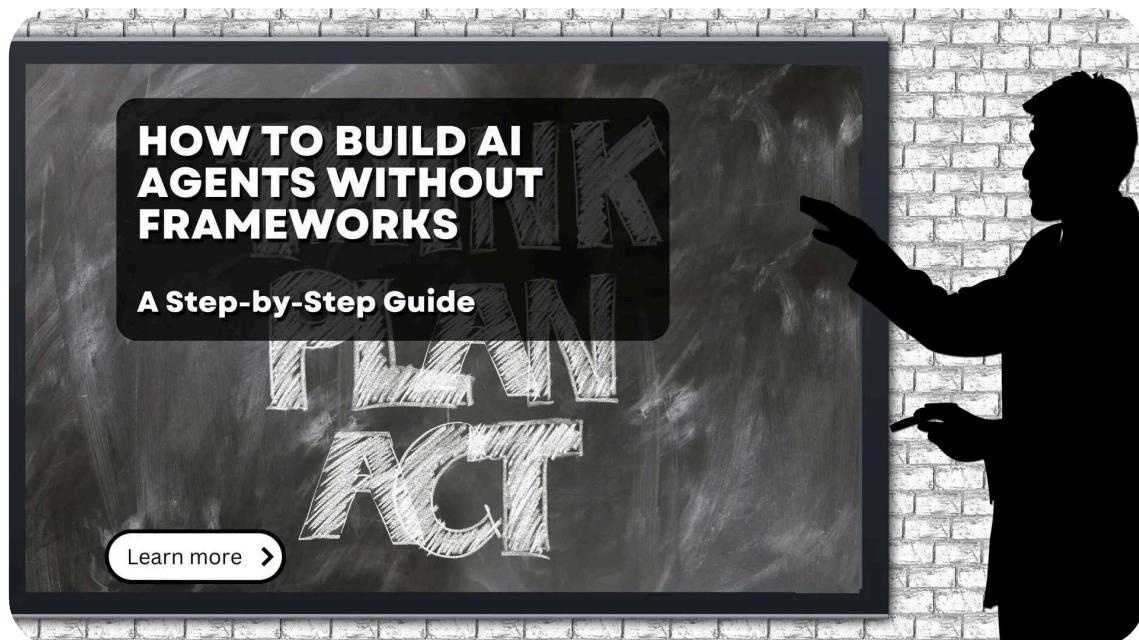


# How to Build AI Agents Without Frameworks: A Step-by-Step Guide



AI agents are often overcomplicated. At their core, they're just language models with the ability to use tools and remember context. While frameworks like [LangChain](#) or [AutoGPT](#) can help you get started quickly, they add layers of abstraction that can make it harder to understand what's actually happening – and harder to customize your agent for specific use cases.

In this tutorial, we'll build an AI agent from scratch using direct API calls to GPT-4o. You'll learn how to create an agent that can query PostgreSQL databases, fetch information from Wikipedia, and maintain conversation context through simple database storage. No frameworks, just clean code and clear concepts.

The goal is to demonstrate that AI agents don't need complex architectures. We'll focus on the essential components: prompt engineering for tool use, storing conversation history in PostgreSQL, and managing the interaction flow between the LLM and its tools.

At the end, you will know exactly, what an agent is, how it works and how to build them. If you then decide to still use a framework - that's fine. However, you will have a much better understanding of what's happening under the hood.

Let's start by understanding what makes something an AI agent, and then build one step by step.

## What Are AI Agents?

An AI agent is fundamentally just a language model that can:

- Understand what tools it has available
- Decide when to use these tools
- Remember previous interactions
- Make decisions based on all this information

## What Are Tools?

Tools (sometimes called "functions" or "capabilities") are simply functions that the agent can call to interact with the outside world. They're just Python functions with:

- A clear description of what they do
- Defined input parameters
- Structured output

Here's a concrete example of a tool:

```
1 def search_wikipedia(query: str) -> str:  
2     """
```

```

3     Searches Wikipedia and returns the first paragraph of the most relevant article
4
5     Args:
6         query: The search term to look up on Wikipedia
7
8     Returns:
9         str: The first paragraph of the most relevant Wikipedia article
10    """
11    return wikipedia.summary(query, sentences=3)
12
13 # This is how we describe the tool to the LLM
14 wikipedia_tool_description = """
15 Tool: search_wikipedia
16 Description: Searches Wikipedia and returns a summary of the topic
17 Input: A search query string
18 Example: To learn about Python programming, use: search_wikipedia("Python programming")
19 """

```

When we give our agent access to this tool, we're essentially telling it: "You can search Wikipedia by calling this function." The LLM then decides, based on the user's question, whether it needs to use this tool to provide a good answer. When we say "the LLM decides," we mean that we prompt the LLM with the available tools and ask it to decide which one to use. The AI model then simply returns a text response with the tool name.

### **Common examples of tools include:**

- Database queries
- Web searches
- API calls
- File operations
- Calculator functions
- Email sending capabilities

To bring both the concepts of agents and tools together, we see that there is no magic involved. It's just AI model interactions and a bunch of python functions.

Bringing this together, a pseudo-code for an AI agent might look like this:

```

1 # This is what an agent interaction basically looks like
2 def agent_interaction(user_input, available_tools, conversation_history):
3     # 1. We tell the LLM what it can do
4     prompt = f"""You have access to these tools:
5         - query_database: Executes SQL queries
6         - search_wikipedia: Searches Wikipedia articles
7
8     Previous conversation:
9     {conversation_history}
10
11     User question: {user_input}
12
13     Decide if you need to use a tool to answer. If yes, specify which tool and
14
15     # 2. The LLM decides what to do
16     llm_response = gpt4_call(prompt)
17
18     # 3. We execute the tool if needed and get the final answer
19     if "use query_database" in llm_response:
20         # Execute database query and get result
21         data = execute_query(extract_query(llm_response))
22         return get_final_answer(data)
23
24     return llm_response

```



This is the core concept of an agent. Everything else – complex planning, multiple tool calls, sophisticated memory systems – is just an extension of this basic pattern.

### **Common misconceptions about agents:**

- They don't need complex orchestration frameworks
- They don't need to run autonomously (though they can)
- They don't need sophisticated planning algorithms (though these can help for complex tasks)

**The key to building effective agents is not in the framework you use, but in:**

- Writing clear prompts that help the LLM understand its capabilities
- Implementing reliable tools that the agent can use

- Maintaining relevant context through some form of memory storage

In the following sections, we'll implement this pattern with real code, building an agent that can actually query databases and search Wikipedia. You'll see that the complexity comes from the specific requirements of your use case, not from the agent architecture itself.

## Why Build Agents Without Frameworks?

Frameworks like LangChain and AutoGPT are often recommended for building AI agents. While they can be useful for quick prototypes, there are compelling reasons to avoid them:

### Understanding What Really Happens

```

1 # With a framework:
2 agent = Agent.from_toolkit(SQLToolkit, WikiToolkit)
3 result = agent.run("Find users who like pizza")
4
5 # Without a framework - what's actually happening:
6 def process_query(user_input: str):
7     prompt = f"""You have these tools:
8     1. SQL queries: Execute database queries
9     2. Wikipedia search: Look up information
10
11     User question: {user_input}
12
13     Think step by step:
14     1. Do I need any tools to answer this?
15     2. If yes, which tool and how should I use it?
16     """
17
18     response = gpt4_call(prompt)
19     if "SQL" in response:
20         query = extract_sql(response)
21         data = execute_query(query)
22     if "Wikipedia" in response:
23         data = search_wikipedia(topic)
24     return get_final_answer(data)

```

The framework-free version is more explicit. You can see exactly:

- What the LLM is being told
- When tools are being called
- How decisions are being made

## Flexibility and Control

Frameworks often make assumptions about:

- How tools should be structured
- What the prompt format should be
- How memory should work

When you need to customize any of these, you often end up fighting against the framework rather than working with it.

## Simpler Debugging

When something goes wrong with a framework:

```
1 # Framework error:  
2 AgentError: Tool execution failed: SQLToolkit.query failed with error:  
3 Invalid chain configuration...  
4  
5 # Without framework:  
6 try:  
7     data = execute_query(query)  
8 except SQLException as e:  
9     print(f"SQL error: {e}") # Clear what went wrong
```

**Note:** This is more or less the most prominent reason why you should build agents without frameworks. Frameworks abstract away the core of any AI Agent - the LLM interaction. Without adding additional infrastructure, you don't know the exact prompts used, the specific LLM answers replied and the tools used. This makes debugging a nightmare.

## Lower Learning Overhead

Building agents requires understanding:

- LLM prompting
- Tool integration
- Memory management

Frameworks add another layer:

- Framework-specific concepts
- Configuration options
- Integration patterns

Why learn both when you only need the fundamentals?

## Performance and Resources

Frameworks often:

- Load unnecessary components
- Add network calls
- Include unused features

Direct implementation lets you:

- Load only what you need
- Optimize critical paths
- Control resource usage

**Note:** Arguably, raw performance is not the most important thing for AI agents, as the AI model interaction in itself is quite slow. So loading some additional modules does not add up that much.

## When to Use Frameworks

Frameworks make sense when:

- You're building a quick prototype
- You need pre-built tools
- You're learning agent concepts initially

But for production systems or when you need full control, direct implementation is often better.

Let's move on to building our agent from scratch, where you'll see how straightforward the direct implementation can be.

## Building an AI Agent Step by Step

Before we dive into the code, let's understand what we'll be building in this first step. We'll start by defining our tools - the capabilities we want to give our agent.

### Step 1: Defining Tools

In this step, we'll create two tools:

1. A PostgreSQL query tool that allows the agent to query your database and search for information stored in the database.
2. A Wikipedia search tool that lets the agent look up information in Wikipedia.

We can see that both tools are quite different, so depending on the user's question, the agent should be able to decide to use one of these tools.

When defining tools for OpenAI's function calling, we need three components for each tool:

- A schema that tells the LLM what the tool does and how to use it. This is basically just a json description of the python function. It will get clearer once you see the example.
- The actual implementation of the tool (the Python function)
- Error handling to ensure our agent behaves gracefully when things go wrong

The schema is particularly important because it acts as documentation for the LLM. Think of it as writing API documentation - the better you describe your tools, the better the LLM will use them. Important are the `name` (which should match the function name), `description` (what the tool does) and `parameters` (what function parameter values we need).

For our database tool, we'll:

- Allow only SELECT queries for safety
  - Return results as JSON for easy parsing
  - Include error handling for database connection issues

For the Wikipedia tool, we'll:

- Limit results to three sentences to keep responses concise
  - Handle disambiguation pages and missing articles
  - Return structured error messages when searches fail

Here's how we implement these tools:

```
1 import json
2 import psycopg2
3 import wikipedia
4 from typing import Dict, List, Any
5 from openai import OpenAI
6
7 # First, we define how the LLM should understand our tools
8 tools = [
9     {
10         "type": "function",
11         "function": {
12             "name": "query_database",
13             "description": """Execute a PostgreSQL SELECT query and return the results. Only SELECT queries are allowed for security reasons. Returns data in JSON format.""",
14             "parameters": {
15                 "type": "object",
16                 "properties": {
17                     "query": {
18                         "type": "string",
19                         "description": """The SQL SELECT query to execute. Must start with SELECT and should not be a multi-select query."}}},
20             }
21         }
22     }
```

```

23                                     contain any data modification command:
24                         }
25                     },
26                     "required": ["query"],
27                     "additionalProperties": False
28                 },
29                 "strict": True
30             }
31         },
32     {
33         "type": "function",
34         "function": {
35             "name": "search_wikipedia",
36             "description": """Search Wikipedia and return a concise summary.
37                             Returns the first three sentences of the most
38                             relevant article."""
39             "parameters": {
40                 "type": "object",
41                 "properties": {
42                     "query": {
43                         "type": "string",
44                         "description": "The topic to search for on Wikipedia"
45                     }
46                 },
47                 "required": ["query"],
48                 "additionalProperties": False
49             },
50             "strict": True
51         }
52     }
53 ]
54
55 # Now implement the actual tool functions
56 def query_database(query: str) -> str:
57     """
58     Execute a PostgreSQL query and return results as a JSON string.
59     Includes safety checks and error handling.
60     """
61     # Security check: Only allow SELECT queries
62     if not query.lower().strip().startswith('select'):
63         return json.dumps({
64             "error": "Only SELECT queries are allowed for security reasons."
65         })
66
67     try:
68         # You should use environment variables for these in production
69         conn = psycopg2.connect(
70             "dbname=your_db user=your_user password=your_pass host=localhost"
71         )

```

```
72
73     with conn.cursor() as cur:
74         cur.execute(query)
75
76         # Get column names from cursor description
77         columns = [desc[0] for desc in cur.description]
78
79         # Fetch results and convert to list of dictionaries
80         results = []
81         for row in cur.fetchall():
82             results.append(dict(zip(columns, row)))
83
84         return json.dumps({
85             "success": True,
86             "data": str(results),
87             "row_count": len(results)
88         })
89
90     except psycopg2.Error as e:
91         return json.dumps({
92             "error": f"Database error: {str(e)}"
93         })
94     except Exception as e:
95         return json.dumps({
96             "error": f"Unexpected error: {str(e)}"
97         })
98     finally:
99         if 'conn' in locals():
100             conn.close()
101
102 def search_wikipedia(query: str) -> str:
103     """
104     Search Wikipedia and return a concise summary.
105     Handles disambiguation and missing pages gracefully.
106     """
107     try:
108         # Try to get the most relevant page summary
109         summary = wikipedia.summary(
110             query,
111             sentences=3,
112             auto_suggest=True,
113             redirect=True
114         )
115
116         return json.dumps({
117             "success": True,
118             "summary": summary,
119             "url": wikipedia.page(query).url
120         })
121
```

```

121
122     except wikipedia.DisambiguationError as e:
123         # Handle multiple matching pages
124         return json.dumps({
125             "error": "Disambiguation error",
126             "options": e.options[:5], # List first 5 options
127             "message": "Topic is ambiguous. Please be more specific."
128         })
129     except wikipedia.PageError:
130         return json.dumps({
131             "error": "Page not found",
132             "message": f"No Wikipedia article found for: {query}"
133         })
134     except Exception as e:
135         return json.dumps({
136             "error": "Unexpected error",
137             "message": str(e)
138         })

```

## Step 2: Creating the Agent Class - Detailed Design

Before implementing the code, let's understand exactly what we're building and how it works:

### Core Concept

An AI agent is essentially a loop of:

1. Receiving user input
2. Deciding whether to use tools
3. Using tools if needed
4. Formulating a response

### Key Components Design

#### 1. Message Management

- We need to maintain a conversation history as a list of messages
- Each message has a specific role: 'user', 'assistant', or 'tool'
- The history provides context for the LLM to make better decisions

- Think of it like a chat application's state management

## 2. Tool Execution System

- Works like a command dispatcher:
  1. LLM decides to use a tool
  2. We receive a structured function call
  3. We map this to our actual Python functions
  4. We handle the execution and any errors
  5. We format the result for the LLM

## 3. Main Processing Loop

```

1 User Input → LLM Decision → [Tool Use?] → Final Response
2                                         ↓
3                                         Execute Tool(s)
4                                         ↓
5                                         Feed Results back to LLM

```

## Error Handling Strategy

We need three layers of error handling:

1. Tool-level errors (e.g., database connection failed)
2. LLM-level errors (e.g., API issues)
3. General processing errors (e.g., JSON parsing)

## Message Flow Example

```

1 User: "How many users are in our database?"
2 |
3 |→ LLM Thinks: "I need to query the database"
4 |   |
5 |   |→ Calls query_database tool
6 |   |   |
7 |   |   |→ Returns user count
8 |   |
9 |   |→ Formulates final response with data
10 |

```

```
11 ↳ Final Response: "There are 1,234 users in the database"
```

This design allows for:

- Clear separation of concerns
- Easy debugging and monitoring
- Flexible extension with new tools
- Robust error handling
- Maintainable codebase

Note that we implement the general agent structure as part of a loop. Meaning that we send back the results of our tool calls to the LLM which then decides whether to formulate a final answer or call another tool.

## Implementation

```
1 from typing import List, Dict, Optional, Any
2 from openai import OpenAI
3 import json
4
5 class Agent:
6     def __init__(self, system_prompt: Optional[str] = None):
7         """
8             Initialize an AI Agent with optional system prompt.
9
10            Args:
11                system_prompt: Initial instructions for the AI
12            """
13            # Initialize OpenAI client - expects OPENAI_API_KEY in environment
14            self.client = OpenAI()
15
16            # Initialize conversation history
17            self.messages = []
18
19            # Set up system prompt if provided, otherwise use default
20            default_prompt = """You are a helpful AI assistant with access to a database and Wikipedia. Follow these rules:
21            1. When asked about data, always check the database first
22            2. For general knowledge questions, use Wikipedia
23            3. If you're unsure about data, query the database to verify
24            4. Always mention your source of information
25            5. If a tool returns an error, explain the error to the user clearly
26
27            """
```

```
28
29         self.messages.append({
30             "role": "system",
31             "content": system_prompt or default_prompt
32         })
33
34     def execute_tool(self, tool_call: Any) -> str:
35         """
36             Execute a tool based on the LLM's decision.
37
38             Args:
39                 tool_call: The function call object from OpenAI's API
40
41             Returns:
42                 str: JSON-formatted result of the tool execution
43         """
44
45         try:
46             function_name = tool_call.function.name
47             function_args = json.loads(tool_call.function.arguments)
48
49             # Execute the appropriate tool. Add more here as needed.
50             if function_name == "query_database":
51                 result = query_database(function_args["query"])
52             elif function_name == "search_wikipedia":
53                 result = search_wikipedia(function_args["query"])
54             else:
55                 result = json.dumps({
56                     "error": f"Unknown tool: {function_name}"
57                 })
58
59             return result
60
61         except json.JSONDecodeError:
62             return json.dumps({
63                 "error": "Failed to parse tool arguments"
64             })
65         except Exception as e:
66             return json.dumps({
67                 "error": f"Tool execution failed: {str(e)}"
68             })
69
70     def process_query(self, user_input: str) -> str:
71         """
72             Process a user query through the AI agent.
73
74             Args:
75                 user_input: The user's question or command
76
77             Returns:
```

```
77     str: The agent's response
78     """
79     # Add user input to conversation history
80     self.messages.append({
81         "role": "user",
82         "content": user_input
83     })
84
85     try:
86         max_iterations = 5
87         current_iteration = 0
88
89         while current_iteration < max_iterations: # Limit to 5 iteration:
90             current_iteration += 1
91             completion = self.client.chat.completions.create(
92                 model="gpt-4o",
93                 messages=self.messages,
94                 tools=tools, # Global tools list from Step 1
95                 tool_choice="auto" # Let the model decide when to use tools
96             )
97
98             response_message = completion.choices[0].message
99
100            # If no tool calls, we're done
101            if not response_message.tool_calls:
102                self.messages.append(response_message)
103                return response_message.content
104
105            # Add the model's thinking to conversation history
106            self.messages.append(response_message)
107
108            # Process all tool calls
109            for tool_call in response_message.tool_calls:
110                try:
111                    print("Tool call:", tool_call)
112                    result = self.execute_tool(tool_call)
113                    print("Tool executed.....")
114                except Exception as e:
115                    print("Execution failed.....")
116                    result = json.dumps({
117                        "error": f"Tool execution failed: {str(e)}"
118                    })
119
120                    print(f"Tool result custom: {result}")
121
122                    self.messages.append({
123                        "role": "tool",
124                        "tool_call_id": tool_call.id,
125                        "content": str(result)
```

```

126         })
127         print("Messages:", self.messages)
128
129     # If we've reached max iterations, return a message indicating th:
130     max_iterations_message = {
131         "role": "assistant",
132         "content": "I've reached the maximum number of tool calls (5)
133     }
134     self.messages.append(max_iterations_message)
135     return max_iterations_message["content"]
136
137 except Exception as e:
138     error_message = f"Error processing query: {str(e)}"
139     self.messages.append({
140         "role": "assistant",
141         "content": error_message
142     })
143     return error_message
144
145 def get_conversation_history(self) -> List[Dict[str, str]]:
146     """
147     Get the current conversation history.
148
149     Returns:
150         List[Dict[str, str]]: The conversation history
151     """
152     return self.messages

```

## (Optional) Step 3: Adding Memory to Our Agent

**Note:** Please consider whether you need memory at all. Does your application really need to store and retrieve conversation histories? More often than not, the answer is no! But if you do need memory, here's how you can add it.

Memory in AI agents is often overcomplicated. Let's understand what we're actually trying to achieve:

### What is Agent Memory?

At its core, agent memory is just:

1. Storing conversations and their context
2. Retrieving relevant information when needed
3. Summarizing past interactions when they become too long

## Memory Types We'll Implement

### 1. Conversation History

- Store: User inputs, agent responses, tool usage results
- Purpose: Keep track of the current conversation flow

### 2. Summary Memory

- What: Periodic summaries of longer conversations
- Why: Prevent context window overflow
- How: Create summaries after X messages or on demand

## Database Schema

```
1 CREATE TABLE conversations (
2     id SERIAL PRIMARY KEY,
3     session_id UUID,
4     user_input TEXT,
5     agent_response TEXT,
6     tool_calls JSONB,    -- Store tool usage
7     timestamp TIMESTAMPTZ DEFAULT NOW()
8 );
9
10 CREATE TABLE conversation_summaries (
11     id SERIAL PRIMARY KEY,
12     session_id UUID,
13     summary TEXT,
14     start_time TIMESTAMPTZ,
15     end_time TIMESTAMPTZ,
16     message_count INTEGER
17 );
```

## Memory Flow

1 User Input → Store in DB —————

```

2 |           |
3 |           ▼
4 |             Check Context Length
5 |           |
6 |           ▼
7 |             If too long, Summarize
8 |           |
9 |           ▼
10|             Load Relevant Context
11|           |
12|           ▼
13|             Send to LLM with
14|               Next User Input

```

## Implementation of Memory for an AI Agent

```

1 from datetime import datetime
2 import uuid
3 from typing import List, Dict, Optional
4 from dataclasses import dataclass
5 import psycopg2
6 from psycopg2.extras import Json, UUID_adapter
7
8 @dataclass
9 class MemoryConfig:
10     """Configuration for memory management"""
11     max_messages: int = 20 # When to summarize
12     summary_length: int = 2000 # Max summary length in words
13     db_connection: str = "dbname=your_db user=your_user password=your_pass"
14
15 class AgentMemory:
16     def __init__(self, config: Optional[MemoryConfig] = None):
17         self.config = config or MemoryConfig()
18         self.session_id = uuid.uuid4()
19         self.setup_database()
20
21     def setup_database(self):
22         """Create necessary database tables if they don't exist"""
23         queries = [
24             """
25             CREATE TABLE IF NOT EXISTS conversations (
26                 id SERIAL PRIMARY KEY,
27                 session_id UUID NOT NULL,
28                 user_input TEXT NOT NULL,
29                 agent_response TEXT NOT NULL,
30                 tool_calls JSONB,
31                 timestamp TIMESTAMPTZ DEFAULT NOW()

```

```

32         );
33         """
34         """
35
36             CREATE TABLE IF NOT EXISTS conversation_summaries (
37                 id SERIAL PRIMARY KEY,
38                 session_id UUID NOT NULL,
39                 summary TEXT NOT NULL,
40                 start_time TIMESTAMPTZ NOT NULL,
41                 end_time TIMESTAMPTZ NOT NULL,
42                 message_count INTEGER NOT NULL
43             );
44         """
45     ]
46
47     with psycopg2.connect(self.config.db_connection) as conn:
48         with conn.cursor() as cur:
49             for query in queries:
50                 cur.execute(query)
51
52     def store_interaction(self,
53                         user_input: str,
54                         agent_response: str,
55                         tool_calls: Optional[List[Dict]] = None):
56         """Store a single interaction in the database"""
57         query = """
58             INSERT INTO conversations
59                 (session_id, user_input, agent_response, tool_calls)
60             VALUES (%s, %s, %s, %s)
61         """
62
63         with psycopg2.connect(self.config.db_connection) as conn:
64             with conn.cursor() as cur:
65                 cur.execute(query, (
66                     self.session_id,
67                     user_input,
68                     agent_response,
69                     Json(tool_calls) if tool_calls else None
70                 ))
71
72     def create_summary(self, messages: List[Dict]) -> str:
73         """Create a summary of messages using the LLM"""
74         client = OpenAI()
75
76         # Prepare messages for summarization
77         summary_prompt = f"""
78             Summarize the following conversation in less than {self.config.summary}
79             Focus on key points, decisions, and important information discovered in
80             Conversation:
```

```
81     {messages}
82     """
83
84     response = client.chat.completions.create(
85         model="gpt-4o",
86         messages=[{"role": "user", "content": summary_prompt}]
87     )
88
89     return response.choices[0].message.content
90
91     def store_summary(self, summary: str, start_time: datetime, end_time: datetime):
92         """Store a conversation summary"""
93         query = """
94             INSERT INTO conversation_summaries
95                 (session_id, summary, start_time, end_time, message_count)
96             VALUES (%s, %s, %s, %s, %s)
97             """
98
99         with psycopg2.connect(self.config.db_connection) as conn:
100             with conn.cursor() as cur:
101                 cur.execute(query, (
102                     self.session_id,
103                     summary,
104                     start_time,
105                     end_time,
106                     message_count
107                 ))
108
109     def get_recent_context(self) -> str:
110         """Get recent conversations and summaries for context"""
111         # First, get the most recent summary
112         summary_query = """
113             SELECT summary, end_time
114             FROM conversation_summaries
115             WHERE session_id = %s
116             ORDER BY end_time DESC
117             LIMIT 1
118             """
119
120         # Then get conversations after the summary
121         conversations_query = """
122             SELECT user_input, agent_response, tool_calls, timestamp
123             FROM conversations
124             WHERE session_id = %s
125             AND timestamp > %s
126             ORDER BY timestamp ASC
127             """
128
129         with psycopg2.connect(self.config.db_connection) as conn:
```

```
130         with conn.cursor() as cur:
131             # Get latest summary
132             cur.execute(summary_query, (self.session_id,))
133             summary_row = cur.fetchone()
134
135             if summary_row:
136                 summary, last_summary_time = summary_row
137                 # Get conversations after the summary
138                 cur.execute(conversations_query, (self.session_id, last_si
139             else:
140                 # If no summary exists, get recent conversations
141                 cur.execute("""
142                     SELECT user_input, agent_response, tool_calls, timestamp
143                     FROM conversations
144                     WHERE session_id = %s
145                     ORDER BY timestamp DESC
146                     LIMIT %s
147                     """, (self.session_id, self.config.max_messages))
148
149                 conversations = cur.fetchall()
150
151             # Format context
152             context = []
153             if summary_row:
154                 context.append(f"Previous conversation summary: {summary}")
155
156             for conv in conversations:
157                 user_input, agent_response, tool_calls, _ = conv
158                 context.append(f"User: {user_input}")
159                 if tool_calls:
160                     context.append(f"Tool Usage: {tool_calls}")
161                     context.append(f"Assistant: {agent_response}")
162
163             return "\n".join(context)
164
165     def check_and_summarize(self):
166         """Check if we need to summarize and do it if necessary"""
167         query = """
168             SELECT COUNT(*)
169             FROM conversations
170             WHERE session_id = %s
171             AND timestamp > (
172                 SELECT COALESCE(MAX(end_time), '1970-01-01'::timestamptz)
173                 FROM conversation_summaries
174                 WHERE session_id = %s
175             )
176         """
177
178         with psycopg2.connect(self.config.db_connection) as conn:
```

```
179     with conn.cursor() as cur:
180         cur.execute(query, (self.session_id, self.session_id))
181         count = cur.fetchone()[0]
182
183         if count >= self.config.max_messages:
184             # Get messages to summarize
185             cur.execute("""
186                 SELECT user_input, agent_response, tool_calls, timestamp
187                 FROM conversations
188                 WHERE session_id = %s
189                 ORDER BY timestamp ASC
190                 LIMIT %s
191             """ , (self.session_id, count))
192
193             messages = cur.fetchall()
194             if messages:
195                 # Create and store summary
196                 summary = self.create_summary(messages)
197                 self.store_summary(
198                     summary,
199                     messages[0][3], # start_time
200                     messages[-1][3], # end_time
201                     len(messages)
202                 )
203
204 # Update Agent class to use memory.
205 class MemoryAgent(Agent):
206     def __init__(self, memory_config: Optional[MemoryConfig] = None):
207         self.client = OpenAI()
208         self.memory = AgentMemory(memory_config)
209         self.messages = []
210
211         # Initialize with system prompt
212         self.messages.append({
213             "role": "system",
214             "content": "You are a helpful AI assistant..."
215         })
216
217     def process_query(self, user_input: str) -> str:
218         try:
219             # Check if we need to summarize before adding new messages
220             self.memory.check_and_summarize()
221
222             # Get context (including summaries) from memory
223             context = self.memory.get_recent_context()
224
225             # Add context to messages if it exists
226             if context:
227                 self.messages = [
```

```
228         self.messages[0], # Keep system prompt
229     {
230         "role": "system",
231         "content": f"Previous conversation context:\n{context}"
232     }
233 ]
234
235     # Process the query as before...
236     response = super().process_query(user_input)
237
238     # Store the interaction in memory
239     tool_calls = None
240     if hasattr(self, 'last_tool_calls'):
241         tool_calls = self.last_tool_calls
242
243     self.memory.store_interaction(
244         user_input=user_input,
245         agent_response=response,
246         tool_calls=tool_calls
247     )
248
249     return response
250
251 except Exception as e:
252     error_message = f"Error processing query: {str(e)}"
253     self.memory.store_interaction(
254         user_input=user_input,
255         agent_response=error_message
256     )
257     return error_message
258
259 def execute_tool(self, tool_call: Any) -> str:
260     # Store tool calls for memory
261     if not hasattr(self, 'last_tool_calls'):
262         self.last_tool_calls = []
263
264     result = super().execute_tool(tool_call)
265
266     # Store tool call information
267     self.last_tool_calls.append({
268         'tool': tool_call.function.name,
269         'arguments': tool_call.function.arguments,
270         'result': result
271     })
272
273     return result
```

---

Before we continue...

## Interested in how to train your very own Large Language Model?

We prepared a well-researched guide for how to use the latest advancements in Open Source technology to fine-tune your own LLM. This has many advantages like:

- Cost control
- Data privacy
- Excellent performance - adjusted specifically for your intended use

[Get your free LLM training guide](#)

---

## Step 4: Using Your AI Agent

After all the setup, using the agent is surprisingly straightforward. Here's how:

```
1 from your_agent import Agent, MemoryConfig # assuming we saved our code in you
2
3 # Create an agent
4 agent = Agent(
5     memory_config=MemoryConfig(
6         max_messages=10, # summarize after 10 messages
7         db_connection="postgresql://user:pass@localhost/your_db"
8     )
9 )
10
11 # Simple question-answer interaction
12 def chat_with_agent(question: str):
13     print(f"\nUser: {question}")
14     print(f"Assistant: {agent.process_query(question)}")
15
16 if __name__ == "__main__":
```

```
17     chat_with_agent("How many users do we have in our database?")
18     chat_with_agent("What's the average age of users from New York?")
19     chat_with_agent("Compare our user base in California with PostgreSQL's popul
20     chat_with_agent("What were the numbers you just mentioned?")
21     chat_with_agent("What are the most interesting trends you can find in our t
```

That's it! The agent will:

- Use the database when needed
- Search Wikipedia when appropriate
- Remember previous context
- Create summaries automatically
- Handle all the complexity behind the scenes - but at the same time giving us direct debugging access

The output looks clean and natural:

```
1 User: How many users do we have in our database?
2 Assistant: According to the database, we currently have 10,432 registered users
3
4 User: What's the average age of users from New York?
5 Assistant: I've queried the database and found that users from New York have an
6
7 User: Compare our user base in California with PostgreSQL's popularity there
8 Assistant: Let me check both sources. According to our database, we have 2,345 u
```



## Enhanced Database Tool with Schema Awareness

Arguably, the above PostgreSQL database tool is a little simple isn't it? From where should the AI model know what database columns are available? You are right with that - we left out this crucial detail for more clarity. We'll modify our database tool to include schema information, which helps the LLM understand what data is available and how to query it correctly:

```

1 def get_database_schema() -> str:
2     """Retrieve the database schema information"""
3     schema_query = """
4     SELECT
5         t.table_name,
6         array_agg(
7             c.column_name || ' ' ||
8             c.data_type ||
9             CASE
10                 WHEN c.is_nullable = 'NO' THEN ' NOT NULL'
11                 ELSE ''
12             END
13         ) as columns
14     FROM information_schema.tables t
15     JOIN information_schema.columns c
16         ON c.table_name = t.table_name
17     WHERE t.table_schema = 'public'
18         AND t.table_type = 'BASE TABLE'
19     GROUP BY t.table_name;
20 """
21
22     try:
23         with psycopg2.connect("dbname=your_db user=your_user password=your_pa:
24             with conn.cursor() as cur:
25                 cur.execute(schema_query)
26                 schema = cur.fetchall()
27
28                 # Format schema information
29                 schema_str = "Database Schema:\n"
30                 for table_name, columns in schema:
31                     schema_str += f"\n{table_name}\n"
32                     for col in columns:
33                         schema_str += f" - {col}\n"
34
35             return schema_str
36     except Exception as e:
37         return f"Error fetching schema: {str(e)}"
38
39 # Updated tool definition with schema information;
40 # Note that we could also dynamically fetch the schema when needed
41 tools = [
42     {
43         "type": "function",
44         "function": {
45             "name": "query_database",
46             "description": """Execute a PostgreSQL SELECT query and return the
47             Available tables and their schemas:
48

```

```

49         users
50             - id SERIAL PRIMARY KEY
51             - email VARCHAR(255) NOT NULL
52             - age INTEGER
53             - location VARCHAR(100)
54             - signup_date DATE NOT NULL
55             - last_login TIMESTAMP WITH TIME ZONE
56             - job_industry VARCHAR(100)
57
58     user_activity
59         - id SERIAL PRIMARY KEY
60         - user_id INTEGER REFERENCES users(id)
61         - activity_date DATE NOT NULL
62         - activity_type VARCHAR(50) NOT NULL
63     """
64     "parameters": {
65         "type": "object",
66         "properties": {
67             "query": {
68                 "type": "string",
69                 "description": "The SQL SELECT query to execute. Must
70                     }
71             },
72             "required": ["query"],
73             "additionalProperties": False
74         },
75         "strict": True
76     }
77 }
78 ]
79
80 # Updated database query function
81 def query_database(query: str) -> str:
82     """
83     Execute a PostgreSQL query with schema awareness
84     """
85     if not query.lower().strip().startswith('select'):
86         return json.dumps({
87             "error": "Only SELECT queries are allowed for security reasons.",
88             "schema": get_database_schema() # Return schema for reference
89         })
90
91     try:
92         with psycopg2.connect("dbname=your_db user=your_user password=your_pa:
93             with conn.cursor() as cur:
94                 cur.execute(query)
95
96                 # Get column names from cursor description
97                 columns = [desc[0] for desc in cur.description]

```

```

98
99          # Fetch results and convert to list of dictionaries
100         results = []
101         for row in cur.fetchall():
102             results.append(dict(zip(columns, row)))
103
104         return json.dumps({
105             "success": True,
106             "data": str(results),
107             "row_count": len(results),
108             "columns": columns
109         })
110
111     except Exception as e:
112         return json.dumps({
113             "error": str(e),
114             "schema": get_database_schema() # Return schema on error for help
115         })
116
117 # Usage example:
118 agent = Agent()
119
120 # The agent now knows about the schema and can write better queries
121 response = agent.process_query(
122     "What's the age distribution of users by location and industry?"
123 )

```

Key improvements:

1. Schema information in tool description
2. Column types and constraints included
3. Schema returned on errors for better error handling
4. Detailed result structure with column names

This helps the LLM:

- Write correct queries on the first try
- Understand available data
- Handle errors better
- Provide more detailed responses

# Further Reading

- [DSPy: Automatic prompt tuning and creating](#)
- [Free and open source ChatGPT alternative: LibreChat](#)
- [How to: RAG with Azure AI search](#)

More information on our  
managed RAG solution?

To Pondhouse AI

More tips and tricks on how to  
work with AI?

To our Blog



We know data.

## ⌚ Location

Pondhouse Data OG  
Litscherweg 10/21  
6800 Feldkirch  
Austria

## 📞 Phone

+43 650 9820448

## 🔗 Links

[Impressum](#)  
[Data Protection](#)  
[Article-Archive](#)

---

© 2025 Pondhouse Data OG. All rights reserved.