



API

Tool Use

Enable LLMs to interact with external functions and APIs.

Tool use enables LLMs to request calls to external functions and APIs through the `/v1/chat/completions` endpoint, via LM Studio's REST API (or via any OpenAI client). This expands their functionality far beyond text output.

Quick Start

1. Start LM Studio as a server

To use LM Studio programmatically from your own code, run LM Studio as a local server.

You can turn on the server from the "Developer" tab in LM Studio, or via the `lms` CLI:

```
lms server start
```



Install `lms` by running `npx lmstudio install-cli`

This will allow you to interact with LM Studio via an OpenAI-like REST API. For an intro to LM Studio's OpenAI-like API, see [Running LM Studio as a server](#).

2. Load a Model

You can load a model from the "Chat" or "Developer" tabs in LM Studio, or via the `lms` CLI:

```
lms load
```



3. Copy, Paste, and Run an Example!

- `curl`



- [Advanced Agent Example](#)

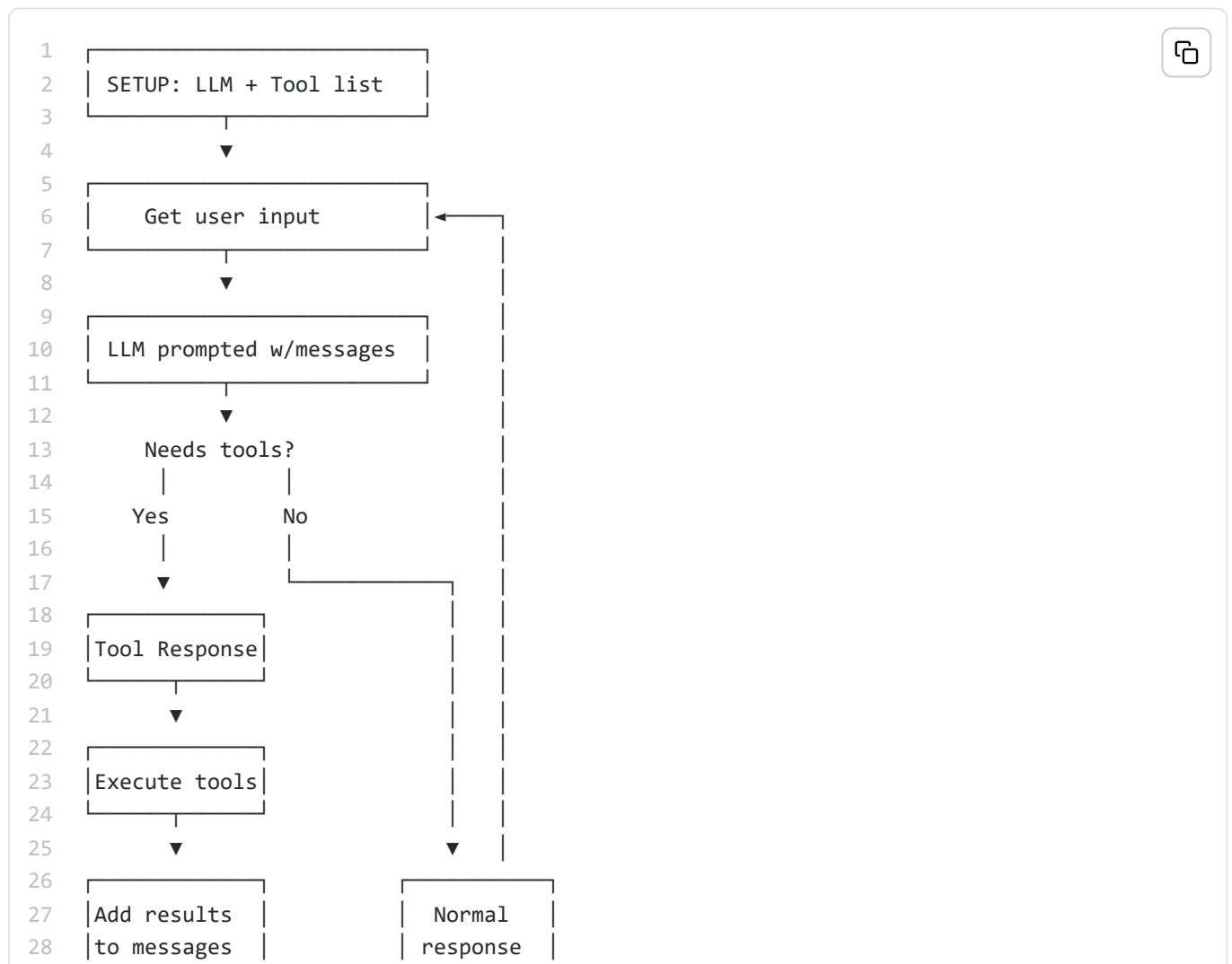
Tool Use

What really is "Tool Use"?

Tool use describes:

- LLMs output text requesting functions to be called (LLMs cannot directly execute code)
- Your code executes those functions
- Your code feeds the results back to the LLM.

High-level flow





Documentation

LM Studio supports tool use through the `/v1/chat/completions` endpoint when given function definitions in the `tools` parameter of the request body. Tools are specified as an array of function definitions that describe their parameters and usage, like:

It follows the same format as OpenAI's [Function Calling](#) API and is expected to work via the OpenAI client SDKs.

We will use [lmstudio-community/Qwen2.5-7B-Instruct-GGUF](#) as the model in this example flow.

1. You provide a list of tools to an LLM. These are the tools that the model can *request* calls to. For example:

```
1 // the list of tools is model-agnostic
2 [
3   {
4     "type": "function",
5     "function": {
6       "name": "get_delivery_date",
7       "description": "Get the delivery date for a customer's order",
8       "parameters": {
9         "type": "object",
10        "properties": {
11          "order_id": {
12            "type": "string"
13          }
14        },
15        "required": ["order_id"]
16      }
17    }
18  ]
```

This list will be injected into the `system` prompt of the model depending on the model's chat template. For `Qwen2.5-Instruct`, this looks like:

```
1 <|im_start|>system
2 You are Qwen, created by Alibaba Cloud. You are a helpful assistant.
3
4 # Tools
5
6 You may call one or more functions to assist with the user query.
7
8 You are provided with function signatures within <tools></tools> XML tags:
9 <tools>
10 {"type": "function", "function": {"name": "get_delivery_date", "description": "Get th
```



Documentation

Important: The model can only *request* calls to these tools because LLMs *cannot* directly call functions, APIs, or any other tools. They can only output text, which can then be parsed to programmatically call the functions.

2. When prompted, the LLM can then decide to either:

- (a) Call one or more tools

```
1 User: Get me the delivery date for order 123
2 Model: <tool_call>
3 {"name": "get_delivery_date", "arguments": {"order_id": "123"}}
4 </tool_call>
```



- (b) Respond normally

```
1 User: Hi
2 Model: Hello! How can I assist you today?
```



3. LM Studio parses the text output from the model into an OpenAI-compliant `chat.completion` response object.

- If the model was given access to `tools`, LM Studio will attempt to parse the tool calls into the `response.choices[0].message.tool_calls` field of the `chat.completion` response object.
- If LM Studio cannot parse any **correctly formatted** tool calls, it will simply return the response to the standard `response.choices[0].message.content` field.
- **Note:** Smaller models and models that were not trained for tool use may output improperly formatted tool calls, resulting in LM Studio being unable to parse them into the `tool_calls` field. This is useful for troubleshooting when you do not receive `tool_calls` as expected. Example of an improperly formatting Qwen2.5-Instruct tool call:

```
1 <tool_call>
2 [{"name": "get_delivery_date", function: "date"}]
3 </tool_call>
```



Note that the brackets are incorrect, and the call does not follow the `name, argument` format.



Documentation

2. The result of the tool call

To the `messages` array to send back to the model

```

1 # pseudocode, see examples for copy-paste snippets
2 if response.has_tool_calls:
3     for each tool_call:
4         # Extract function name & args
5         function_to_call = tool_call.name      # e.g. "get_delivery_date"
6         args = tool_call.arguments            # e.g. {"order_id": "123"}
7
8         # Execute the function
9         result = execute_function(function_to_call, args)
10
11        # Add result to conversation
12        add_to_messages([
13            ASSISTANT_TOOL_CALL_MESSAGE,      # The request to use the tool
14            TOOL_RESULT_MESSAGE              # The tool's response
15        ])
16 else:
17     # Normal response without tools
18     add_to_messages(response.content)

```

5. The LLM is then prompted again with the updated messages array, but without access to tools. This is because:

- The LLM already has the tool results in the conversation history
- We want the LLM to provide a final response to the user, not call more tools

```

1 # Example messages
2 messages = [
3     {"role": "user", "content": "When will order 123 be delivered?"},
4     {"role": "assistant", "function_call": {
5         "name": "get_delivery_date",
6         "arguments": {"order_id": "123"}
7     }},
8     {"role": "tool", "content": "2024-03-15"},
9 ]
10 response = client.chat.completions.create(
11     model="lmstudio-community/qwen2.5-7b-instruct",
12     messages=messages
13 )

```

The `response.choices[0].message.content` field after this call may be something like:

Your order #123 will be delivered on March 15th, 2024



Supported Models

Through LM Studio, all models support at least some degree of tool use.

However, there are currently two levels of support that may impact the quality of the experience: Native and Default.

Models with Native tool use support will have a hammer badge in the app, and generally perform better in tool use scenarios.

Native tool use support

"Native" tool use support means that both:

1. The model has a chat template that supports tool use (usually means the model has been trained for tool use)
 - This is what will be used to format the `tools` array into the system prompt and tell them model how to format tool calls
 - Example: [Qwen2.5-Instruct chat template](#)
2. LM Studio supports that model's tool use format
 - Required for LM Studio to properly input the chat history into the chat template, and parse the tool calls the model outputs into the `chat.completion` object

Models that currently have native tool use support in LM Studio (subject to change):

- Qwen
 - GGUF [lmstudio-community/Qwen2.5-7B-Instruct-GGUF](#) (4.68 GB)
 - MLX [mlx-community/Qwen2.5-7B-Instruct-4bit](#) (4.30 GB)
- Llama-3.1, Llama-3.2
 - GGUF [lmstudio-community/Meta-Llama-3.1-8B-Instruct-GGUF](#) (4.92 GB)
 - MLX [mlx-community/Meta-Llama-3.1-8B-Instruct-8bit](#) (8.54 GB)
- Mistral
 - GGUF [bartowski/Ministral-8B-Instruct-2410-GGUF](#) (4.67 GB)



Documentation

1. The model does not have chat template that supports tool use (usually means the model has not been trained for tool use)
2. LM Studio does not currently support that model's tool use format

Under the hood, default tool use works by:

- Giving models a custom system prompt and a default tool call format to use
- Converting `tool` role messages to the `user` role so that chat templates without the `tool` role are compatible
- Converting `assistant` role `tool_calls` into the default tool call format

Results will vary by model.

You can see the default format by running `lms log stream` in your terminal, then sending a chat completion request with `tools` to a model that doesn't have Native tool use support. The default format is subject to change.

► Expand to see example of default tool use format

All models that don't have native tool use support will have default tool use support.

Example using `curl`

This example demonstrates a model requesting a tool call using the `curl` utility.

To run this example on Mac or Linux, use any terminal. On Windows, use [Git Bash](#).

```
1 curl http://localhost:1234/v1/chat/completions \  
2   -H "Content-Type: application/json" \  
3   -d '{  
4     "model": "lmstudio-community/qwen2.5-7b-instruct",  
5     "messages": [{"role": "user", "content": "What dell products do you have under $50 in  
6     "tools": [  
7       {  
8         "type": "function",  
9         "function": {  
10          "name": "search_products",  
11          "description": "Search the product catalog by various criteria. Use this whenever
```





Documentation

```
19         "category": {
20             "type": "string",
21             "description": "Product category to filter by",
22             "enum": ["electronics", "clothing", "home", "outdoor"]
23         },
24         "max_price": {
25             "type": "number",
26             "description": "Maximum price in dollars"
27         }
28     },
29     "required": ["query"],
30     "additionalProperties": false
31 }
32 }
33 }
34 ]
35 }'
```

All parameters recognized by `/v1/chat/completions` will be honored, and the array of available tools should be provided in the `tools` field.

If the model decides that the user message would be best fulfilled with a tool call, an array of tool call request objects will be provided in the response field, `choices[0].message.tool_calls`.

The `finish_reason` field of the top-level response object will also be populated with `"tool_calls"`.

An example response to the above `curl` request will look like:

```
1  {
2    "id": "chatcmpl-gb1t1uqzefudice8ntxd9i",
3    "object": "chat.completion",
4    "created": 1730913210,
5    "model": "lmstudio-community/qwen2.5-7b-instruct",
6    "choices": [
7      {
8        "index": 0,
9        "logprobs": null,
10       "finish_reason": "tool_calls",
11       "message": {
12         "role": "assistant",
13         "tool_calls": [
14           {
15             "id": "365174485",
```





Documentation

```
23     }
24   }
25 ],
26 "usage": {
27   "prompt_tokens": 263,
28   "completion_tokens": 34,
29   "total_tokens": 297
30 },
31 "system_fingerprint": "lmstudio-community/qwen2.5-7b-instruct"
32 }
```

In plain english, the above response can be thought of as the model saying:

"Please call the `search_products` function, with arguments:

- 'dell' for the `query` parameter,
- 'electronics' for the `category` parameter
- '50' for the `max_price` parameter

and give me back the results"

The `tool_calls` field will need to be parsed to call actual functions/APIs. The below examples demonstrate how.

Examples using python

Tool use shines when paired with program languages like python, where you can implement the functions specified in the `tools` field to programmatically call them when the model requests.

Single-turn example

Below is a simple single-turn (model is only called once) example of enabling a model to call a function called `say_hello` that prints a hello greeting to the console:

`single-turn-example.py`



Documentation

```
7 def say_hello(name: str) -> str:
8     print(f"Hello, {name}!")
9
10 # Tell the AI about our function
11 tools = [
12     {
13         "type": "function",
14         "function": {
15             "name": "say_hello",
16             "description": "Says hello to someone",
17             "parameters": {
18                 "type": "object",
19                 "properties": {
20                     "name": {
21                         "type": "string",
22                         "description": "The person's name"
23                     }
24                 },
25                 "required": ["name"]
26             }
27         }
28     }
29 ]
30
31 # Ask the AI to use our function
32 response = client.chat.completions.create(
33     model="lmstudio-community/qwen2.5-7b-instruct",
34     messages=[{"role": "user", "content": "Can you say hello to Bob the Builder?"}],
35     tools=tools
36 )
37
38 # Get the name the AI wants to use a tool to say hello to
39 # (Assumes the AI has requested a tool call and that tool call is say_hello)
40 tool_call = response.choices[0].message.tool_calls[0]
41 name = eval(tool_call.function.arguments)["name"]
42
43 # Actually call the say_hello function
44 say_hello(name) # Prints: Hello, Bob the Builder!
45
```

Running this script from the console should yield results like:

```
1 → % python single-turn-example.py
2 Hello, Bob the Builder!
```



Play around with the name in



Documentation

Multi-turn example

Now for a slightly more complex example.

In this example, we'll:

1. Enable the model to call a `get_delivery_date` function
2. Hand the result of calling that function back to the model, so that it can fulfill the user's request in plain text

► [multi-turn-example.py](#) (click to expand)

Running this script from the console should yield results like:

```
1  → % python multi-turn-example.py
2
3  Model response requesting tool call:
4
5  ChatCompletion(id='chatcmpl-wwpstqqu94go4hvclqnpwn', choices=[Choice(finish_reason='tool_c
6
7  get_delivery_date function returns delivery date:
8
9  2024-11-19 13:03:17.773298
10
11 Final model response with knowledge of the tool call result:
12
13 Your order number 1017 is scheduled for delivery on November 19, 2024, at 13:03 PM.
```

Advanced agent example

Building upon the principles above, we can combine LM Studio models with locally defined functions to create an "agent" - a system that pairs a language model with custom functions to understand requests and perform actions beyond basic text generation.

The agent in the below example can:

1. Open safe urls in your default browser
2. Check the current time
3. Analyze directories in your file system

► [agent-chat-example.py](#) (click to expand)



Documentation

```

5 You: What time is it?
6
7 Assistant: The current time is 14:11:40 (EST) as of November 6, 2024.
8
9 You: What time is it now?
10
11 Assistant: The current time is 14:13:59 (EST) as of November 6, 2024.
12
13 You: Open lmstudio.ai
14
15 Assistant: The link to lmstudio.ai has been opened in your default web browser.
16
17 You: What's in my current directory?
18
19 Assistant: Your current directory at `/Users/matt/project` contains a total of 14 files an
20
21 - Files without an extension: 3
22 - `.mjs` files: 2
23 - `.ts` (TypeScript) files: 3
24 - Markdown (`.md`) file: 1
25 - JSON files: 4
26 - TOML file: 1
27
28 The total size of these items is 1,566,990,604 bytes.
29
30 You: Thank you!
31
32 Assistant: You're welcome! If you have any other questions or need further assistance, fee
33
34 You:

```

Streaming

When streaming through `/v1/chat/completions (stream=true)`, tool calls are sent in chunks.

Function names and arguments are sent in pieces via

`chunk.choices[0].delta.tool_calls.function.name` and
`chunk.choices[0].delta.tool_calls.function.arguments`.

For example, to call `get_current_weather(location="San Francisco")`, the streamed `ChoiceDeltaToolCall` in each `chunk.choices[0].delta.tool_calls[0]` object will look like:

```

1 ChoiceDeltaToolCall(index=0, id='814890118', function=ChoiceDeltaToolCallFunction(argu
2 ChoiceDeltaToolCall(index=0, id=None, function=ChoiceDeltaToolCallFunction(arguments='{',
3 ChoiceDeltaToolCall(index=0, id=None, function=ChoiceDeltaToolCallFunction(arguments='loca

```



Documentation

signature for execution.

The below example shows how to create a simple tool-enhanced chatbot through the `/v1/chat/completions` streaming endpoint (`stream=true`).

► `tool-streaming-chatbot.py` (click to expand)

You can chat with the bot by running this script from the console:

```
1 → % python tool-streaming-chatbot.py
2 Assistant: Hi! I am an AI agent empowered with the ability to tell the current time (Type
3
4 You: Tell me a joke, then tell me the current time
5
6 Assistant: Sure! Here's a light joke for you: Why don't scientists trust atoms? Because th
7
8 Now, let me get the current time for you.
9
10 **Calling Tool: get_current_time**
11
12 The current time is 18:49:31. Enjoy your day!
13
14 You:
```

Community

Chat with other LM Studio users, discuss LLMs, hardware, and more on the [LM Studio Discord server](#).

This page's source is available on [GitHub](#)