

A Comparison of Assistant and Non-Assistant Tool Calling Pipelines

Introduction

Part 1: Which System to Use: Non-Assistant or Assistant?

When to Use Non-Assistant Tool Calling Pipelines:

When to Use Assistant Tool Calling Pipelines:

Part 2: A brief introduction to two pipelines

Non-assistant tool calling pipelines

Assistant tool calling pipelines

Part 3: Comparison between two pipelines

Part 4: Differences in snap settings

Stop condition of Pipeloop

Part 5: Assistant's two built-in tools

File search

Code Interpreter

Part 6: Key Differences Summarized

Introduction

At a high level, the logic behind assistant tool calling and non-assistant tool calling is fundamentally the same: the model instructs the user to call specific function(s) in order to answer the user's query. The user then executes the function and returns the result to the model, which uses it to generate an answer. This process is identical for both.

However, since the assistant specifies the function definitions and access to tools as part of the Assistant configuration within the OpenAI or Azure OpenAI dashboard rather than within your pipelines, there will be major differences in the pipeline configuration. Additionally submitting tool responses to an Assistant comes with significant changes and challenges since the Assistant owns the conversational history rather than the pipeline.

This article focuses on contrasting these differences. For a detailed understanding of assistant pipelines and non-assistant pipelines, please refer to the following article:

Non-assistant pipelines: [Introducing Tool Calling Snaps and LLM Agent Pipelines](#)

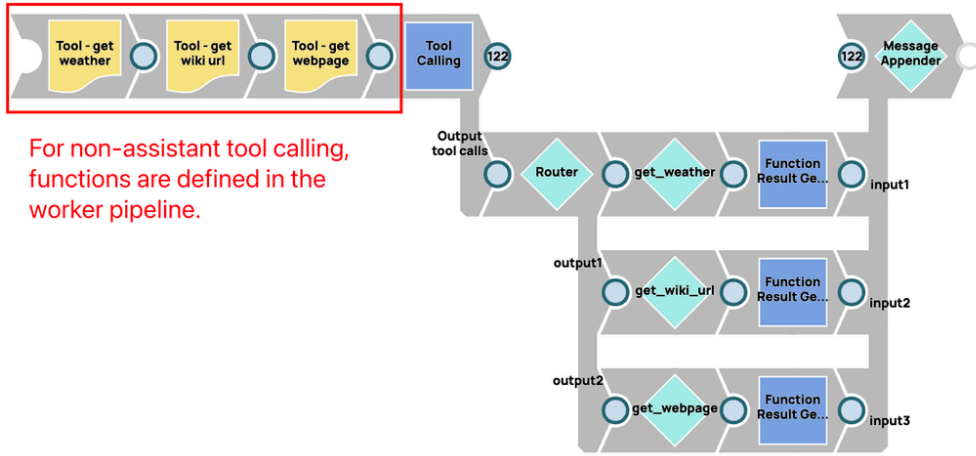
Assistant pipelines: [Introducing Assistant Tool Calling Pipelines](#)

Part 1: Which System to Use: Non-Assistant or Assistant?

When to Use Non-Assistant Tool Calling Pipelines:

Non-Assistant Tool Calling Pipelines offer greater flexibility and control over the tool calling process, making them suitable for the following specific scenarios.

- **When preferring a "run-time" approach:** Non-Assistant pipelines exhibit greater flexibility in function definition, offering a more "runtime" approach. You can dynamically adjust the available functions by simply adding or removing Function Generator snaps within the pipeline.



For non-assistant tool calling, functions are defined in the worker pipeline.

Non-assistant Tool Calling

In contrast, Assistant Tool Calling Pipelines necessitate a "design-time" approach. All available functions must be pre-defined within the Assistant configuration, requiring modifications to the Assistant definition in the OpenAI/Azure OpenAI dashboard.

ASSISTANT

Playground ↗

Name

Test Assistant

asst_nwIrrRaBwD6E6xa7EmnD0y2fx

System instructions

You are a helpful assistant...

Model

gpt-4o

TOOLS

File search ⓘ

+ Files

Vector store for Test-assistant
vs_hovZYr4f8Y1dxvk7L9rohHEo

15 KB

Code interpreter ⓘ

+ Files

Functions ⓘ

{f} get_weather
{f} get_wiki_url
{f} get_webpage

For assistant tool calling, functions are defined in the dashboard before running a pipeline.

+ Functions

Assistant Tool Calling

- **When wanting detailed chat history:** Non-Assistant pipelines provide a comprehensive history of the interaction between the model and the tools in the output message list. The message list within the Non-Assistant pipeline preserves every model response and the results of each function execution. This detailed logging allows for thorough debugging, analysis, and auditing of the tool calling process.

```

▼ [
  ▼ {
    ▶ "output": {"message": {...}},
    "stopReason": end_turn,
    ▶ "usage": {"inputTokens": 2138, "outputTokens": 435, "totalTokens": 2573},
    ▶ "metrics": {"latencyMs": 10329},
    ▼ "messages": [
      ▶ {"role": "user", "content": [...]}, user's original prompt
      ▶ {"role": "assistant", "content": [...]}, model requests tool call
      ▶ {"content": [...], "role": "user"}, user provides tool execution result
      ▶ {"role": "assistant", "content": [...]},
      ▶ {"content": [...], "role": "user"},
      ▶ {"role": "assistant", "content": [...]}, repeat...
      ▶ {"content": [...], "role": "user"},
      ▶ {"role": "assistant", "content": [...]} model's final response
    ],
    ▶ "original": {"tools": [...], "original": {...}},
    ▶ "original_0": {"messages": [...], "original": {...}}
  }
]

```

Non-assistant Tool Calling

In contrast, Assistant pipelines maintain a more concise message history, focusing on key steps and omitting some intermediate details. While this can simplify the overall view of the message list, it can also make it more difficult to trace the exact sequence of events or diagnose issues that may arise during tool execution in child pipelines.

```

▼ [
  ▼ {
    ▼ "messages": [
      ▶ {"id": "msg_sRghgW4fm5Ds9CbtCZIOhno", "object": "thread.message", "created_at": 1736805916, "assi... },
      ▶ {"id": "msg_LqjnaB9potXT55fogEmh0920", "object": "thread.message", "created_at": 1736805872, "assist... }
    ],
    ▶ "run": {"id": "run_N7qeXD13RCElel6WLX10MGIX", "object": "thread.run", "created_at": 1736805872, "assistant_id"... },
    ▶ "original": {"run": [...], "original": [...], "tool_outputs": [...]},
    ▶ "original_0": {"run": [...]}
  }
]

```

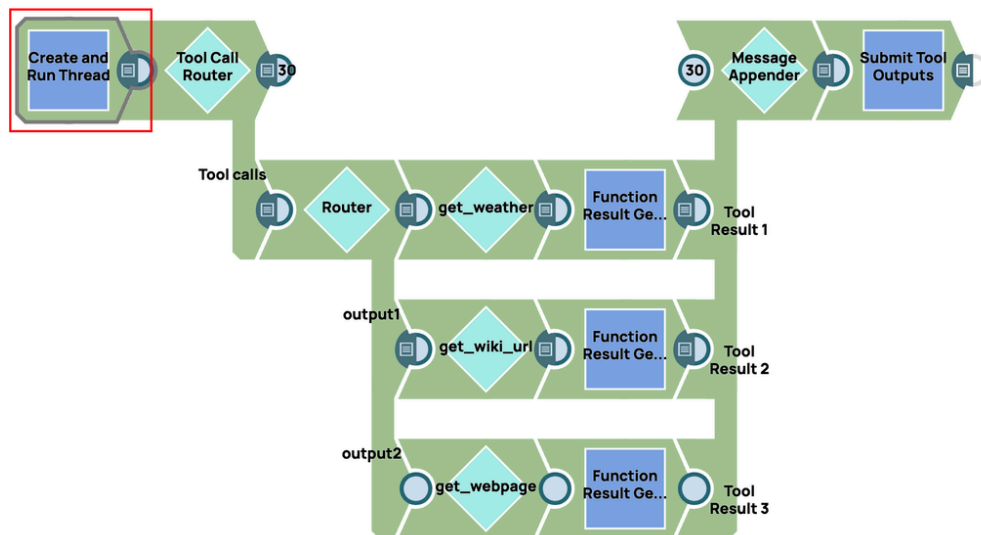
Assistant Tool Calling

- When needing easier debugging and iterative development:** Non-Assistant pipelines facilitate more granular debugging and iterative development. You can easily simulate individual steps of the agent by making calls to the model with specific function call histories. This allows for more precise control and experimentation during development, enabling you to isolate and address issues more effectively. For example, by providing three messages, we can "force" the model to call the second tool, allowing us to inspect the tool calling process and its result against our expectations.

The screenshot shows a 'JSON Generator1 output0' window with a preview of JSON data. The JSON includes fields like 'totalTokens', 'latencyMs', and 'messages'. The 'messages' array contains three entries: a user prompt, an assistant request for 'tool 1', and the result of 'tool 1'. To the right, a pipeline diagram shows a 'Tool Call Router' leading to three parallel paths for 'get_weather', 'get_wiki_url', and 'get_webpage', each followed by a 'Function Result Generator' and an 'Input' node. A red box highlights the 'get_wiki_url' path.

In contrast, debugging and iterating with Assistant pipelines can be more cumbersome. Since Assistants manage the conversation history internally, to simulate a specific step, you often need to replay the entire interaction from the beginning, potentially requiring multiple iterations to reach the desired state. This internal management of history makes it less straightforward to isolate and debug specific parts of the interaction.

To simulate calling the third tool, we need to start a new thread from scratch and then call tool1 and tool2, repeating the preceding process. The current thread cannot be reused.



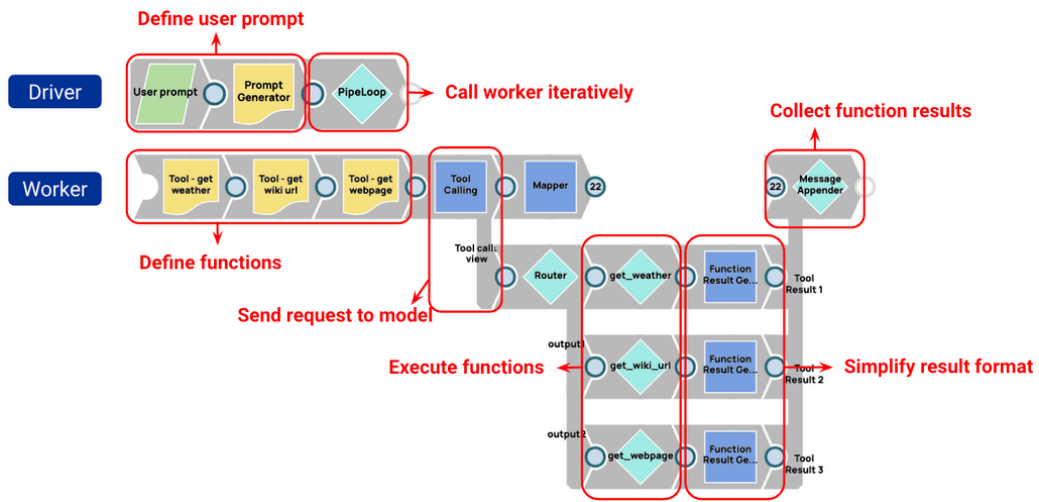
When to Use Assistant Tool Calling Pipelines:

Assistant Tool Calling Pipelines also offer a streamlined approach to integrating LLMs with external tools, prioritizing ease of use and built-in functionalities. Consider using Assistant pipelines in the following situations:

- **For simplified pipeline design:** Assistant pipelines reduce pipeline complexity by eliminating the need for Tool Generator snaps. In Non-Assistant pipelines, these snaps are essential for dynamically generating tool definitions within the pipeline itself. With Assistant pipelines, tool definitions are configured beforehand within the Assistant settings in the OpenAI/Azure OpenAI dashboard. This pre-configuration results in shorter, more manageable pipelines, simplifying development and maintenance.
- **When leveraging built-in tools is required:** If your use case requires functionalities like searching external files or executing code, Assistant pipelines offer these capabilities out-of-the-box through their built-in File Search and Code Interpreter tools (see Part 5 for more details). These tools provide a convenient and efficient way to extend the LLM's capabilities without requiring custom implementation within the pipeline.

Part 2: A brief introduction to two pipelines

Non-assistant tool calling pipelines

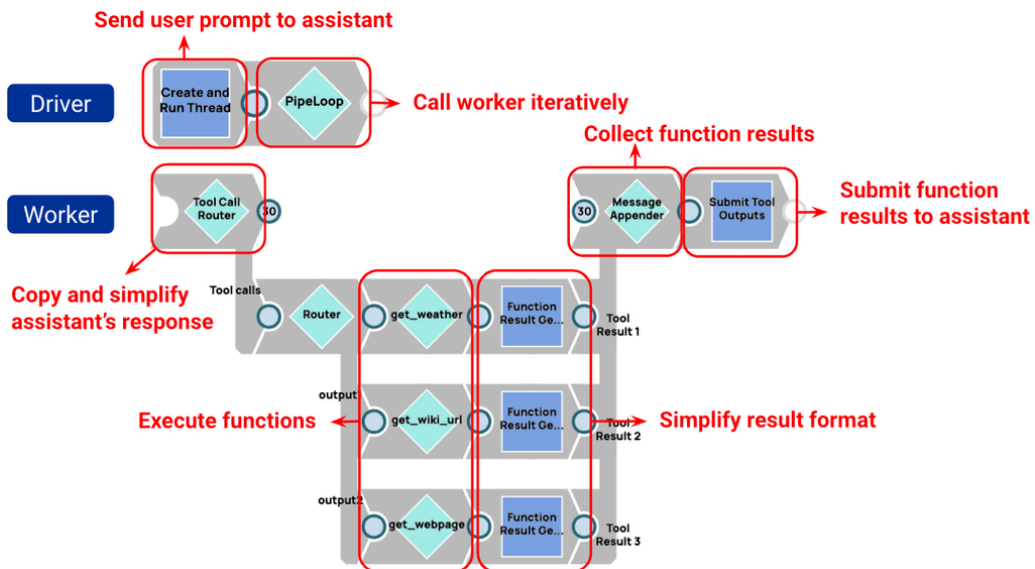


Non-assistant tool calling pipelines

Key points:

- Functions are defined in the worker.
- The worker pipeline's Tool Calling snap manages all model interactions.
- Function results are collected and sent to the model in the next iteration via the Tool Calling snap.

Assistant tool calling pipelines



Assistant tool calling pipelines

Key points:

- No need to define functions in any pipeline. Functions are pre-defined in the assistant.
- Two snaps : interact with the model: **Create and Run Thread**, and **Submit Tool Outputs**.
- Function results are collected and sent to the model immediately during the current iteration.

Part 3: Comparison between two pipelines

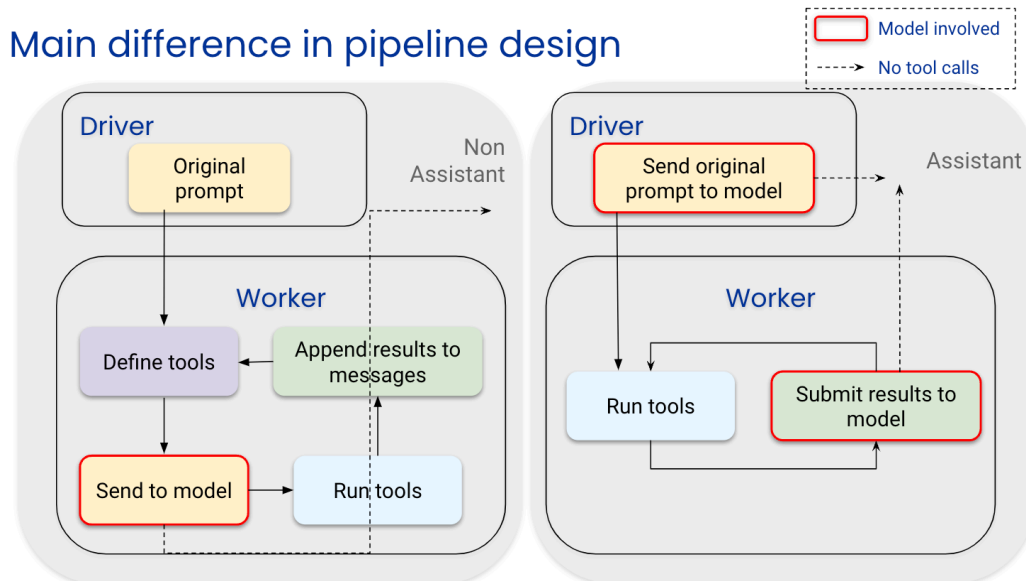
Here are two primary reasons why the assistant and non-assistant pipelines differ, listed in decreasing order of importance:

1. Distinct methods of submitting tool results:

- For non-assistant pipelines, tool results are appended to the message history list and subsequently forwarded to the model during the **next** iteration.
Non-assistant pipelines exhibit a "**while-loop**" behavior, where the worker interacts with the model at the beginning of the iteration, and **while any tools need to be called**, the worker **executes those tool(s)**.
- In contrast, for assistants, tool results are specifically sent to a dedicated endpoint designed to handle tool call results within the **current** iteration.
The assistant pipelines operate more like a "**do-while-loop**." The driver initiates the interaction by sending the prompt to the model. Subsequently, the worker **execute the tool(s) first** and **interacts with the model at the end of the iteration to deliver tool results**.

2. Predefined and stored tool definitions for assistants:

- Unlike non-assistant pipelines, assistants have the capability to predefine and store function definitions. This eliminates the need for the three `Function Generator` snaps to repeatedly transmit tool definitions to the model with each request. Consequently, the worker pipeline for assistants appears shorter.



Due to the aforementioned differences, non-assistant pipelines have only one interaction point with the model, located in the worker.

In contrast, assistant pipelines involve two interaction points: the driver sends the initial prompt to the model, while the worker sends tool results back to the model.

Part 4: Differences in snap settings

Stop condition of Pipeline

A key difference in snap settings lies in the stop condition of the pipeline.

i Assistant pipeline's stop condition: `$run.required_action == null`.

Non-assistant pipeline's stop condition: `$finish_reason != "tool_calls"`.

Assistant's output

```

    "required_action": {
      "type": "submit_tool_outputs",
      "submit_tool_outputs": {
        "tool_calls": [
          {
            "id": "call_xj74QONQ9XhpLSGMyhd0WltV",
            "type": "function",
            "function": {
              "name": "get_webpage",
              "arguments": "{\"url\":\"https://en.wikipedia.org/wiki/San_Francisco_CA\"}"
            }
          }
        ]
      }
    }
  },
},
],
),

```

Example when tool calls are required

```

"cancelled_at": null,
"failed_at": null,
"completed_at": 1731345620,
"required_action": null,
"last_error": null,
"model": "gpt-4o",
"instructions": "",
"tools": [{...}, {...}, {...}, {...}, {...}],
"tool_resources": {},
"metadata": {},
"temperature": 1.0

```

Example when tool calls are NOT required

Non-assistant's output

```

"choices": [
  {
    "index": 0,
    "message": {
      "role": "assistant",
      "content": null,
      "tool_calls": [
        { "id": "call_IWfxGMVOAS2cgA9bb5aWsQNL", "type": "function", "function": {...} },
        { "id": "call_AkOuBx6mGuEPu0EH1yy7sY4j", "type": "function", "function": {...} }
      ],
      "refusal": null
    },
    "logprobs": null,
    "finish_reason": "tool_calls"
  }
]

```

Example when tool calls are required

```

"choices": [
  {
    "index": 0,
    "message": { "role": "assistant",
    "logprobs": null,
    "finish_reason": "stop"
  }
]

```

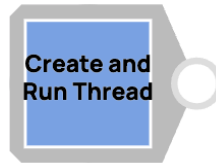
Example when tool calls are NOT required

Part 5: Assistant's two built-in tools

The assistant not only supports all functions that can be defined in non-assistant pipelines but also provides two special built-in functions, **file search** and **code interpreter**, for user convenience.

If the model determines that either of these tools is required, it will automatically call and execute the tool within the assistant without requiring manual user intervention.

i You don't need a tool call pipeline to experiment with file search and code interpreter. A simple create and run thread snap is sufficient.



File search

File Search augments the Assistant with knowledge from outside its model, such as proprietary product information or documents provided by your users. OpenAI automatically parses and chunks your documents, creates and stores the embeddings, and use both vector and keyword search to retrieve relevant content to answer user queries.

Example

Prompt: What is the number of federal fires between 2018 and 2022?

The assistant's response is as below:

Assistant's response

```
1  [
2  {
3    "messages": [
4      {
5        "id": "msg_cyvIQG7htmHnwTTbkrES3ms",
6        "object": "thread.message",
7        "created_at": 1731106910,
8        "assistant_id": "asst_nwIrRaBwD6E6xa7EmnD0y2fx",
9        "thread_id": "thread_ciR3mFR1jEcXK07pX06jCRgM",
10       "run_id": "run_6lXt9zvpXLYxfgIff7GV8Nz2",
11       "role": "assistant",
12       "content": [
13         {
14           "type": "text",
15           "text": {
16             "value": "The number of federal fires between 2018 and 2022 is as follows:\n\n- 2018:
17             12,500\n- 2019: 10,900\n- 2020: 14,400\n- 2021: 14,000\n- 2022: 11,700 [4:1twildfire_stats.pdf] .",
18             "annotations": [
19               {
20                 "type": "file_citation",
21                 "text": " [4:1twildfire_stats.pdf] ",
22                 "start_index": 140,
23                 "end_index": 164,
24                 "file_citation": {
25                   "file_id": "file-fJGINZ4R7XlIGtjfvv0W71CH"
26                 }
27               }
28             ]
29           }
30         }
31       ]
32     }
33   ]
34 }
```



```

26     }
27   ]
28 }
29 }
30 ],
31 "attachments": [],
32 "metadata": {}
33 },
34 {
35   "id": "msg_LBP4fengd7G1Qnu7ZkfqvM2W",
36   "object": "thread.message",
37   "created_at": 1731106907,
38   "assistant_id": null,
39   "thread_id": "thread_ciR3mFR1jEcXK07pX06jCRgM",
40   "run_id": null,
41   "role": "user",
42   "content": [
43     {
44       "type": "text",
45       "text": {
46         "value": "What is the number of federal fires between 2018 and 2022",
47         "annotations": []
48       }
49     }
50   ],
51   "attachments": [],
52   "metadata": {}
53 }
54 ],
55 "run": {...}
56 }
57 ]

```

The assistant's response is correct. As the answer to the prompt is in the first row of a table on the first page of `wildfire_stats.pdf`, a document accessible to the assistant via a vector store.

wildfire_stats.pdf
Page 1 of 3

(78,600).

Table 1. Annual Wildfires and Acres Burned

	2018	2019	2020	2021	2022
Number of Fires (thousands)					
Federal	12.5	10.9	14.4	14.0	11.7
FS	5.6	5.3	6.7	6.2	5.9
DOI	7.0	5.3	7.6	7.6	5.8
Other	0.1	0.2	<0.1	0.2	0.1
Nonfederal	45.6	39.6	44.6	45.0	57.2
Total	58.1	50.5	59.0	59.0	69.0
Acres Burned (millions)					
Federal	4.6	3.1	7.1	5.2	4.0
FS	2.3	0.6	4.8	4.1	1.9
DOI	2.3	2.3	2.3	1.0	2.1
Other	<0.1	<0.1	<0.1	<0.1	<0.1
Nonfederal	4.1	1.6	3.1	1.9	3.6
Total	8.8	4.7	10.1	7.1	7.6

Source: National Interagency Coordination Center (NICC) Wildland Fire Summary and Statistics annual reports.
Notes: FS = Forest Service; DOI = Department of the Interior. Column totals may not sum precisely due to rounding.

<https://crsreports.congress.gov>

Acres Burned Since 1960

Year	Acres burned (millions)	Number of Fires
2015	10.13	68.2
2020	10.12	59.0
2017	10.03	71.5
2006	9.87	96.4
2007	9.33	67.8

Source: NICC Wildland Fire Summary and Statistics annual reports.
Note: Number of fires in thousands.


The number of fires and acreage burned are indicators of the annual level of wildfire activity. These numbers may not be indicative of fire's impact on human development or communities, since many fires occur in large, relatively undeveloped areas. Acreage burned also does not indicate the severity of a wildfire, the degree of impact upon forests or soils, or other ecological effects.

Most wildfires are human-caused (89% of the average number of wildfires from 2018 to 2022). Wildfires caused by lightning tend to be slightly larger and to burn more acreage (53% of the average acreage burned from 2018 to 2022) than human-caused fires.



Answer to the prompt

VECTOR STORE

Vector store for Test-assistant

ID	vs_hovZYr4f8Y1dxvk7L9xohHEo
Usage this month	0 KB hours · \$0.1 / GB per day
Size	15 KB
Last active	Nov 8, 2024, 3:01 PM
Expiration policy	Never 
Expires	Never
Created	Jun 12, 2024, 12:19 PM

Files attached

FILE	UPLOADED	
 wildfire_stats.pdf	10/8/2024, 3:51 PM	

Used by

[+ Create assistant](#)


Resource	ID
Test Assistant	asst_nwlrRaBwD6E6xa7EmnD0y2fx

The file is from a vector store used by the assistant.

Code Interpreter

Code Interpreter allows Assistants to write and run Python code in a sandboxed execution environment. This tool can process files with diverse data and formatting, and generate files with data and images of graphs. Code Interpreter allows your Assistant to run code iteratively to solve challenging code and math problems. When your Assistant writes code that fails to run, it can iterate on this code by attempting to run different code until the code execution succeeds.

Example

 **Prompt:** Find the number of federal fires between 2018 and 2022 and use Matplotlib to draw a line chart.

* Matplotlib is a python library for creating plots.

The assistant's response is as below:

Assistant's response

```
1  [
2  {
3    "messages": [
4      {
5        "id": "msg_lzBiM0J4sC0Zji510f1N0jjm",
6        "object": "thread.message",
7        "created_at": 1731108369,
8        "assistant_id": "asst_nwlrRaBwD6E6xa7EmnD0y2fx",
9        "thread_id": "thread_3q9AV6ivrYqqzsexv1rzMFSV",
10       "run_id": "run_DbjZQbBVQgoVge74PRbyGh44",
11       "role": "assistant",
12       "content": [
13         {
14           "type": "image_file",
15           "image_file": {
16             "file_id": "file-CLH0iYRfuWD45DsN6M4b8ga9"
17         }
18       },
19       {
20         "type": "text",
```

```

21     "text": {
22         "value": "Here is the line chart showing the number of federal fires from 2018 to 2022. As
you can see, there is a fluctuation in the number of fires over these years, with a peak in 2020.",
23         "annotations": []
24     }
25 }
26 ],
27 "attachments": [],
28 "metadata": {}
29 },
30 {
31     ...,
32     "content": [
33         {
34             "type": "text",
35             "text": {
36                 "value": "The number of federal fires between 2018 and 2022 was as follows:\n\n- 2018: 12.5
thousand fires\n- 2019: 10.9 thousand fires\n- 2020: 14.4 thousand fires\n- 2021: 14.0 thousand fires\n-
2022: 11.7 thousand fires [4:0†wildfire_stats.pdf] .\n\nI will now create a line chart using Matplotlib to
represent this data.",
37                 "annotations": [
38                     {
39                         "type": "file_citation",
40                         "text": " [4:0†wildfire_stats.pdf] ",
41                         "start_index": 206,
42                         "end_index": 230,
43                         "file_citation": {
44                             "file_id": "file-fJGINZ4R7XlIGtjfvv0W71CH"
45                         }
46                     }
47                 ]
48             }
49         }
50     ],
51     ...
52 },
53 {
54     ...,
55     "content": [
56         {
57             "type": "text",
58             "text": {
59                 "value": "Find the number of federal fires between 2018 and 2022 and use Matplotlib to draw a
line chart.",
60                 "annotations": []
61             }
62         }
63     ],
64     ...
65 }
66 ],
67 "run": {...}
68 }
69 ]

```

From the response, we can see that the assistant indicated it used file search to find 5 years of data and then generated an image file. This file can be downloaded from the assistant's dashboard under storage-files. Simply add a file extension like .png to see the image.

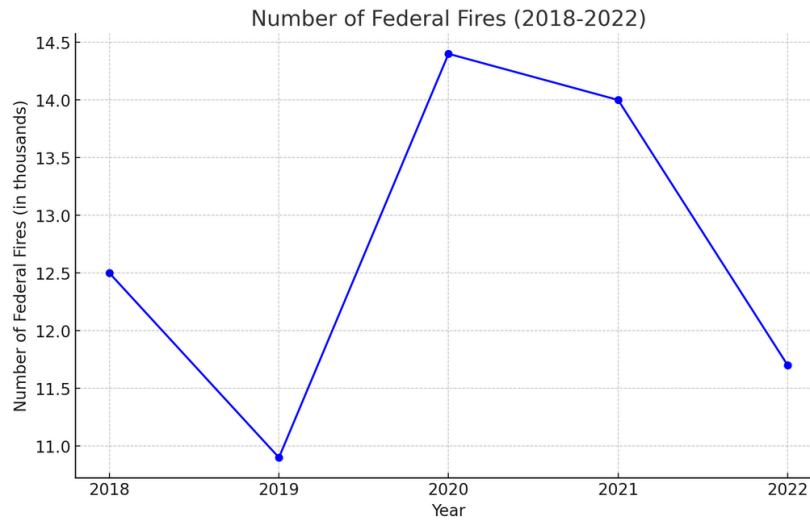


Image file generated by assistant

Part 6: Key Differences Summarized

Feature	Non-Assistant Tool Calling Pipelines	Assistant Tool Calling Pipelines
Function Definition	Defined within the worker pipeline using Function Generator snaps.	Pre-defined and stored within the Assistant configuration in the OpenAI/Azure OpenAI dashboard.
Tool Result Submission	Appended to the message history and sent to the model in the <i>next</i> iteration.	Sent to a dedicated endpoint within the <i>current</i> iteration.
Model Interaction Points	One (in the worker pipeline).	Two (driver sends initial prompt, worker sends tool results).
Built-in Tools	None.	File Search and Code Interpreter.
Pipeline Complexity	More complex pipeline structure due to function definition within the pipeline.	Simpler pipeline structure as functions are defined externally.