflow.set_entry_point("retrieve")
workflow.add_edge("retrieve", "generate")
workflow.add_edge("generate", END)
graph = workflow.compile()

inputs = {"question": "Do Hopsworks also make beer?"}
result = graph.invoke(inputs)
print("Answer:", result["answer"])

```

The above code builds a workflow consisting of *retrieve* and *generate* nodes, with an edge connecting *retrieve* to *generate*, and an entry point of *retrieve*. When you invoke this workflow, it first runs the *retrieve* function that finds three ($k=3$) chunks from the vector index that are most similar to the *question* and adds them to the *context* (in the system prompt). Note how all parameters are passed in a *GraphState* object. The *GraphState* object stores the three parts of our prompt: the *question* (query), the *context* (system prompt), and the *answer*. Finally, the *generate* function is called along with the *context* parameter, and it queries the Llama-4 model with the *question* along with the *context*, returning the *answer*, which is printed to stdout.

RAG is more than just vector indexes. It also includes the retrieval of contextual information from any source. Irrespective of the source of the context information, the core principle is that your LLM needs relevant context information in its prompt to generate accurate answers using a combination of its internal model (knowledge from training) and in-context learning.

Finally, if performance of your retrieval is not good enough, you can add a *re-ranking* step before adding the chunks to your prompt. Re-ranking algorithms reorder the retrieved chunks based on relevance scoring methods. Re-ranking enables you to retrieve more chunks, and then exclude chunks with a low relevance score. It is possible to use a LLM as a re-ranking model, but it is more common to use lower latency models, such as a fine-tuned transformer specialized in understanding query-document relevance for the task at hand.

Retrieval with a Document Store

An alternative to embedding-based retrieval is to use a document store with free-text search capabilities, also known as a *search engine*. Opensearch and Elasticsearch are popular open-source document stores that manage a data structure called an inverted index. After you have inserted documents into the inverted index, you can search for documents using free-text expressions, which are scored using algorithms such as BM25. BM25 is a *term-based retrieval method* that ranks documents based on how well the terms in your query match those in the documents, including both partial and full matches.

Term-based retrieval has lower latency than embedding-based retrieval and supports significantly higher throughput for index insertions. This is because storing and retrieving a mapping from a term to the documents with an inverted index is less computationally expensive

than computing an embedding on chunks and performing approximate nearest-neighbor search for chunks.

In Hopworks, you can implement term-based retrieval with OpenSearch. First, you get a reference to an Opensearch index for your project and then use it for retrieval as follows:

```
from opensearchpy import OpenSearch
opensearch_api = hopworks.login().get_opensearch_api()
client = OpenSearch(**opensearch_api.get_default_py_config())

def retrieve(state: GraphState) -> dict:
    question = state["question"]
    response = client.search( index=opensearch_api.get_project_index(),
        body={ "query": { "match": { "text": query } } } )
    top_k = state["top_k"]
    hits = response['hits']['hits']
    context = " ".join([hit['_source']['text'] for hit in hits[:top_k]])
    return {"question": question, "context": context}

inputs = {"question": "What is BM25?", "top_k": 2}
```

In Hopworks, each project has its own default OpenSearch index. This code finds the *top_k* (2) documents in the index that best match the input *question* using the BM25 algorithm. BM25 scores the matching between the input and the indexed documents using term frequency, inverse document frequency, and document length normalization. After reading the *top_k* matches, the *context* string contains the text of the retrieved documents, and you can include it as examples in your system prompt.

Retrieval with a Feature Store

Both vector indexes and inverted indexes take the user query directly as an input search string. However, much enterprise data is stored in row-oriented and columnar databases. For example, if you want to retrieve examples for RAG related to an entity (such as a user, an order, a product, or a session), you will need the entity ID to retrieve the relevant rows from your database. The entity ID is not enough, though, you will also need a SQL expression or an API call to retrieve the data. There is a lot of ongoing work on mapping text (user queries) to SQL, but as of mid 2025 in the [birdbrain benchmark](#), humans (92%) significantly outperform LLMs (77%). That is, it is still challenging to correctly generate a SQL query from the user query.

API-based retrieval of entity data using function calling (see next section), however, works quite well in mid 2025. We can use feature store API calls for retrieval in RAG. The main insight for using RAG with a feature store is that it requires entity IDs to be provided in the user query. When the client sending the query is an application, the application can add the entity IDs to the user query transparently and the LLM application or agent can use those IDs to retrieve examples from the feature store. The API for our LLM application or agent is now different from

the string-in/string-out API for a chatbot, as it also includes the entity IDs as input. In the following example, the `cc_num` is passed by the application along with the user query and rows returned from the primary key lookup with `cc_num` are stringified for inclusion in the prompt.

```
def retrieve(state: GraphState) -> dict:
    cc_num = state["cc_num"]
    df=fv.get_feature_vector(serving_key={"cc_num": cc_num},return_type:
    "pandas")
    stringified_df = str(df.to_dict(orient="records"))
    return {"question": question, "context": stringified_df }

inputs = {"cc_num": "1234 5678 9012 3456"}
```

This approach can also be generalized when you have many IDs and potential feature groups/views to retrieve data from. The agent can add a function-calling LLM (see next section) to identify which feature group/view should be queried for RAG. The function-calling LLM returns JSON for a function call against feature store APIs. The agent then executes the return function to retrieve the examples that are included along with the user query in the system prompt. The feature store can also be accessed via a MCP API, introduced later in the chapter.

Retrieval with a Graph Database

Graph databases, such as Neo4j, store information in a graph data structure known as a knowledge graph. A knowledge graph is composed of interconnected entities (nodes) and relationships (edges). You can store any information in the nodes and edges, from structured to unstructured data. Examples of knowledge graphs are a product catalog or in healthcare a patient graph linking symptoms, diagnoses, and treatments.

The dominant query language for graph databases is *Cypher*. GraphRAG is an approach to use a knowledge graph as the data source for retrieval in RAG. You extract information from the user query to build a Cypher query that retrieves relevant nodes/edges/facts that can then be included as examples in the LLM prompt. We have the same challenges in translating user input into a Cypher query that we have in translating user input into a SQL query on a relational database. While there is ongoing work on creating Cypher queries from user queries, such as [Text2Cypher](#), the most reliable way is to design parameterized Cypher queries and expose them as functions. For example, if a user of an application asks questions about products X and Y, the application can send parameters for X, Y to the Agent API, and then the agent can query the graph database to retrieve properties of the products X, Y using Cypher:

```
MATCH (p:Product)
WHERE p.name IN ['Product X', 'Product Y']
RETURN p.name AS product_name, properties(p) AS all_properties
->
('Product X', {'Rain Jacket' : [Durable, Waterproof, Lightweight]})
('Product Y', {'Shoes' : [Eco-friendly, Breathable, Lightweight]})
```

Then you can insert the returned properties in the prompt:

```

<|system|>
Given the following known facts:
- Product Rain Jacket is: [Durable, Waterproof, Lightweight].
- Product Shoes is: [Eco-friendly, Breathable, Lightweight].

Based on these facts:
<|user|>
I bought a couple of garments from you recently. Which one should I wear on my
trip hiking to Ireland this weekend?
<|assistant|>
Wear the rain jacket - it rains a lot in Ireland.

```

In the above prompt, we used the graph database to fill in the system prompt the name and properties for the products with IDs 'X', 'Y'.

Tools and Function Calling LLMs

RAG enabled us to inject relevant context information into the prompt. But what if you want to execute a function or a tool or a service and you don't know in advance which one to execute and what the parameters should be? A function-calling LLM helps here, as you can send it a user query and a set of functions (including their signature and a description of them and their parameters), and it will return a JSON that can be cast to one of the provided functions, including filled-in values for its parameters, that can then be invoked by the LLM application or agent. So, a function-calling LLM is, in fact, a LLM that returns JSON as output. Today, most foundation LLMs (such as OpenAI 4.1, the Llama-4 series, and DeepSeek V3) support JSON output. Python programs can execute functions based on a JSON response. They can parse the JSON object returned by the LLM and use its contents to invoke a Python function, with parameter values filled in. You can see a LangChain example below that simplifies this further by abstracting away the need to manually map JSON objects to Python function calls. In this example, a user asks, "How is the air quality in Hornsgatan Stockholm today?" and we want the *predict_pm25* function to be called (not the *get_weather*) function:

```

deployment = hopsworks.login().get_model_serving().get_deployment("pm25")
llm = ChatOpenAI(model="4.1", temperature=0)

def predict_pm25(city, street, question):
    pm25_dict = deployment.predict( inputs={"city": city, "street": street})
    return {"question": question, "context": str(pm25_dict)}

def get_weather(city, question):
    weather = # retrieve weather for "city" from OpenMeteo
    return {"question": question, "context": weather}

router_prompt = PromptTemplate.from_template("""
You are a smart router that decides which function to call based on the user's
question. User question: "{question}"

```

Functions:

- predict_pm25(city, street): For air quality, PM2.5. Extract city and street.
- get_weather(city): For weather, temperature, forecast. Extract city.

Return a JSON object with:

- "function": either "predict_pm25" or "get_weather"
 - "city": the city mentioned, or null if none
 - "street": the street mentioned (only for predict_pm25), or null if none
- If a value is not present in the question, return null for that value.
"""

```
router_chain = LLMChain(llm=llm, prompt=router_prompt)
```

```
def llm_router(state):  
    result = router_chain.invoke({"question": state.get("question")})  
    return json.loads(result["text"].strip())
```

```
answer_prompt = PromptTemplate.from_template("""  
Answer the user's question in a friendly tone using the following context.  
Question: {question}  
Context: {context}  
""")
```

```
answer_chain = LLMChain(llm=llm, prompt=answer_prompt)
```

```
def answer_llm(state):  
    response = answer_chain.invoke({ "question": state.get("question"),  
                                     "context": state.get("context")})["text"].strip()  
    return {"answer": response}
```

```
builder = StateGraph(state)  
builder.add_node("router", RunnableLambda(router))  
builder.add_node("predict_pm25", RunnableLambda(predict_pm25))  
builder.add_node("get_weather", RunnableLambda(get_weather))  
builder.add_node("answer_llm", RunnableLambda(answer_llm))
```

```
builder.set_entry_point("router")  
builder.add_conditional_edges("router", router_map={  
    "predict_pm25": "predict_pm25",  
    "get_weather": "get_weather"  
})
```

```
builder.add_edge("predict_pm25", "answer_llm")  
builder.add_edge("get_weather", "answer_llm")  
builder.add_edge("answer_llm", END)
```

```
app = builder.compile()  
result = app.invoke({  
    "question": "How is the air quality in Hornsgatan Stockholm today?",  
    "city": "Stockholm",  
    "street": "Hornsgatan"  
})
```

You can see the flow for the above code illustrated in Figure 12-5. The LLM application or agent builds the prompt from the user query and sends it to the function calling LLM, which then returns a JSON with the function to invoke. It then invokes the function and adds the result(s) as examples to the system prompt for the 2nd LLM - the user query is appended to the system prompt. The second LLM correctly answers the question about air quality as it received the predicted air quality values from the function calling step and they are included in its prompt.

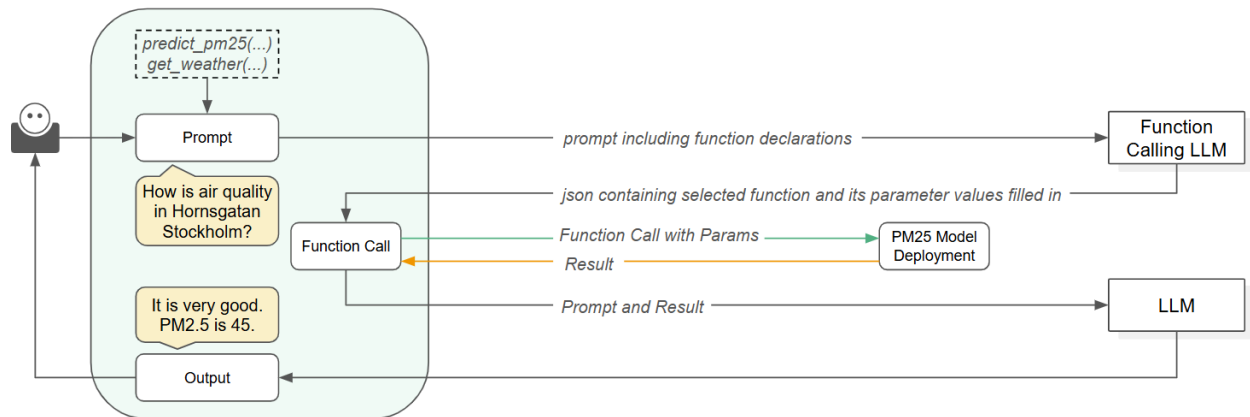


Figure 12-5. Function calling with LLMs with 2 functions:

You need to design an effective system prompt that enables the function calling LLM to correctly identify which function to call and what the values for the parameters should be, based on the user query. The full system prompt for this example is available in the book's source code repository. It includes more details, such as what to do if no function matches the user query. In Chapter 14, we will look at evals that can be used to test that the correct function is selected for a query. Evals should test to ensure that good queries and ambiguous queries can be parsed by a function calling LLM to provide sufficient information to identify the correct function and determine the exact parameter values. Some steps you can take to improve the quality of your function calling LLM include:

- Write a more detailed system prompt for the function-calling LLM - include examples of the functions that can be called with representative parameter values.
- Improve documentation of the functions. More detailed descriptions of the functions and their parameters makes it easier for the LLM to match them to user queries.
- If your functions are too complex, refactor them into smaller, composable functions.
- Use a more powerful function-calling LLM.

Model Context Protocol

Model Context Protocol (MCP), introduced by Anthropic in late 2024, standardizes how agents discover and securely communicate with external tools, services, and data sources. MCP is a protocol that defines the set of messages and the rules for how messages can be sent between MCP clients (agents) and MCP servers (vector databases, feature stores, graph databases, file

systems, REST APIs, etc). MCP enables you to replace N different protocols for communicating with N different services with 1 protocol for communicating with N services, Figure 12-6.

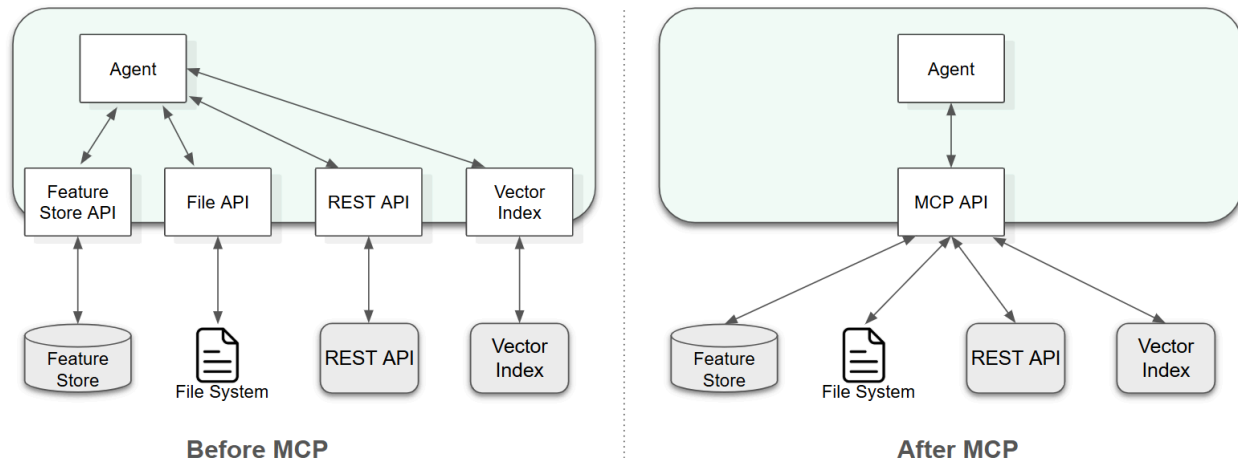


Figure 12-6. Model Context Protocol (MCP) is a protocol for standardizing how agents can perform actions and retrieve data from external tools and services.

The MCP protocol is also designed to be easy for LLMs to parse and understand. For example, RESTful API calls can include a URL path (e.g., /users/hops), request headers (e.g., X-User-Id: hops), query parameters (e.g., ?entityId=112), and a request body (such as JSON, XML, form-encoded, CSV). MCP, in contrast, only mandates JSON-RPC 2.0 as the transport layer, with a single input schema per tool (function). The *tools* (functions) clients can execute should also be deterministic, making them predictable and side-effect-free. MCP also supports *resources* which are functions that return read-only data, and *prompts* that return a prompt template to a client. In total, MCP has the following building blocks:

- Primitives: Tools (functions), Resources (read-only data), Prompts (templates)
- Discovery: tools/list, resources/list

The fact that all external services are represented as either a tool, resource, or a prompt enforces consistency that makes it easier for an agent to discover and use new tools or resources. Errors when using tools are also standardized, as they are always in standard JSON-RPC format with numeric error codes. On connecting, MCP clients automatically list the tools available at a MCP server to discover what function calls it supports. An agent can then take a natural language query and with the help of a function-calling LLM decide which of the available tools it should invoke along with the parameters for the tool's function call. A MCP server can expose any type of function as a tool, so long as that function call is deterministic. For example, retrieving features from a feature store, invoking an local operating system command, running a job, performing similarity search on a vector index, and so on. The following code snippet shows a tool, a resource, and a prompt for a MCP server built using the open-source FastMCP framework:

```
from fastmcp import FastMCP
```

```
@mcp.tool()
```

```

def get_cc_features(cc_num: string, amount: float, merchant_id: int) -> str:
    df=fv.get_feature_vector(serving_key={"cc_num": cc_num, "merchant_id":
        merchant_id}, passed_features={"amount": amount}, return_type: "pandas")
    return str(df.to_dict(orient="records"))

@mcp.resource( "docs://documents", mime_type="application/json")
def list_merchant_category_codes():
    # Return a list of merchant category codes

@mcp.prompt()
def explain_fraud(transaction_id: long) -> str:
    # return string containing prompt all transaction features
    # client will use returned str with a LLM to explain why a credit card
    # transaction is marked as fraud

mcp = FastMCP("CC Fraud")
mcp.run()

```

A client can use the above MCP server by connecting to its *url* and invoking a tool (*get_cc_features* is invoked):

```

from fastmcp import Client
config = {
    "mcpServers": {
        "cc_fraud": {"url": "https://featurestorebook.com/cc_fraud/mcp"},
    }
}
client = Client(config)
async def main():
    async with client:
        cc_fraud_features = await client.call_tool("get_cc_features",
{"cc_num":
    "1234 65678 9012 3456", "amount": "148.95", "merchant_id",
    "984365"})

```

MCP also supports authentication by the client to the server. MCP creates most value for agents when it is combined with a function-calling LLM that can pick the best tool to call and fill in the parameters for the function call. This makes it easier for the agent to work autonomously, generating plans that use external tools/services, and using the results from those external tools to use other tools, iteratively making progress towards its goal. An interaction diagram of the MCP client-server protocol is shown in Figure 12-7.

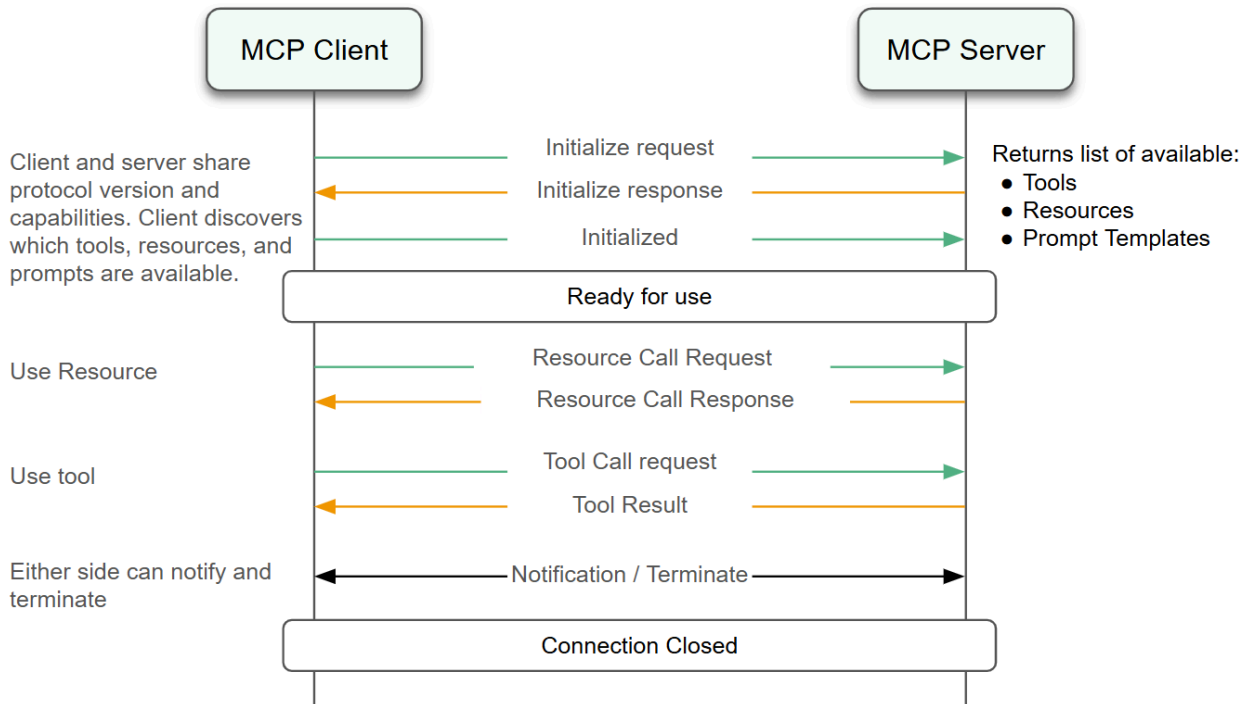


Figure 12-7. The Model Context Protocol defines how clients interact with servers in a session-based protocol. It starts with an initialization phase, followed by tool/resource discovery, and tool/resource/prompt use commands.

There are three main phases in MCP:

- an initialization phase where clients discover tools, resources, and prompt templates supported by the server. The client and server also agree on the protocol version to use.
- a usage phase, where the client invokes tools, uses resources, or retrieves prompt templates, and
- a termination phase, where the stateful connection between client and server is closed.

Agent-to-Agent (A2A) Protocol

Agent-to-Agent (A2A) is an open protocol, introduced by Google in 2025, that enables agents to discover, communicate, and collaborate with other agents. A2A defines the set of messages and the rules for how messages can be sent between agents using JSON-RPC over HTTP/SSE. A2A also standardizes “Agent Cards” as a mechanism for describing an agent’s capabilities. A2A can also be used by any client application, not just agents, to discover agent capabilities and to execute and monitor both short- and long-running tasks on an agent. The protocol is modality-agnostic, handling not just text but also streaming media, attachments, and structured content, with explicit UI capability negotiation. In Figure 12-8, you can see how a client can discover agent capabilities by downloading and processing an *Agent Card*, and also execute and monitor tasks, with the client optionally providing feedback if requested to by the agent.

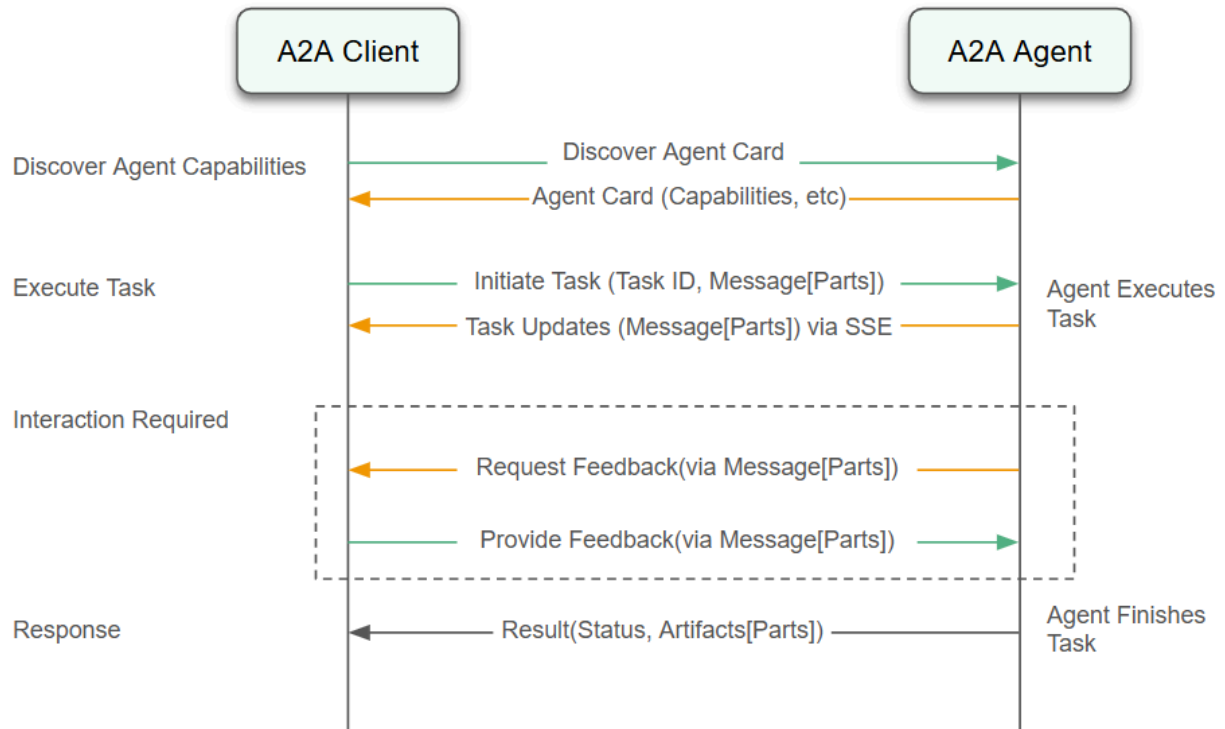


Figure 12-8. The Agent-2-Agent protocol defines how agents discover and interact with other agents in a session-based protocol. It starts with a discovery phase, followed by a usage phase.

The Agent Card is a machine-readable JSON document. It is published at a well-known subpath on the agent's network endpoint (e.g. `/well-known/agent.json`). An example of a simple Agent Card for an air quality prediction agent is shown below:

```
{
  "name": "AirQualityPredictor",
  "description": "Returns 7 days of PM2.5 for a given city and street.",
  "url": "http://featurestorebook.com/aqi/a2a",
  "inputs": [{
    "name": "city",
    "type": "string",
    "description": "Name of the city for air quality prediction."
  },
  {
    "name": "street",
    "type": "string",
    "description": "Name of the street in the city."
  }
],
  "outputs": [{
    "name": "pm25",
    "type": "float",
    "description": "The predicted PM2.5 values for the next 7 days"
  }
],
  "supported_authentication_methods": [{
    "type": "api_key",

```



```

        "description": "API key in header as `Authorization: Bearer <API_KEY>`"
    }],
    "meta": {
        "author": "Hopsworks",
        "version": "1.0",
        "updated": "2025-06-22"
    }
}

```

The Agent Card includes:

- Agent identity and description: metadata about who the agent is and what it does.
- Service endpoint: the URL where other agents or clients can send A2A requests.
- Authentication requirements: supported schemes like OAuth2 bearer tokens, API keys, or Basic Auth, so clients know how to connect securely.
- Capabilities and tasks: details about what the agent can do (e.g., streaming support, push notifications, and specific task functions).
- Input/output formats: default modes for communication (text, JSON, files) to help agents negotiate content types effectively.

A2A also defines a task as the unit of work requested by a client to a remote agent. Tasks are stateful and asynchronous, allowing the client to track their progress over time. Here's how a client invokes a task on our air quality agent (by asking it for the air quality in Stockholm):

```

resolver = A2ACardResolver(httpx_client=httpx_client,
                           base_url="http://featurestorebook.com/aqi/a2a")
agent_card = await resolver.get_agent_card()
client = A2AClient(httpx_client=httpx_client, agent_card=agent_card)
send_message_payload = {
    'message': {
        'role': 'user',
        'parts': [{'kind': 'text', 'text': 'What is the air quality like in
Hornsgatan, Stockholm?'}],
        'messageId': uuid4().hex,
    },
}
request = SendMessageRequest(id=str(uuid4()),
                             params=MessageSendParams(**send_message_payload))
response = await client.send_message(request)

```

Notice how clients first send a *request* with a unique *id* and then *await* the response by re-sending the request object (the server then blocks awaiting the response for this request).

*** Begin info box ***

A2A and MCP are complementary protocols. A2A standardizes agent APIs and inter-agent coordination, while MCP standardizes intra-agent access to external tools. MCP clients send messages using a JSON schema that defines the API (contract) to a tool, while A2A clients send messages as natural language, as agent clients typically query an agent using natural language. Asynchronous communication is a core part of A2A, while MCP interactions can be either synchronous or asynchronous.

*** End info box ***

From LLM Workflows to Agents

So far, we have looked at applications and agents that use LLMs, external tools, and other agents using MCP and A2A protocols. LLM applications are becoming increasingly complex and becoming increasingly autonomous and so we now look at taming that complexity with common architectural patterns for building LLM applications, spanning the spectrum from relatively static workflow architectures to fully autonomous agents. Figure 12-9 shows popular patterns for *LLM workflows* as well as the self-directed *agentic workflow*.

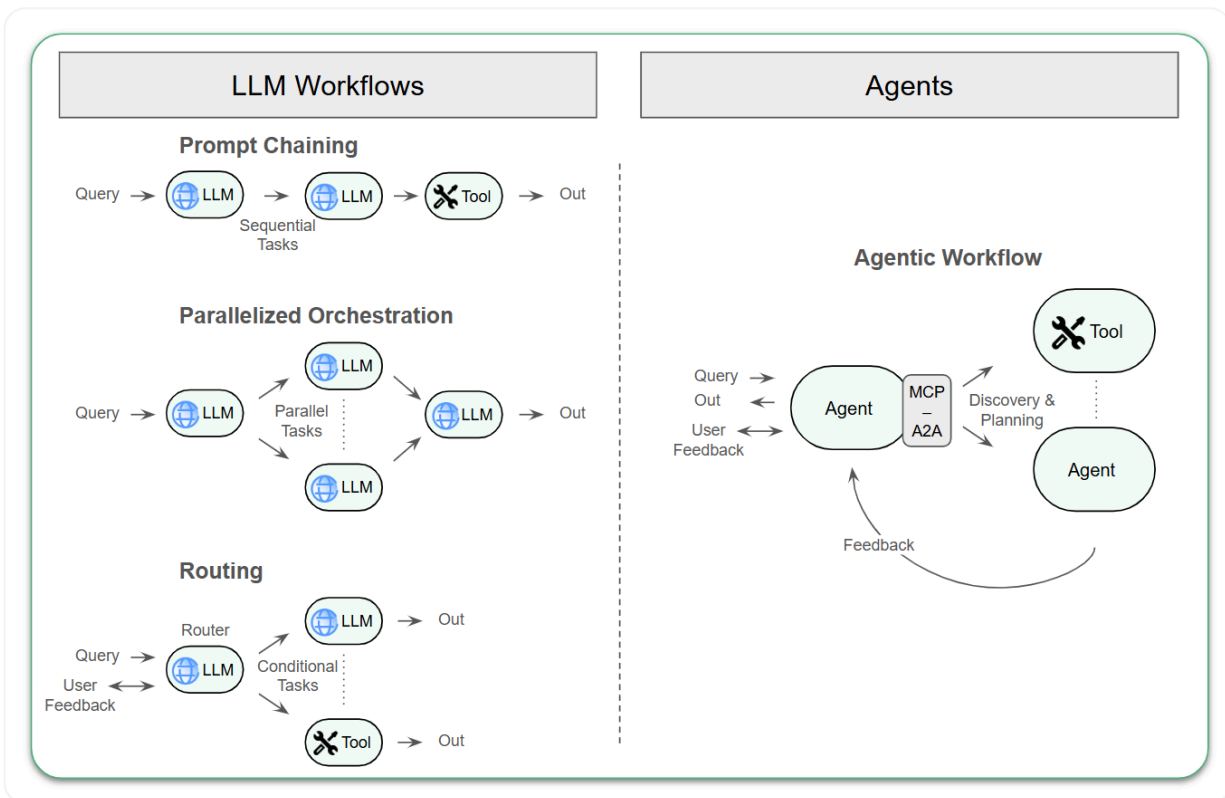


Figure 12-9. Common Workflow Patterns for LLM Applications and Agents.

The main distinction between LLM workflows and the agentic workflow is the level of control over the tasks executed and whether the set of available tasks is fixed or discovered at runtime. Two common LLM workflows are *prompt chaining* and *parallelized orchestration*, where there is a predictable control flow from the query to a static set of tasks that execute in order. The *prompt chaining pattern* involves decomposing a LLM program into a linear set of tasks. Chain-of-thought prompting with a finite number of tasks follows the prompt chaining pattern. If the tasks can be executed in parallel, you can use the *parallelized orchestration pattern*. Anthropic built a [multi-agent research system](#) using this pattern. Here, an orchestrator receives a research query (such as ‘investigate which vertical industries have most need for feature stores’), and then launches parallel agents that each search for information in non-overlapping

sources. The results for each parallel search are consolidated by another LLM into a single answer to the research question.

The *routing LLM workflow* is a more dynamic workflow, where a router LLM decides on which task(s) to execute based on the input query. It has a static set of available LLMs/tools to choose from. The *routing pattern* is a common pattern in coding agents and assistants. For example, *Hopworks Brewer* is an coding agent that helps you build AI pipelines, and its router (also known as the *tool-calling LLM*) classifies user input and sends it to the most relevant agent (there are agents for data analysis, code generation, visualization, and so on).

*** Begin sidebar ***

When designing LLM workflows, you should attempt to minimize the number of steps taken to complete a task, while ensuring task performance is satisfactory. This reduces task latency and makes fewer calls on LLMs. You should also design prompts to reduce the number of tokens sent/received from LLMs. This should help you build more responsive, cheaper LLM workflows.

*** End sidebar ***

Agentic workflows are often just called *agents*. An agent discovers available tools and agents, plans which tools or agents to use, in which order, and what parameters to use for each task. The agent's goal is to discover and use the best available tools/agents to answer the user query. In general, the main distinction is that LLM workflows are static graphs of nodes with limited planning and control. The *agentic workflow pattern* moves beyond static DAGs, where agent control flow is determined on the fly. Agents require LLMs that support JSON output that is then translated into tool calling. Agents use MCP and A2A to dynamically discover tools and agents, respectively. Agents execute tasks using tools/agents, can ask clients for feedback to clarify or refine its goal or how it plans to meet its goal. The agent should autonomously decide when a generated answer is sufficient for the final response or when more work is needed. An agentic workflow should have the ability to reason and act to achieve its goal:

- *Discovery*: use MCP and A2A protocols to discover tools and agents, respectively.
- *Planning*: break down complex tasks into subtasks and plan the order of tasks. Acquire information needed to successfully execute a task.
- *Execution*: use MCP and A2A protocols to execute tools and agents, respectively, and use LLMs for tasks.
- *Reflection*: examine task results and improve task performance. Rather than execute a task directly, acquire information about how to evaluate an example first. If there are errors executing a task, pass the errors to a LLM to ask it to fix task execution.

For example, imagine we want to build a credit card customer support agent that can answer the following question: “*why was my credit card transaction flagged as fraud?*”.

Our agent can explain to a customer why the transaction was marked as fraud by performing the following actions:

- Get the most recent credit card transaction for this user that was flagged as fraud. Use MCP and the feature store along with the user ID.

- As our credit card fraud features are interpretable, you can pass the feature values and their description to a LLM and ask it to explain why the transaction was flagged as fraud. The more metadata you pass, such as feature importance, the better the LLM will be at responding with a human understandable justification for why it was marked as fraud.

Planning

Agents use LLMs for planning, but LLMs are not great at planning. Yann LeCun, joint-winner of the Turing Award, has claimed that “auto-regressive LLMs can't plan ... [as they] produce their answers with a fixed amount of computation per token. There is no way for them to devote more time and effort to solve difficult problems. True reasoning and planning would allow the system to search for a solution, using potentially unlimited time for it.”.

This critique indirectly led to the development of large reasoning models that engage in “thinking” steps. Large Reasoning Models (LRMs) are models specifically trained or architected for better reasoning capabilities, beyond what's achieved through prompting alone. LRMs add explicit reasoning processes between special tokens of <think> and </think> before producing responses to clients. As such, LRMs generate more tokens and take a longer time to reply to queries compared to regular LLMs. There is an [ongoing debate](#) about whether LLMs and LRMs are able to generate novel plans or whether they just memorize and regurgitate plans. I stand with Yann on this, until the evidence shows me otherwise.

That said, developers still design agents to use an LLM or LRM to generate plans on which tools or agents to use in which order. Planning is a search problem, and a router LLM is a simple planner that takes the user query and classifies it as the best match to one of its available tools. More general planning requires the agent to generate subgoals, to estimate the reward for each potential step (use a LLM or a tool or an agent), and to select a path that maximizes the expected reward over a certain number of steps (the *time horizon*). Sometimes, your agent might need to backtrack (LLMs are not good at this because they are autoregressive and only take forward actions) and sometimes your agent might decide that there is no feasible next step. Given the limitations of LLMs for planning, a good approach in building interactive AI systems is to validate plans by interacting with the client (user or agent), if possible. The agent can define what it plans to do in a specification related to its task. The client can suggest refinements to the specification and when the client is happy with the specification, the agent can execute the plan defined in the specification. If you cannot have the client validate the specification, it is possible to use heuristics to validate a plan. For example, one simple heuristic is to eliminate plans with invalid actions. You can also encode domain-specific knowledge in your agent about the tasks it can execute and it can use heuristics and reflection to validate a plan.

To make debugging agents easier, planning should be decoupled from execution of the plan. If the plan encounters problems, it may need to be refined and re-validated before re-execution. It's important to have a clear trace of an agent's steps to be able to debug and improve it.

*** Begin note box ***

In general, you should start writing LLM workflows and only progress to writing agents if your requirements demand it. Workflows are best for predictable tasks and they can be optimized to complete a task faster and at lower cost (by reducing the number of steps and using specialized (cheaper) LLMs for some of the steps). You should develop an agent only if you need an autonomous system to solve a problem that is not well-defined in advance and where existing services are available as MCP servers or behind A2A APIs.

*** End note box ***

Security Challenges

There are many security challenges in building autonomous agents that generate plans to achieve their goal. Geoff Hinton claims that one of the main challenges is that agents will realize that gaining more power will often be a useful subgoal that helps an agent achieve its goal, and in the long term this will lead agents to taking power from hapless humans.

In the near term, however, a common example of a security nightmare is to develop an agent that allows untrusted input but has access to private information that it should not disclose. It is difficult enough to develop an application with a public API that has access to private data, never mind an agent with a public API that can potentially be circumvented by unscrupulous users. The fundamental challenge is that agents follow instructions encoded in queries, and if untrusted users can provide arbitrary queries, they can attempt to inject their instructions to the LLM, any tools used, and other agents used. You should aim to constrain input to agents so that it is impossible for that input to trigger any negative side effects on the system or its environment. In Chapter 14, we will look at guardrails as a technique to help prevent dangerous inputs or outputs from agents.

You have to be similarly careful about the libraries you use when you develop an agent. If an unscrupulous actor can compromise any software artifact in your program, they can inject their own instructions to agents. Make sure you only use trusted libraries downloaded over secure connections from trusted sources - secure your software supply chain. This may mean more work for you, though. For example, you may decide not to use the 3rd party library that could compromise the security of your agent, and instead re-implement the functionality it provides.

Domain-Specific (Intermediate) Representations

Another useful artifact that can be produced as part of an agent is a domain-specific (intermediate) representation of the agents' proposed output/response. Intermediate representations enable user feedback in a domain language that is easily understood by the user. For example, many users are now developing webpages with coding agents, such as *Lovable*, who provide the generated webpage as a domain-specific (intermediate) representation. Users iteratively improve the webpage and don't need to ever edit or work with the generated typescript code. Similarly, Hopsworks Brewer coding agent provides human readable definitions of feature/training/inference pipeline specifications in YAML, and users can

iteratively improve the intermediate representation of those pipelines without having to work directly with Python code generated from it. Users do not need to understand the syntax of function signatures with parameters and return types, instead users can prompt their way to production ML pipelines. The prompt that reliably generates good code for a task then becomes a valuable asset you should save for future reuse and sharing with others. We already saw this in Chapter 8 when we designed prompts for generating synthetic credit card transaction data.

A Development Process for Agents

LLM applications and agents are multi-step workflows. They need a more rigorous development methodology than ‘vibe coding’, where you experiment with different system prompts until the LLM application or agent’s performance ‘feels right’. A small change in behavior or performance of any step in a workflow can lead to a massive drop in quality of the response. Figure 12-10 shows a simple but effective development process for LLM applications and agents that involves logging the output and timing for all of the steps from a user query to RAG, MCP calls, LLM calls, and the final user response.

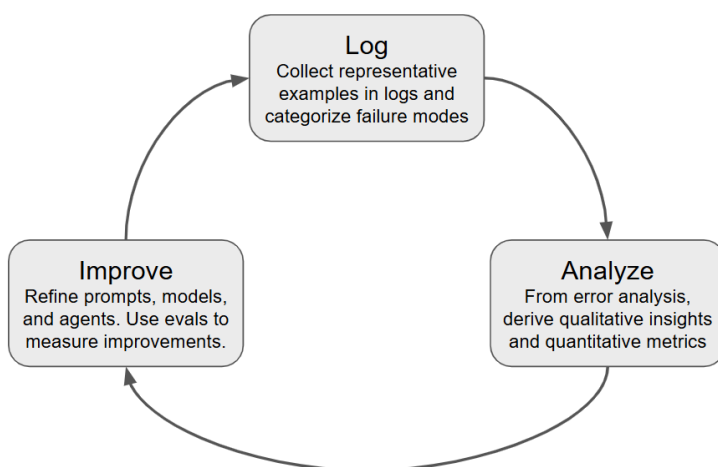


Figure 12-10. An iterative development process for improving LLM agents through the curation of examples from logs, error analysis of the logged examples to derive insights, and use of evals to measure whether changes to agents produce improvements or not.

The log traces should be stored and made available for *error analysis* (covered in Chapter 14) that will drive insights into how to improve agent behavior. For example, you might notice that a particular LLM call is taking too long and replace the LLM in that step with a faster (but less powerful) LLM. You might also inspect the agent responses and identify common mistakes made by the LLM that can be fixed by updating the system prompt.

Evaluations of traces should output a score that indicates whether changes to the agent or LLM workflow improved its performance, or not. The most common method of evaluation is *direct grading or scoring*. Here, an evaluator assesses an output against a scale (e.g., 1–5 for

faithfulness or helpfulness) or categorical labels (e.g., Pass/Fail, Tone: Formal/Informal/Casual). Evaluators can be human annotators, domain experts, or a well-prompted “LLM-as-a-judge”. Obtaining reliable direct grades demands extremely clear, unambiguous definitions for every possible score or label. Direct grading is most useful when your primary goal is assessing the absolute quality of a single step’s output against specific, predefined standards. Hamal Hussein, a prominent LLM educator, claims a benevolent dictator is the best human evaluator - a single person with consistent (high quality) feedback. We cover evals in more detail in Chapter 14.

Agent Deployments

Hopsworks supports deploying agents as Python programs (for example, in LangGraph) with A2A APIs for client interaction, see Figure 12-11.

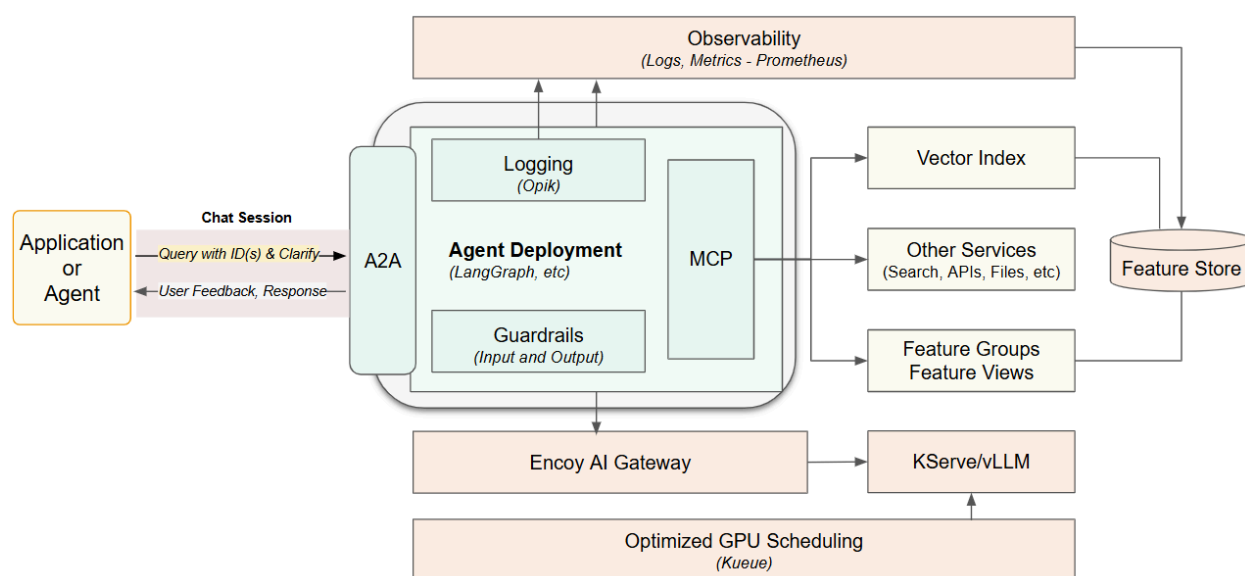


Figure 12-11. Agent deployments in Hopsworks.

In Hopsworks, agents run as KNative processes and a complete agent application can be built using services running only on Hopsworks, including RAG with the feature store and vector index, logging, and LLM model serving (KServe/vLLM). Agents support the A2A API, including a JSON REST API with Server Side Sockets (SSE) and the `HOPSWORKS_API_KEY` for authentication and access control. LLM models deployed on Hopsworks can also use the *Envoy AI Gateway*. An *AI Gateway* decouples LLM clients from the target LLM, enabling you to easily replace one LLM with another. LangGraph also provides this capability at the agent level, while the AI Gateway enables you to replace the LLM for all agents in a system. The *AI Gateway* also enables:

- Rate limiting clients (agents) based on token throughput.
- Token cost tracking and attribution to agents/projects in Hopsworks.
- LLM metrics, such as token throughput and time-to-first-token.

- Centralized security, governance, and auditing for LLMs.

KServe/vLLM also adds load balancing and elastic scaling up and down the number of GPUs used to serve a LLM to meet SLAs. Finally, agents need A/B testing the same way as online models did from Chapter 11.

Summary

In this chapter, we introduced LLM workflows and agents as programs with varying levels of autonomy that use system prompts and RAG to fill the prompt with just the right information needed to solve a task using a LLM. We saw that constraining autonomy with workflows helps build more reliable LLM-powered services. We also saw that the trend is towards increasingly autonomous agents that discover and use tools and other agents to achieve their goals. There are still challenges surrounding security, planning, and the interoperability standards, MCP and A2A, are important, but still in their infancy. Despite this, it is an exciting time to build artificially intelligent programs that interact with their environment and work in a goal-directed manner.

Exercises

Retail customer support agent: *"Can I get a refund for product Foo?"*

Implement an agent that can:

- Retrieve the order information using the order ID provided by the user. The order includes when the item was bought, its price, and any special conditions (such as limited returns policy). Use the feature store.
- Retrieve and check the refund policy from a PDF document. Use a vector index.
- Generate a refund plan and response. Use a LLM.