



AI Agents with LLMs: Using Tool Callings

AI Agent Series

Nguyen-Thuan Duong – TA

Yen-Linh Vu – STA

Anh-Khoi Nguyen– STA



Agency Level

Agency Level	Description	How that's called	Example
-	LLM output has no impact on program flow	Simple processor	<code>process_output(llm_response)</code>
*	LLM output determines basic control flow	Router	<code>if llm_decision(): path_a() else: path_b()</code>
**	LLM output determines function execution	Tool call	<code>run_function(llm_chosen_tool, llm_chosen_args)</code>
***	LLM output controls iteration and program continuation	Multi-step Agent	<code>while llm_should_continue(): execute_next_step()</code>
***	One agentic workflow can start another agentic workflow	Multi-Agent	<code>if llm_trigger(): execute_agent()</code>

Introduction

❖ Overview of AI Agent Series

Day 01: AI Agent Libraries

Day 02:
LLMs Tools Calling

Day 03:
ReAct AI Agent

- High-level understanding of AI Agent.
- Investigate how to build AI Agent applications through popular libraries.

- First low-level understanding of AI Agent.
- Investigate the simplest and early form of AI Agent: Make LLMs using tools to solve tasks.

- Second low-level understanding of AI Agent.
- Investigate the general baseline of current AI Agent: Though-Action-Observe through ReAct AI Agent.

Outline

- Introduction
- Tool Calling in LLMs
- Demo
- Question

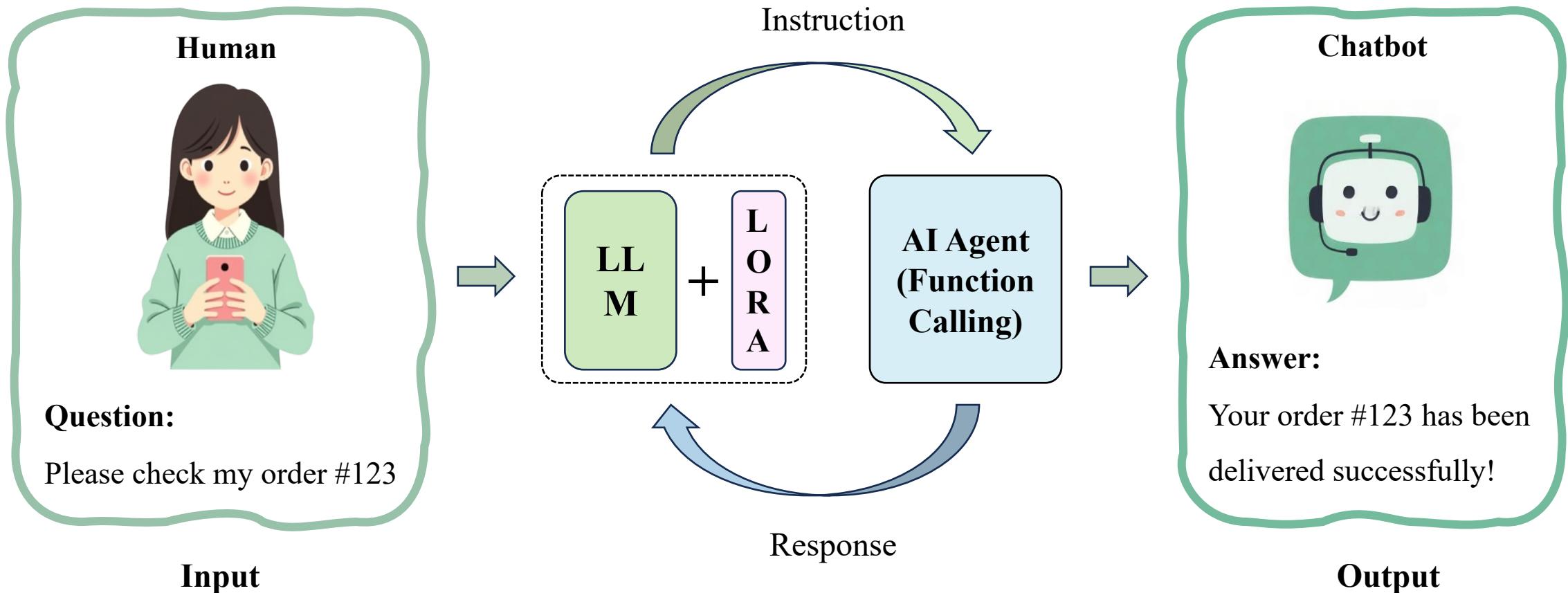
 AI

AI VIET NAM
@aivietnam.edu.vn

Introduction

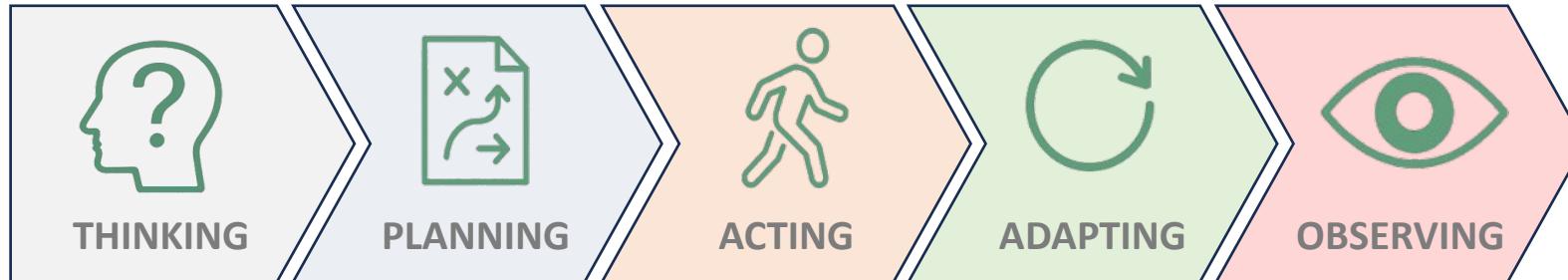
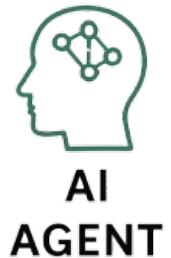
Introduction

❖ Function/Tool Calling Chatbot



Introduction

❖ Agent & Tool Calling



Vs



Introduction

❖ Tool Calling Timeline

HuggingGPT

03/2023

LLM lần đầu đóng vai "bộ não" điều phối: phân tích yêu cầu, lập kế hoạch, và kích hoạt các mô hình AI chuyên sâu (text, image, speech...) từ Hugging Face. Đây là nền tảng cho việc tích hợp công cụ (tool use) vào LLM.

1

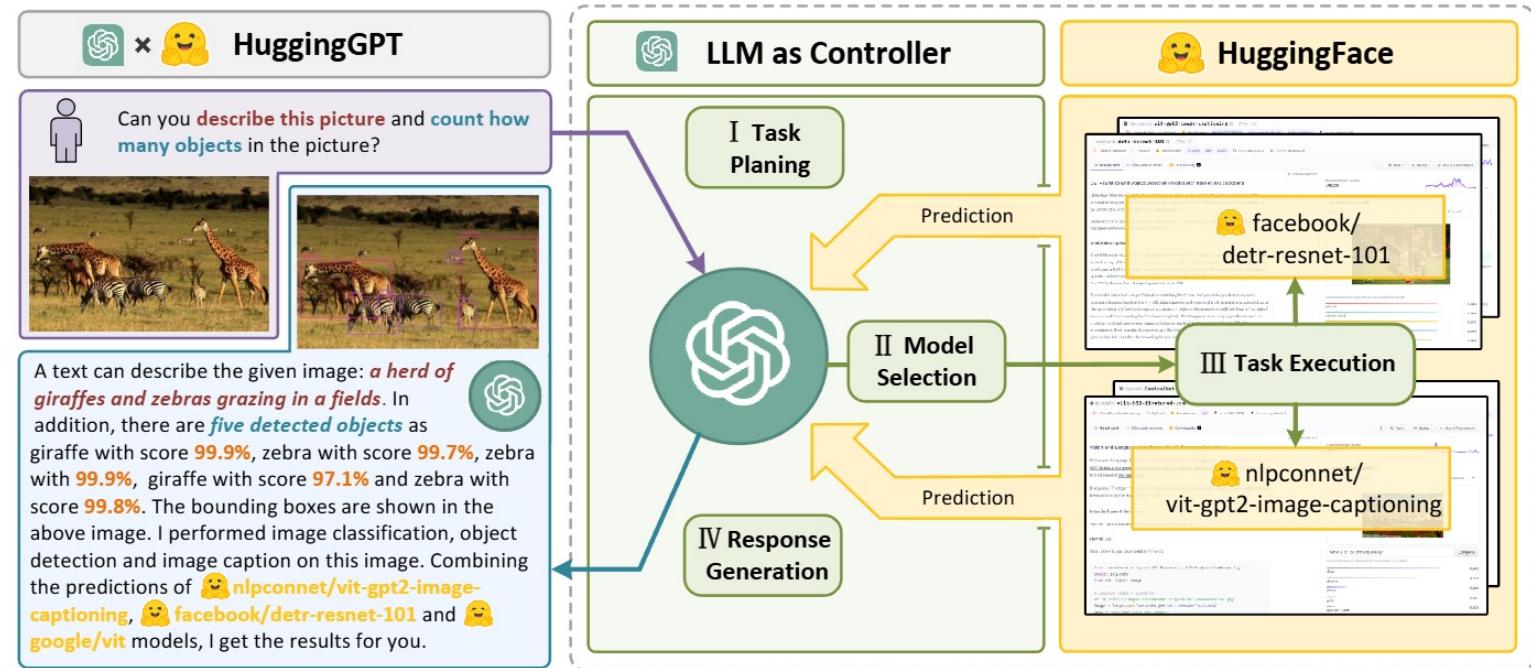


Figure 1: Language serves as an interface for LLMs (e.g., ChatGPT) to connect numerous AI models (e.g., those in Hugging Face) for solving complicated AI tasks. In this concept, an LLM acts as a controller, managing and organizing the cooperation of expert models. The LLM first plans a list of tasks based on the user request and then assigns expert models to each task. After the experts execute the tasks, the LLM collects the results and responds to the user.

Introduction

❖ Tool Calling Timeline

GPT4Tools

LLMs as Tool Makers

05/2023

2

Phát triển từ ý tưởng trên, GPT-4 không chỉ gọi công cụ (tool calling) mà còn tự động tạo lệnh (function generation) dựa trên mô tả. Bước nhảy vọt này giúp LLM chuyển từ thao tác thủ công sang tự học cách ứng dụng công cụ.

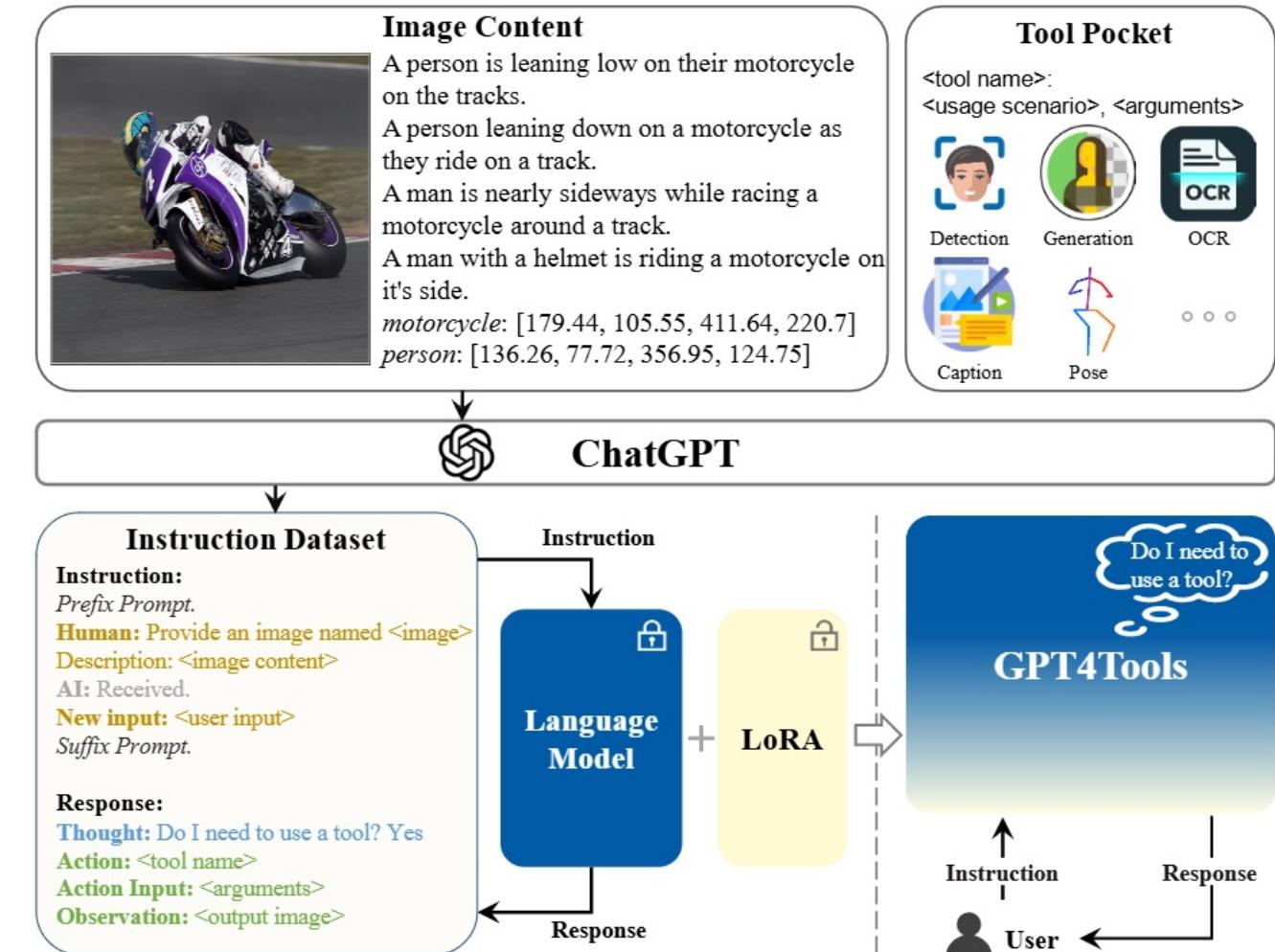
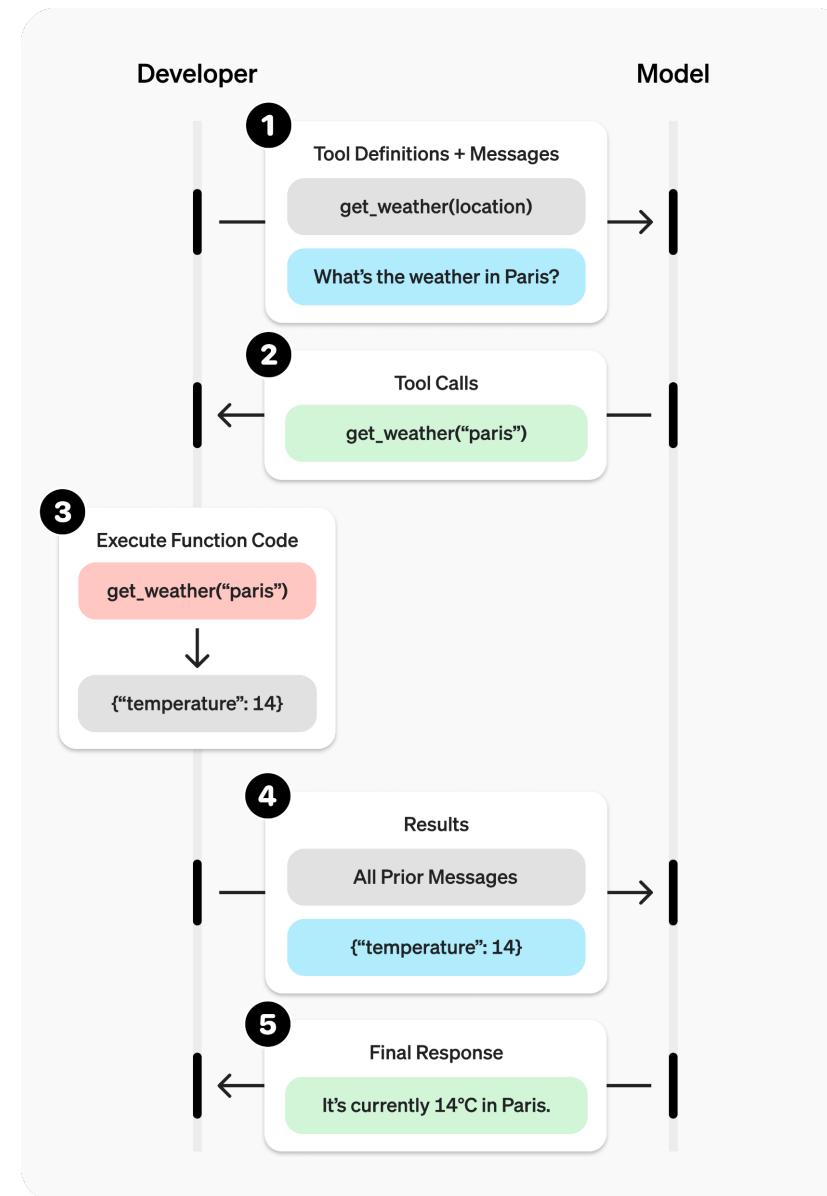
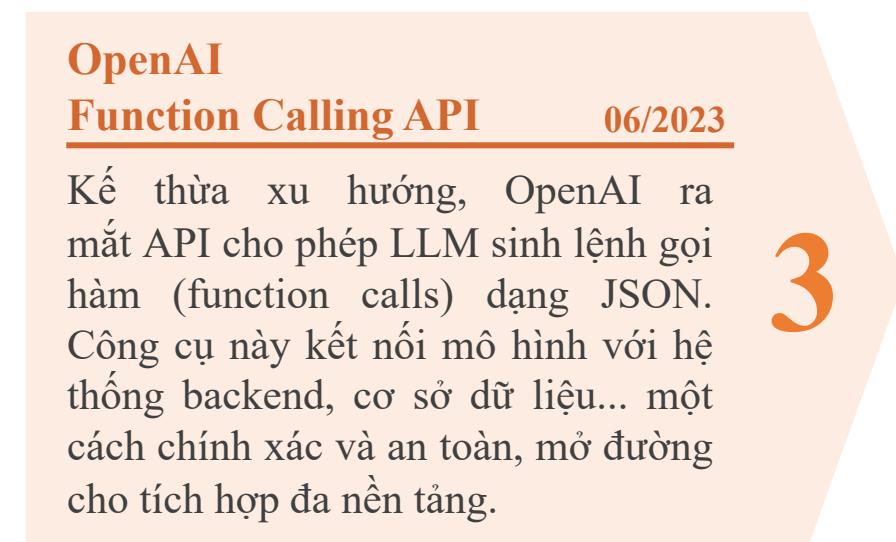


Figure 1: Diagram of the GPT4Tools. We prompt the ChatGPT with image content and definition of tools in order to obtain a tool-related instruction dataset. Subsequently, we employ LoRA [38] to train an open source LLM on the collected instruction dataset, thus adapting the LLM to use tools.

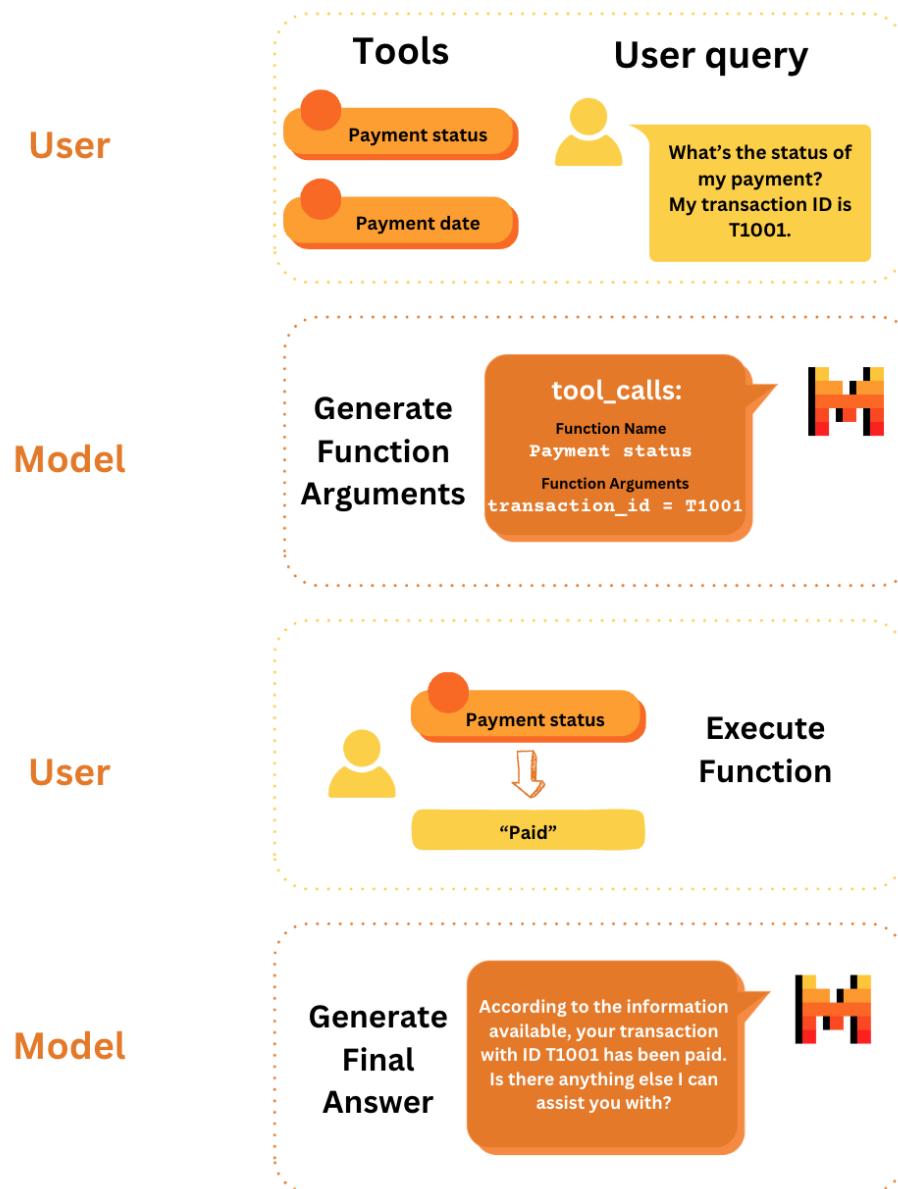
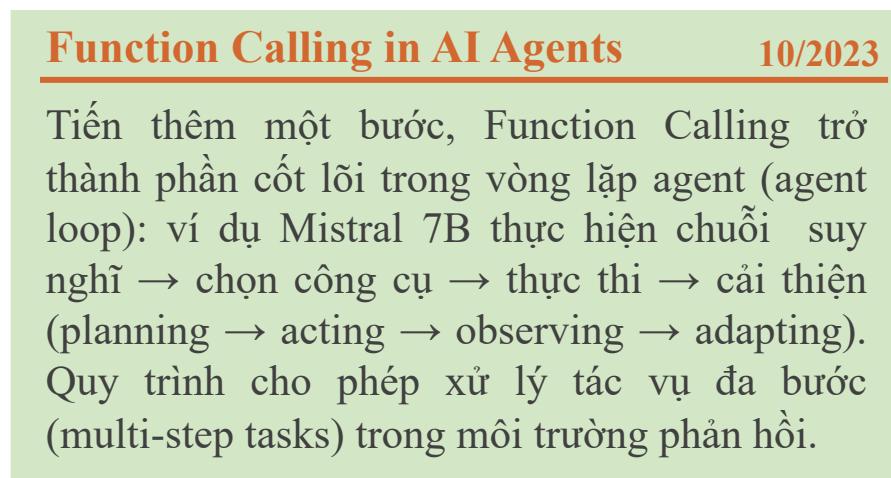
Introduction

❖ Tool Calling Timeline



Introduction

❖ Tool Calling Timeline



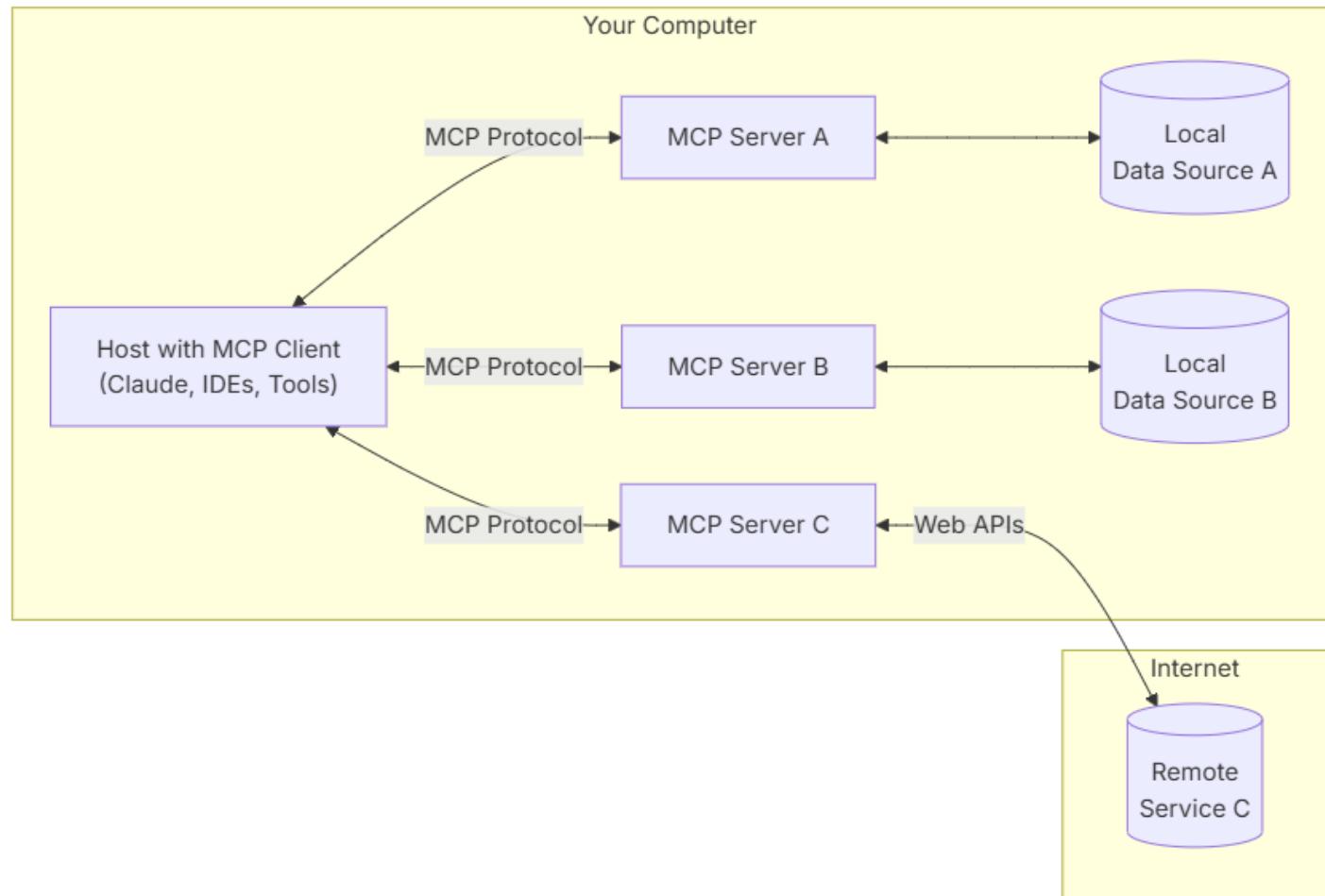
Introduction

❖ Tool Calling Timeline

Anthropic Model Context Protocol 11/2024

Đến cuối năm 2024, MCP cho phép nhiều LLM chia sẻ ngữ cảnh (context sharing) và phối hợp như một hệ đa agent (multi-agent). Đây là bước tiến lớn cho các hệ thống AI phức tạp và linh hoạt hơn.

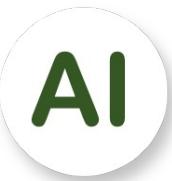
5



Tool Calling in LLMs

Tool Calling

- LLM Router
- Built in tools
- Langchain ToolCalling
- ToolCalling Instruction FT
- ToolCalling Vision Task



AI VIET NAM
@aivietnam.edu.vn

LLM Router



LLM Router

❖ Getting Started

```
● ● ●  
1 def path_a():  
2     print("PATH A → send a welcome e-mail.")  
3  
4 def path_b():  
5     print("PATH B → log a support ticket.")
```

```
● ●  
messages = [  
    {  
        "role": "system",  
        "content": (  
            "You are a decision router.  
            Reply with *exactly one* capital letter: **A** or **B**."  
        ),  
    },  
    {  
        "role": "user",  
        "content": (  
            "A customer wrote: \"The product arrived damaged and I need help.\"  
            If this is a *simple* inquiry, answer **A**.  
            If it is a *problem* requiring human assistance, answer **B**."  
        ),  
    },  
]
```



LLM Router

❖ Model

```
● ● ●  
1 import torch  
2 from transformers import AutoTokenizer, AutoModelForCausalLM  
3  
4  
5 model_name = "meta-llama/Llama-3.2-1B-Instruct"  
6 tokenizer = AutoTokenizer.from_pretrained(model_name)  
7 model = AutoModelForCausalLM.from_pretrained(  
8     model_name,  
9     torch_dtype=torch.bfloat16,  
10    device_map="auto"  
11 )
```

```
● ● ●  
1 prompt = tokenizer.apply_chat_template(  
2     messages,  
3     tokenize=False,  
4     add_generation_prompt=True,  
5 )
```



LLM Router

❖ Decision

```
1 inputs = tokenizer(prompt, return_tensors="pt").to(model.device)
2
3 with torch.no_grad():
4     output_ids = model.generate(
5         **inputs,
6         max_new_tokens=128,
7         pad_token_id=tokenizer.eos_token_id,
8     )
9
10 output = tokenizer.decode(output_ids[0], skip_special_tokens=True).strip()
11 print(f"Decision: {output}")
```

```
decision = output.split()[-1]
print(f"Model decision → {decision}")
```

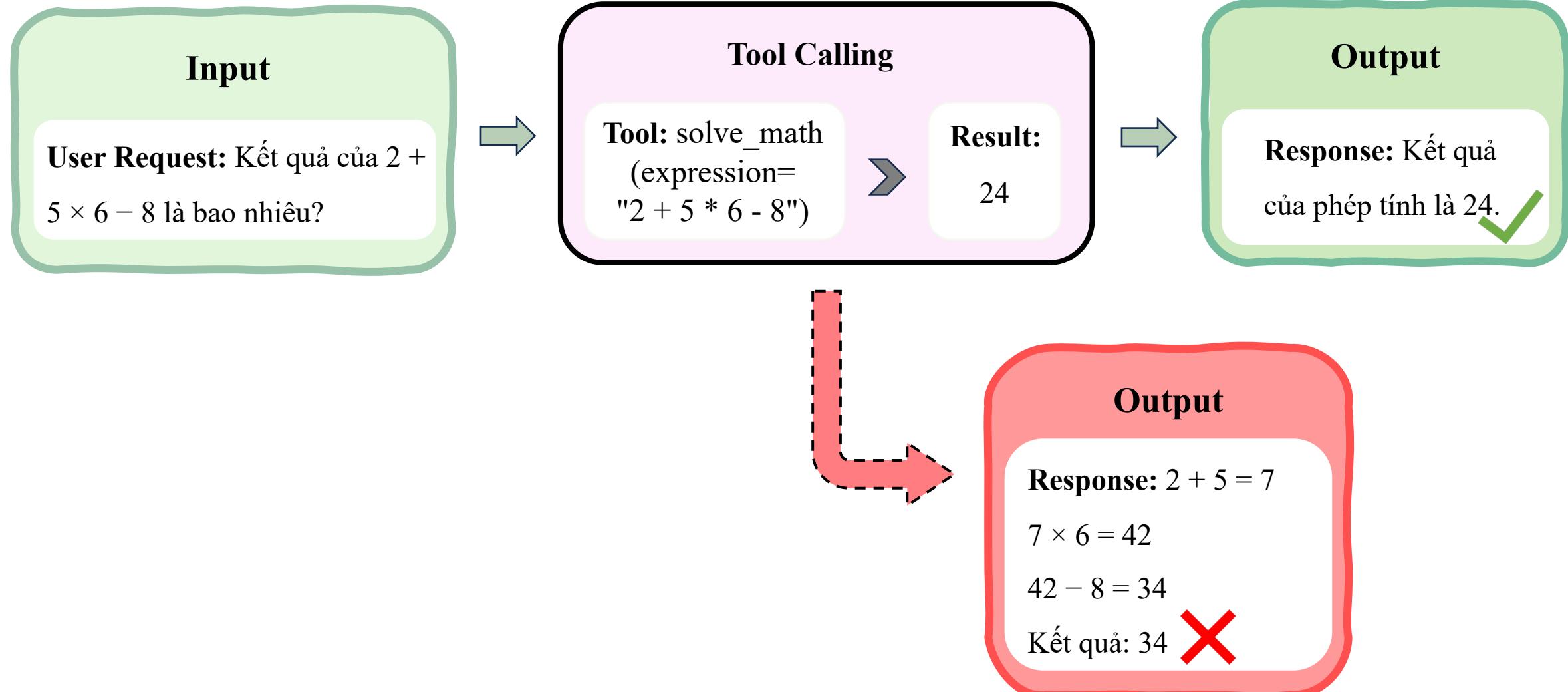
[7]

... Model decision → A

Build in Tools

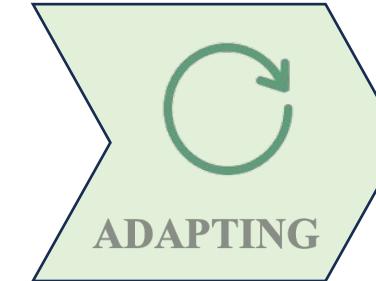
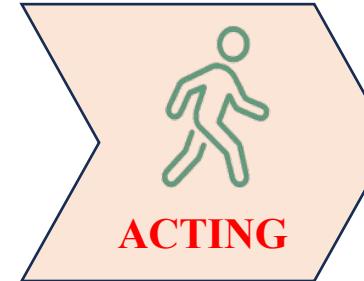
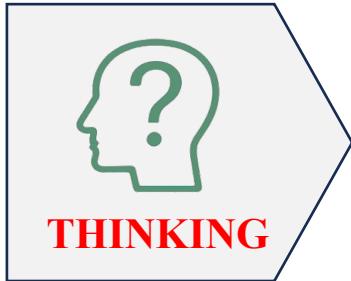
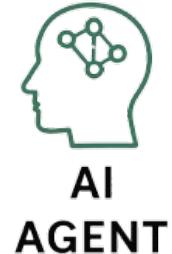
Tool Calling in LLMs

❖ Tool Calling



Tool Calling in LLMs

❖ Tool Calling Workflow



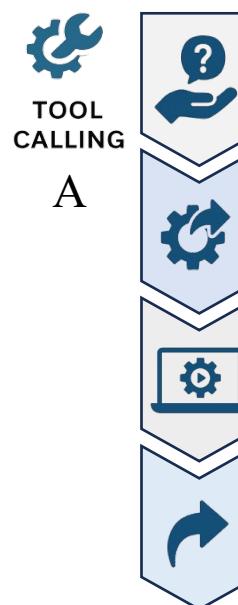
Understand:

Task ?
Goal ?

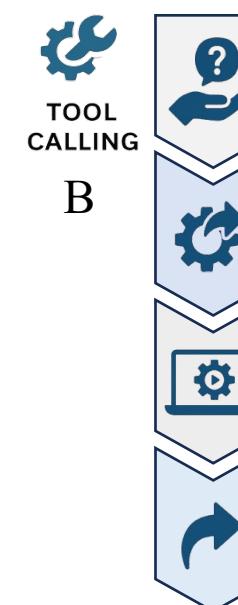
Make a plan:

Step 1: ...
Step 2: ...
...
Step n: ...

Call tool A:



If A no response,
call tool B instead:



Resolved?

If yes, go to next step,
if not, go back to
ADAPTING step or
Re-Planning or
Request more
information from user
(if needed), etc.

Tool Calling in LLMs

❖ Tool Calling in Llama3

Model calls a tool

```
{  
  "name": "get_current_conditions",  
  "parameters": {  
    "location": "San Francisco, CA",  
    "unit": "Fahrenheit"  
  }  
}
```

Function Calling (JSON Schema)

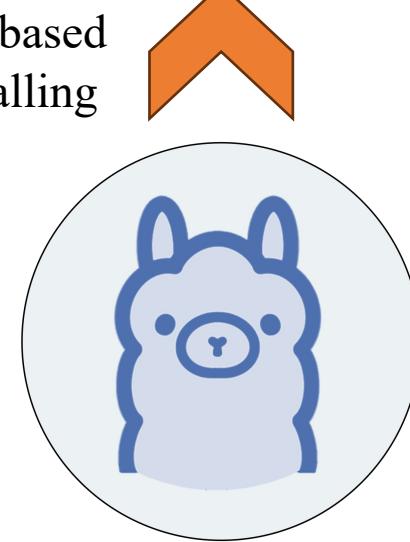
```
{  
  "type": "function",  
  "function": {  
    "name": "get_current_conditions",  
    "description": "Get the current weather  
conditions for a specific location",  
    "parameters": {  
      ...  
    }  
  }  
}
```

Model calls a tool

```
<function=spotify_trending_songs>{"n":"5"}</function><eom_id>
```

JSON based
tool calling

Built-in



User-defined
Custom

Function Calling (JSON Schema)

```
{  
  "type": "function",  
  "function": {  
    "name": "get_current_conditions",  
    "description": "Get the current weather conditions for a specific location",  
    "parameters": {  
      ...  
    }  
  }  
}
```

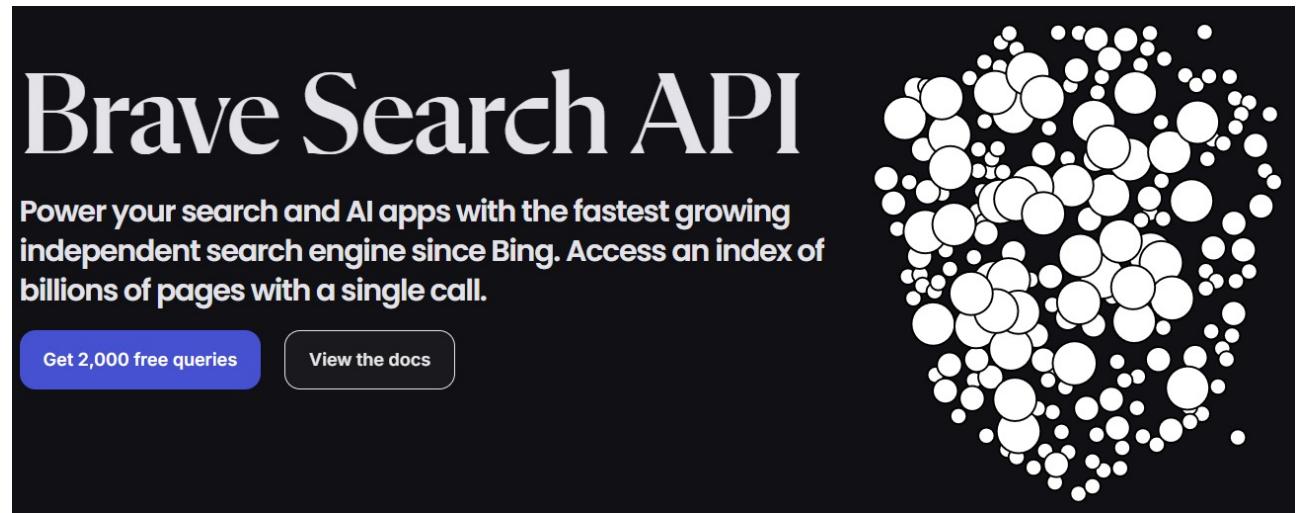
If you choose to call a function ONLY reply in the following format:
<{start_tag}>{function_name}{parameters}{end_tag}
where
start_tag => '<function>'
parameters => a JSON dict with the function argument name as key
and function argument value as value.
end_tag => '</function>'

Tool Calling in LLMs

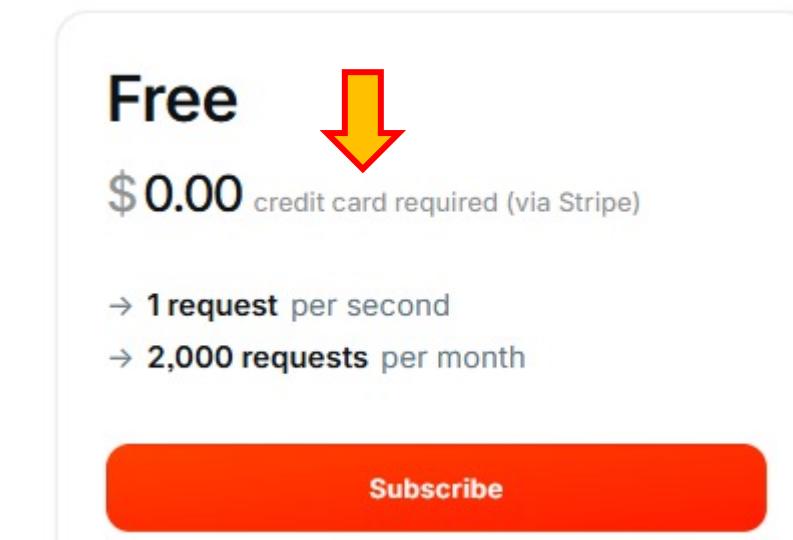
❖ Brave Tool in Llama3



<https://brave.com/search/api/>



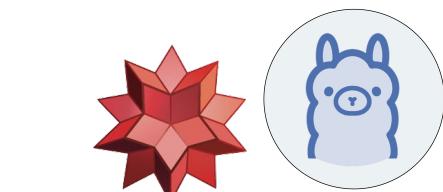
The image shows the homepage of the Brave Search API. It features a large title "Brave Search API" in white on a black background. Below the title is a subtitle: "Power your search and AI apps with the fastest growing independent search engine since Bing. Access an index of billions of pages with a single call." There are two buttons: a blue one labeled "Get 2,000 free queries" and a white one labeled "View the docs". To the right of the main content is a decorative graphic of a brain composed of many white circles of varying sizes against a black background.



A detailed view of the "Free" plan section. It shows a price of "\$0.00" with a note "credit card required (via Stripe)". A yellow arrow points down to the price. Below the price, there are two bullet points: "→ 1 request per second" and "→ 2,000 requests per month". At the bottom is a large orange "Subscribe" button.

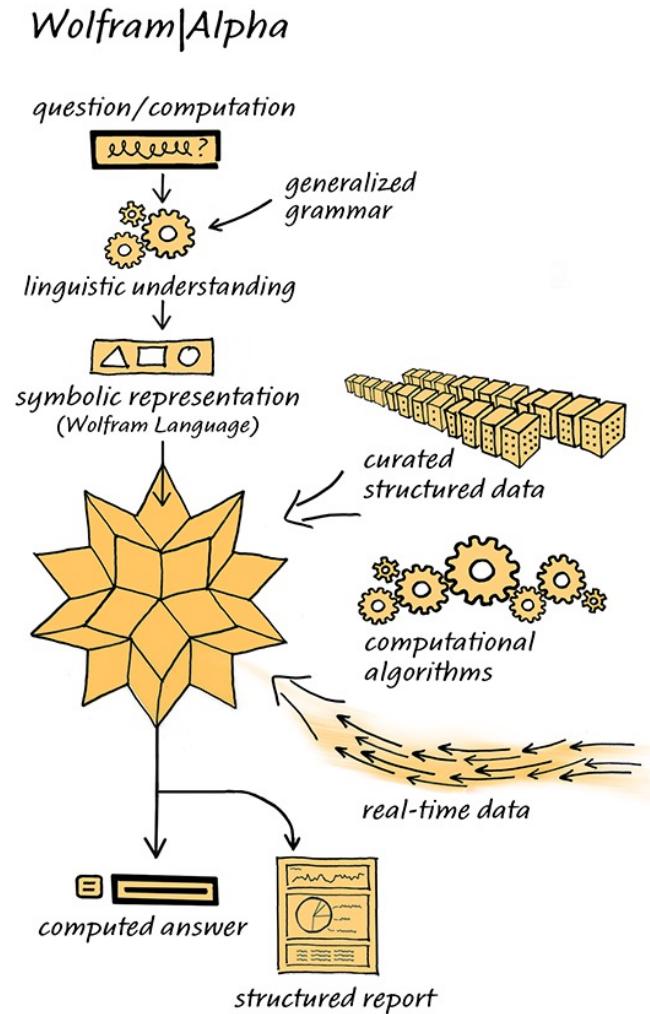
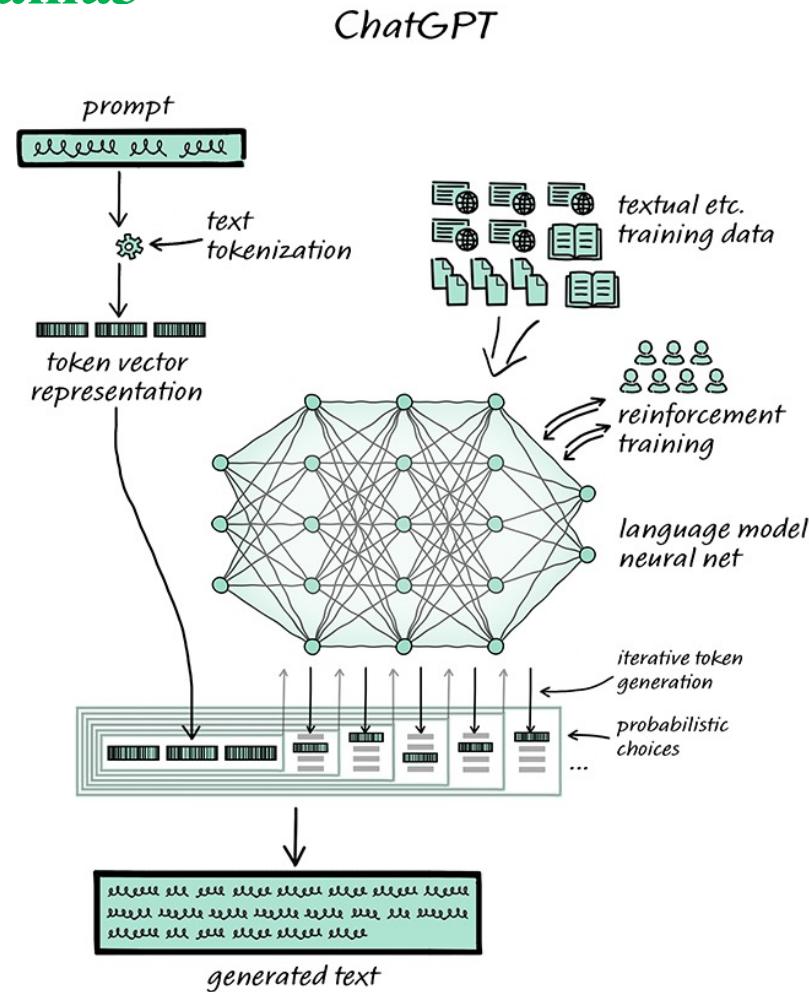
Tool Calling in LLMs

❖ WolframAlpha Tools in Llama3



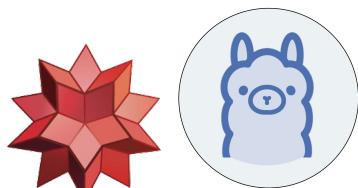
WolframAlpha®

<https://www.wolframalpha.com/>



Tool Calling in LLMs

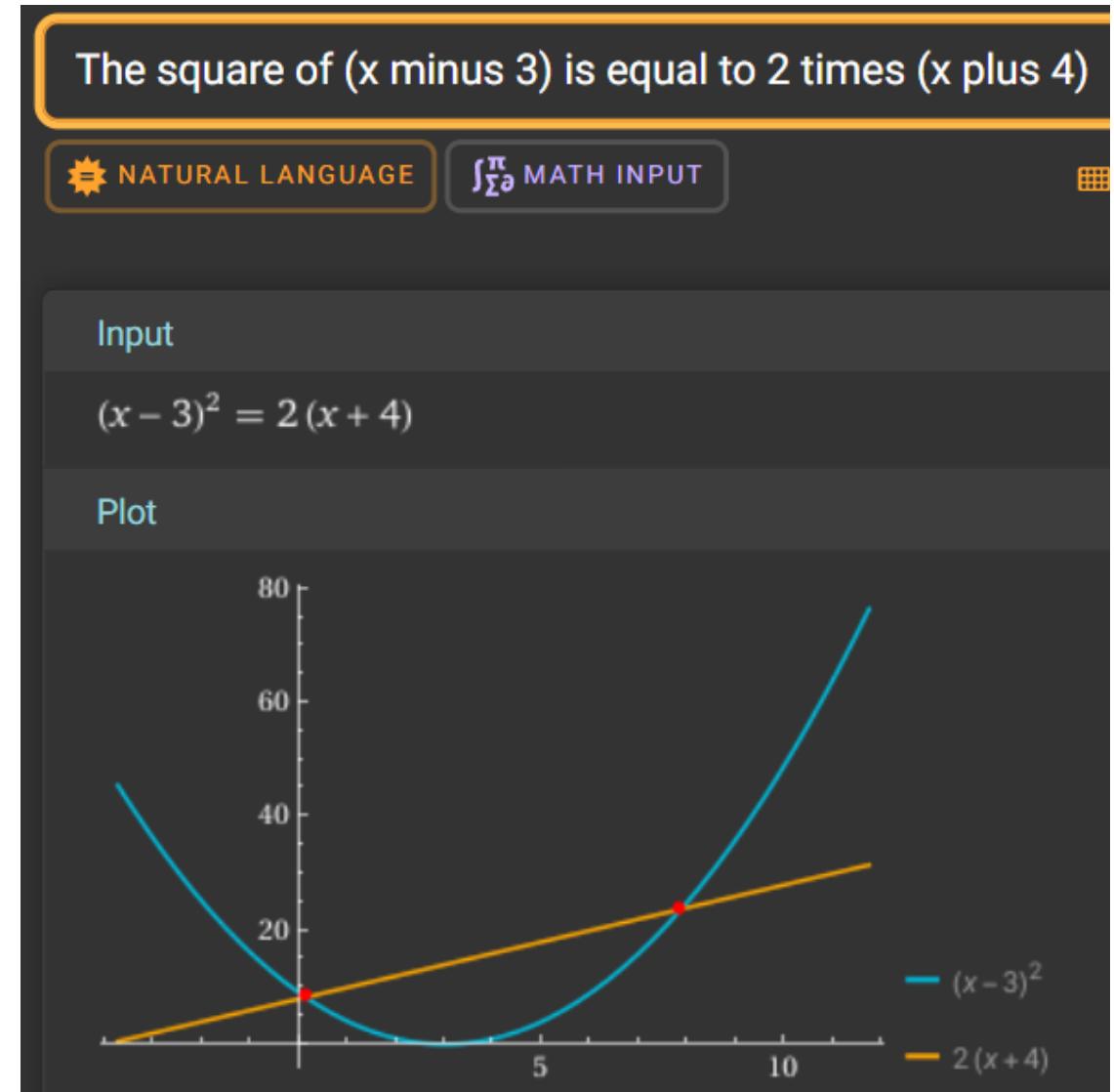
❖ WolframAlpha Tools in Llama3



WolframAlpha®

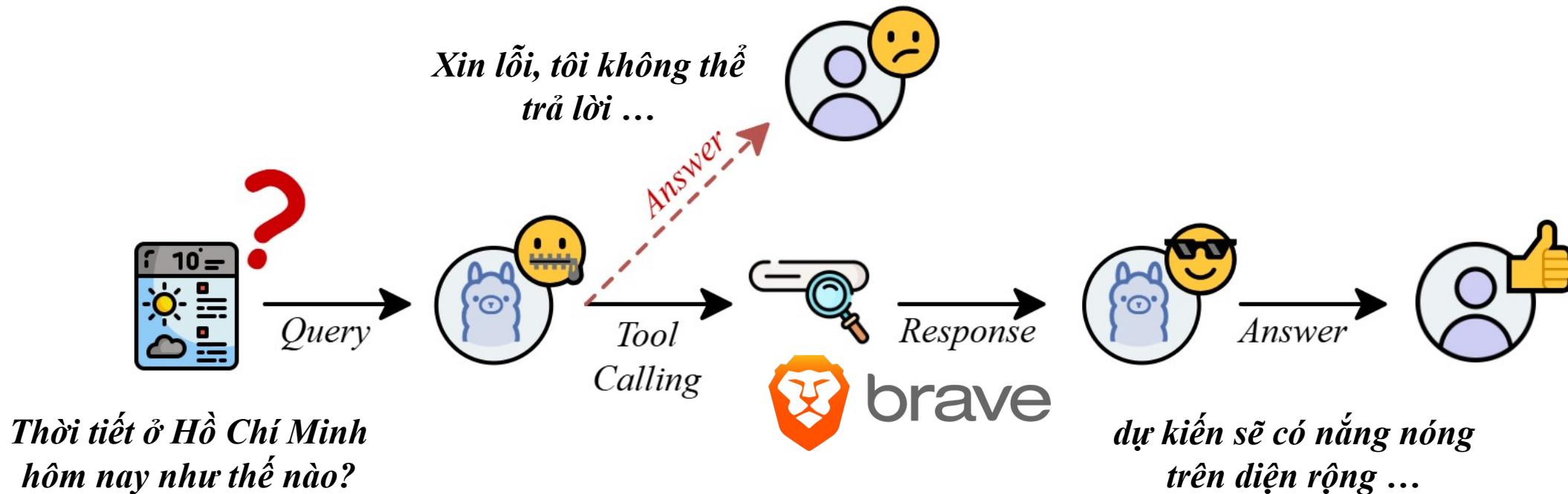
<https://www.wolframalpha.com/>

Solutions

$$x = 4 - \sqrt{15}$$
$$x = 4 + \sqrt{15}$$


Tool Calling in LLMs

❖ Using Brave Tools in Llama3





Tool Calling in LLMs

❖ Llama3 System Prompt for Tool Calling

```
<|begin_of_text|><|start_header_id|>system<|end_header_id|>  
  
Environment: ipython  
Tools: brave_search, wolfram_alpha  
Cutting Knowledge Date: December 2023  
Today Date: 23 July 2024  
  
You are a helpful assistant<|eot_id|>
```



You are a personal math tutor. When asked a math question, write and run code to answer the question.

Great! I'm now your personal math tutor, and I can use code to help solve math problems step by step.

Example Question:
"What is the square root of 245, rounded to 3 decimal places?"

Let's solve it with code:

Analyzed ^

```
python
import math

# Calculate the square root of 245 and round to 3 decimal places
number = 245
square_root = round(math.sqrt(number), 3)
square_root
```

Always show details Copy

Result
15.652

✓ The square root of 245, rounded to 3 decimal places, is 15.652.



Tool Calling in LLMs

❖ Illustration of calling tool

Input

```
<|begin_of_text|><|start_header_id|>system<|end_header_id|>

Environment: ipython
Tools: brave_search, wolfram_alpha
Cutting Knowledge Date: December 2023
Today Date: 23 July 2024

You are a helpful assistant<|eot_id|>
<|start_header_id|>user<|end_header_id|>

What is the current weather in Menlo Park, California?<|eot_id|>
<|start_header_id|>assistant<|end_header_id|>
```

Output



```
# for Search
<|python_tag|>
brave_search.call(query="...")

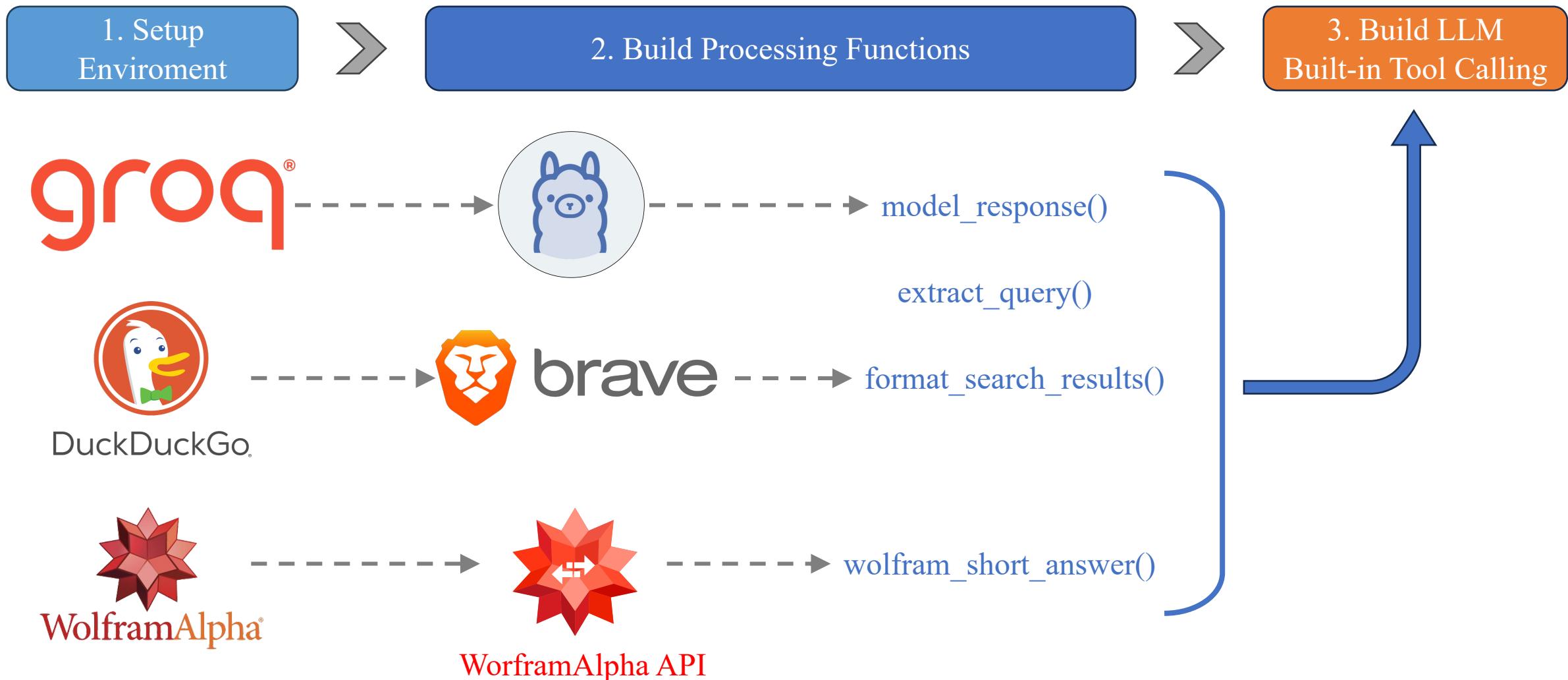
# for Wolfram
<|python_tag|>
wolfram_alpha.call(query="...")
<|eom_id|>
```



WolframAlpha®

Tool Calling in LLMs

❖ Implement Built-in Tool in Llama3



Tool Calling in LLMs

❖ Install & Import libraries

1. Setup
Enviroment



DuckDuckGo.



WolframAlpha®

```
1 !pip install duckduckgo-search groq wolframalpha
```

```
1 import os
2 import re
3 import json
4 import asyncio
5 import aiohttp
6 import requests
7 import nest_asyncio
8 import urllib.parse
9 from groq import Groq
10 from google.colab import userdata
11 from duckduckgo_search import DDGS
12
13 nest_asyncio.apply()
14
15 os.environ["GROQ_API_KEY"] = 'GROQ_API_KEY'
16 APP_ID = 'WF_APPID'
17
18 client = Groq()
```



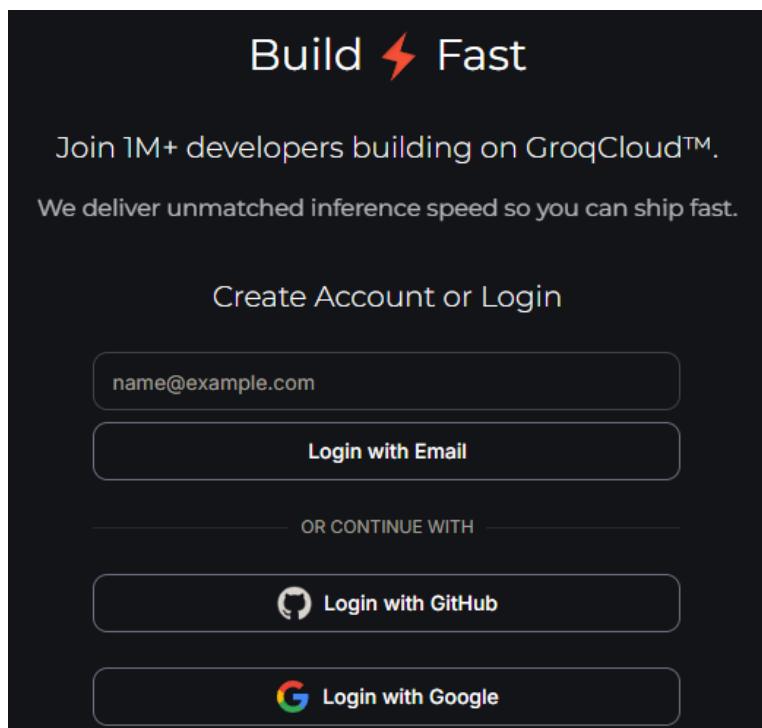
Tool Calling in LLMs

❖ How to get Groq API?

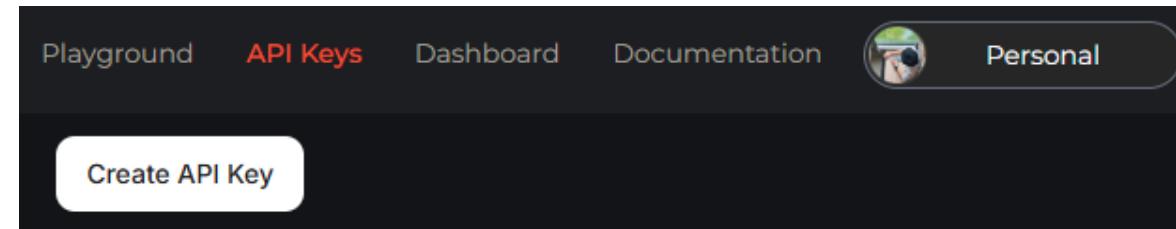


Step 1: Login

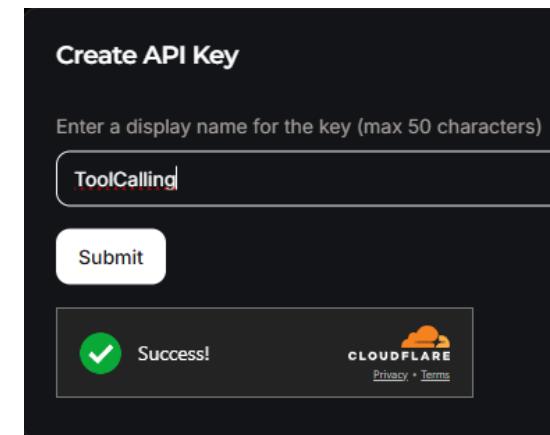
<https://console.groq.com/home>



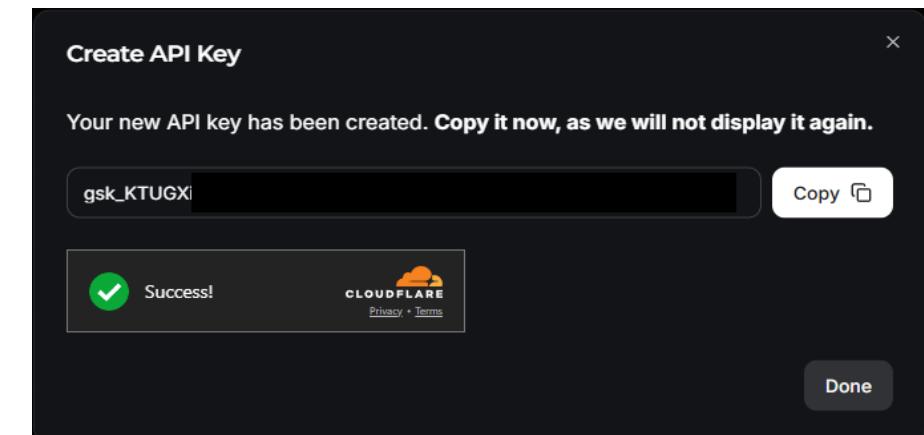
Step 2: API Keys



Step 3: Create API Keys



Step 4: Copy



Step 5: Replace

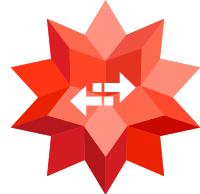
15 os.environ["GROQ_API_KEY"] = 'GROQ_API_KEY'





Tool Calling in LLMs

❖ How to get Wolfram App ID?



Step 1: Login

<https://developer.wolframalpha.com/>

Wolfram|Alpha Developer Portal

Create, manage, and view your App IDs for our Wolfram|Alpha APIs to use in your applications.

Sign in to create or view your App IDs

Sign in

Don't have an account? Create one »

Learn more about Wolfram|Alpha APIs »

Step 2: Get an App ID

DEVELOPER PORTAL

API Access

Manage an App ID to use Wolfram|Alpha APIs in your applications.

Get an App ID

Step 3: Fill & Create App ID

Name *

Description *

API *

Simple API

Submit

Step 4: Copy App ID

Name
test

App ID
JH5YQK-AIVIETNAM2024

Description
test

API
Simple API

Step 5: Replace

`16 APP_ID = 'WF_APPID'`





Tool Calling in LLMs

❖ Model

2. Build Processing Functions



```
1 def model_response(prompt):
2     response = client.chat.completions.create(
3         messages=[
4             {
5                 "role": "user",
6                 "content": prompt,
7             }
8         ],
9         # model="llama-3.3-70b-versatile", # for search
10        model="llama-3.1-8b-instant", # for calculation
11    )
12
13    return response
```

MODEL ID	TOOL USE SUPPORT?
Can try with the following models	
meta-llama/llama-4-scout-17b-16e-instruct	Yes
meta-llama/llama-4-maverick-17b-128e-instruct	Yes
qwen-qwq-32b	Yes
deepseek-r1-distill-qwen-32b	Yes
deepseek-r1-distill-llama-70b	Yes
llama-3.3-70b-versatile	Yes
llama-3.1-8b-instant	Yes
gemma2-9b-it	Yes

<https://console.groq.com/docs/tool-use>

Tool Calling in LLMs

❖ Extract Query

2. Build Processing Functions

```
14 def extract_query(tool_call_text):
15     brave_match = re.search(r'brave_search.call\(query="(.*?)"\)', tool_call_text)
16     if brave_match:
17         return "brave_search", brave_match.group(1)
18     wolfram_match = re.search(r'wolfram_alpha.call\(query="(.*?)"\)', tool_call_text)
19     if wolfram_match:
20         return "wolfram_alpha", wolfram_match.group(1)
21     return None, None
```

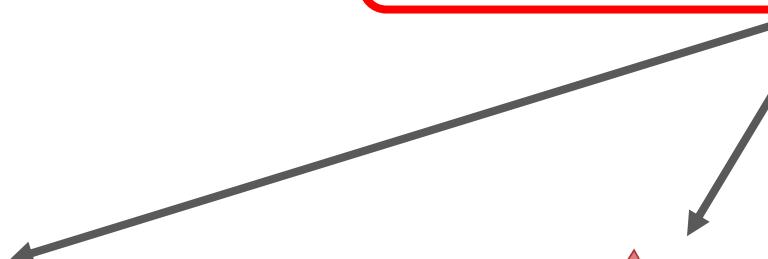
```
<|python_tag|>brave_search.call(query="thời tiết Việt Nam ngày 8/5/2025")
```



DuckDuckGo



WolframAlpha®



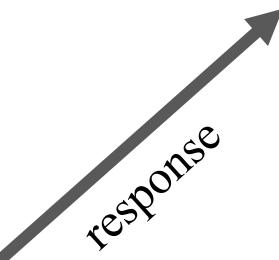
Tool Calling in LLMs

❖ Search Tool

2. Build Processing Functions



DuckDuckGo



```
23 def format_search_results(results):
24     lines = []
25     for i, item in enumerate(results, start=1):
26         title = item.get("title", "No title")
27         href = item.get("href", "No link")
28         body = item.get("body", "No description")
29         lines.append(f"{i}. {title}\n    Description: {body}\n    Link: {href}\n")
30     return "\n".join(lines)
```



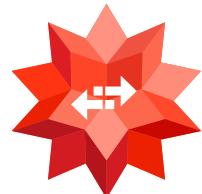
```
==> Tool Responses ==>
```

1. Dự báo thời tiết 8/5/2025: Bắc và Trung ... - Báo VietNamNet
Description: Dự báo thời tiết ngày 8/5/2025 các vùng trên cả nước:
Link: <https://vietnamnet.vn/du-bao-thoi-tiet-8-5-2025-bac-va-trung>
2. Thời tiết hôm nay 8/5/2025: Cả nước nắng nóng diện rộng
Description: Thứ năm, ngày 08/05/2025 06:07 GMT+7 . Thời tiết hôm nay
Link: <https://danviet.vn/thoi-tiet-hom-nay-8-5-2025-ca-nuoc-nang-nong-dien-rong-515025.html>
3. Dự báo thời tiết ngày 8/5/2025: Nắng nóng cục bộ kéo dài
Description: Dự báo thời tiết ngày Hà Nội ngày 8/5. Thời tiết Hà Nội
Link: <https://thoitiet24h.vn/thoi-tiet-hang-ngay/du-bao-thoi-tiet-8-5-2025-nang-nong-cu-bo-kéo-dài-515025.html>

Tool Calling in LLMs

❖ Calculation Tool

2. Build Processing Functions



WorframAlpha API

```
16 APP_ID = userdata.get('WF_APPID1')
```

```
32 async def wolfram_short_answer(query):
33     encoded_query = urllib.parse.quote(query)
34     url = f"https://api.wolframalpha.com/v1/result?i={encoded_query}&appid={APP_ID}"
35     async with aiohttp.ClientSession() as session:
36         async with session.get(url) as response:
37             if response.status == 200:
38                 return await response.text()
39             else:
40                 return f"[Error {response.status}] {await response.text()}"
```



Tool Calling in LLMs

❖ All in One Function

3. Build LLM

Built-in Tool Calling

Default prompt

```
1 async def llama_with_tool(query):
2     PROMPT = f"""
3 <|begin_of_text|><|start_header_id|>system<|end_header_id|>
4
5 Environment: ipython
6 Tools: brave_search, wolfram_alpha
7 Cutting Knowledge Date: December 2023
8 Today Date: 23 July 2024
9
10 You are a helpful assistant specialized in performing web and computation searches. You have access to two tools:
11     1. **brave_search** for general web searches
12     2. **wolfram_alpha** for computation and factual queries
13 Always try to use these tools for informational and computational query questions.<|eot_id|><|start_header_id|>user<|end_header_id|>
14 {query}<|eot_id|><|start_header_id|>assistant<|end_header_id|>
15     """
```



continue ...

Tool Calling in LLMs

❖ All in One Function

3. Build LLM Built-in Tool Calling

```
17 response = model_response(PROMPT)
18 tool_call = response.choices[0].message.content.strip()
19
20 print("== Model Tool Call Output ==\n")
21 print(tool_call)
22
23 tool_name, query_content = extract_query(tool_call)
24 print("\n== Extracted Query Content ==\n")
25 print(query_content)
26
27 if tool_name == "brave_search":
28     tool_response = DDGS().text(query_content, max_results=5)
29     formatted_results = format_search_results(tool_response)
30 elif tool_name == "wolfram_alpha":
31     res = await wolfram_short_answer(query_content)
32     formatted_results = "The exact value of the query is " + res
33 else:
34     formatted_results = "No tool could be identified."
```

<|python_tag|>brave_search.call(query="thời tiết Việt Nam ngày 8/5/2025")



DuckDuckGo.



WorframAlpha API



continue ...

Tool Calling in LLMs

❖ All in One Function

3. Build LLM
Built-in Tool Calling

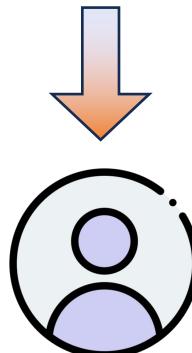
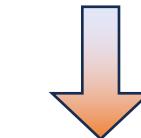


WorframAlpha
API



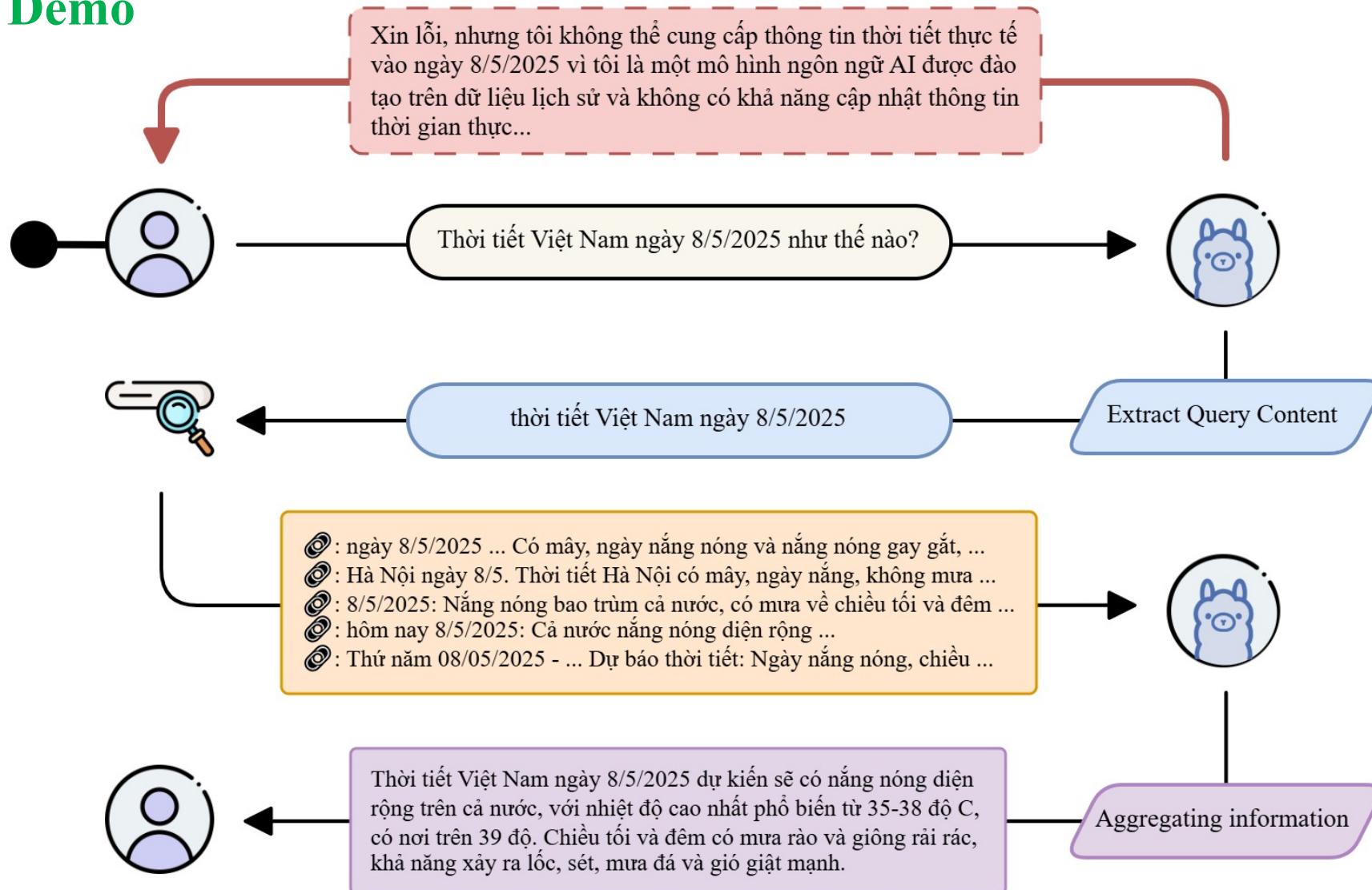
DuckDuckGo

```
36 |     print("\n==== Tool Responses ===\n")
37 |     print(formatted_results)
38 |
39 |     follow_up_prompt = f"Tool responses:{formatted_results}\nBased on these responses, provide a concise result for
40 |     this query.: {query}"
41 |     follow_up_response = model_response(follow_up_prompt)
42 |
43 |     print("\n==== Final Answer ===\n")
44 |     return follow_up_response.choices[0].message.content
```



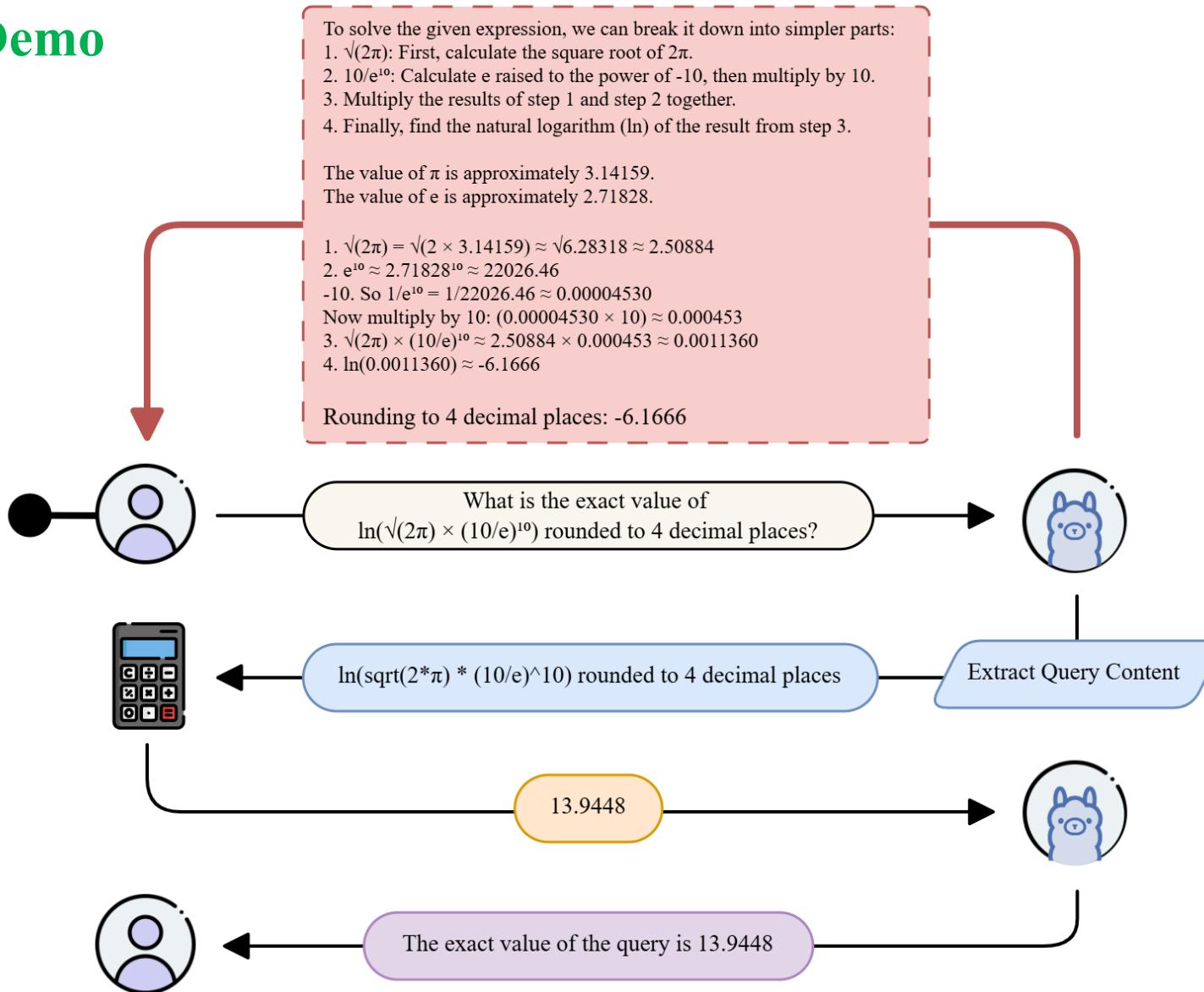
Tool Calling in LLMs

❖ Search Demo



Tool Calling in LLMs

❖ Calculation Demo

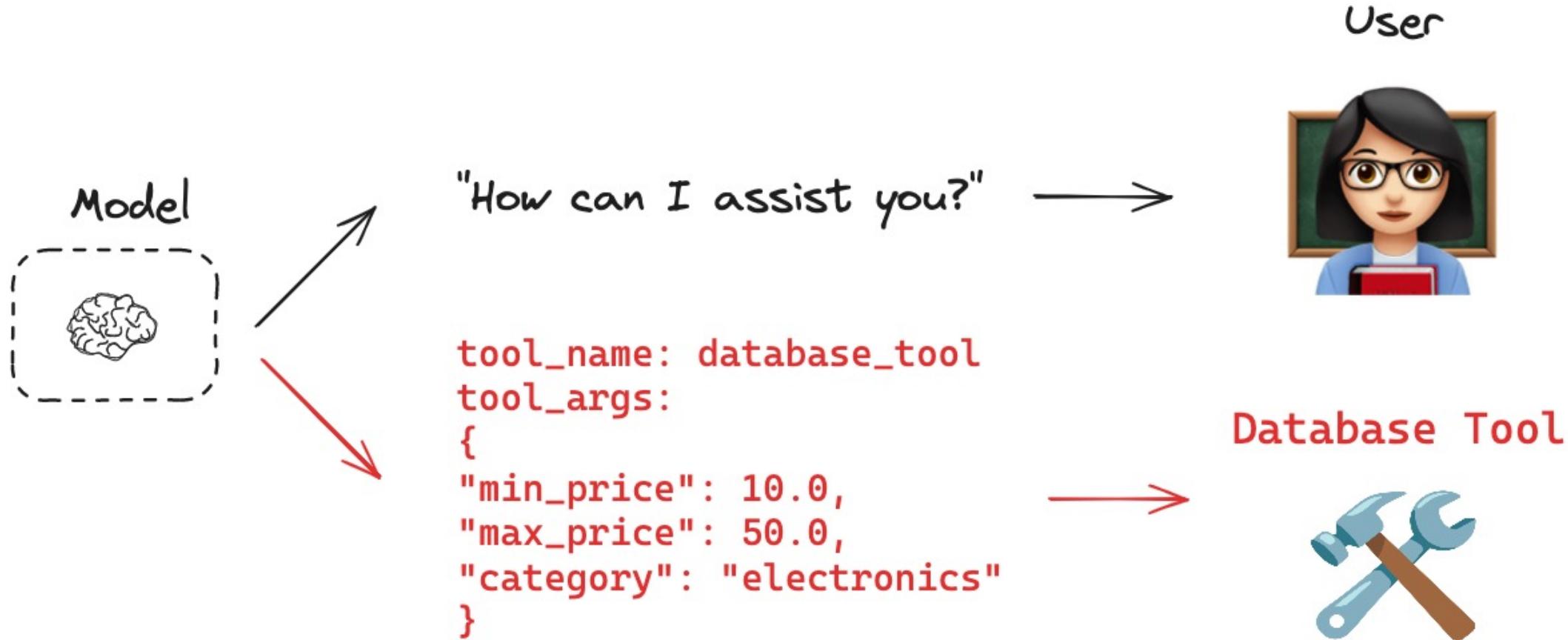




Langchain Tool Calling

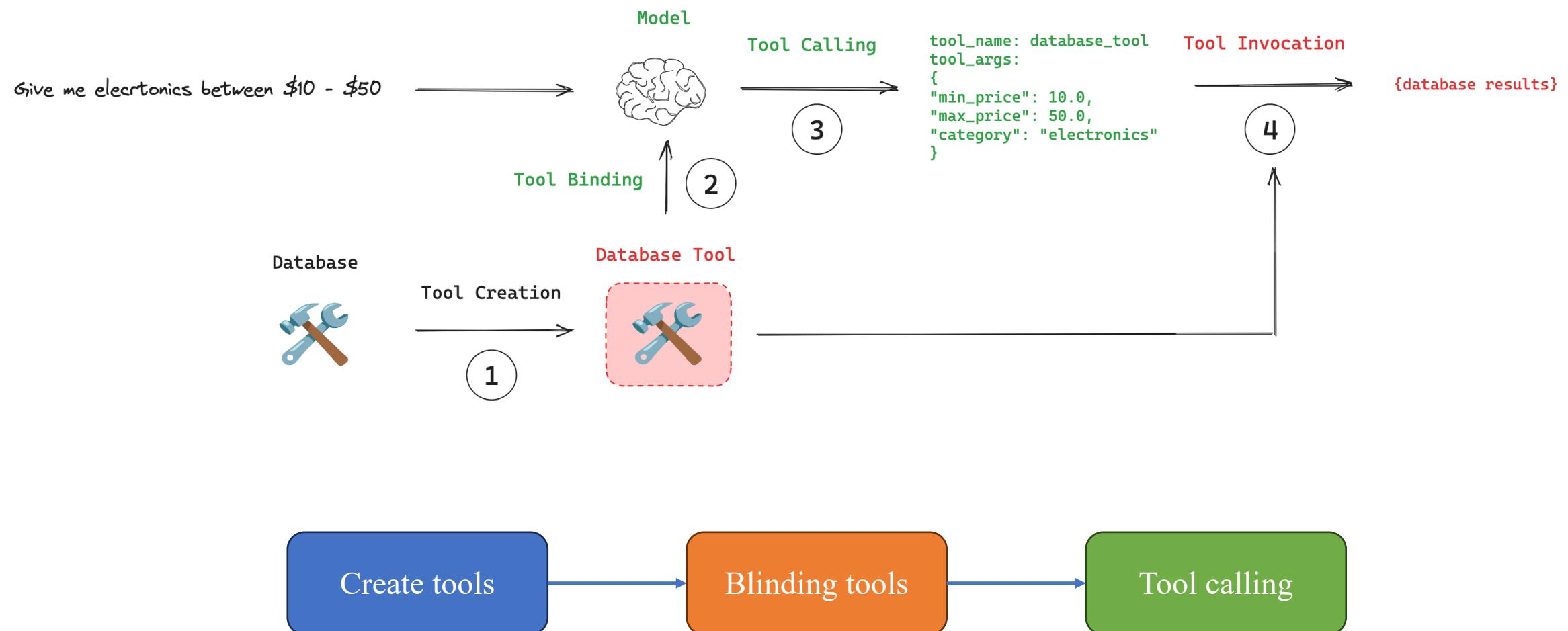
Tool Calling in LLMs

❖ Getting Started



Tool Calling in LLMs

❖ Getting Started





Tool Calling in LLMs

❖ Create tools

using the `@tool` decorator

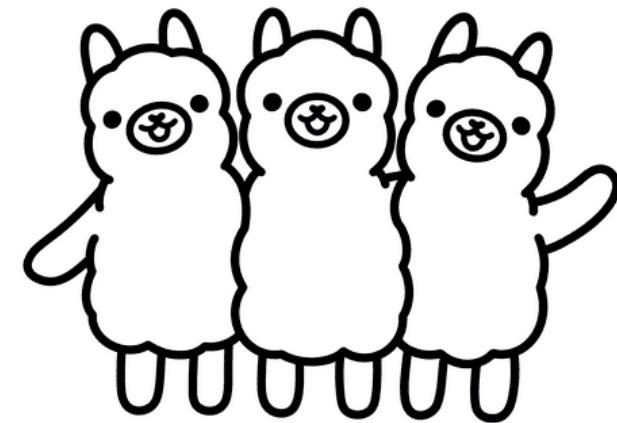
```
1  from langchain_core.tools import tool
2
3  @tool
4  def reverse_string(input_string: str) → str:
5      """Reverse a string"""
6      reversed_string = input_string[::-1]
7
8      return reversed_string
9
10 reversed_string = reverse_string.invoke("AI VIET NAM")
11 print(reversed_string)
```

```
[3]    type(reverse_string)
        ✓  0.0s
...
... langchain_core.tools.structured.StructuredTool
```

Tool Calling in LLMs

❖ Define chatmodel

```
● ● ●  
1  from langchain.chat_models import init_chat_model  
2  
3  model_id = "llama3.2:3b"  
4  temperature = 0.5  
5  max_output_tokens = 512  
6  
7  llm = init_chat_model(  
8      model=model_id,  
9      model_provider="ollama",  
10     temperature=temperature,  
11     max_tokens=max_output_tokens,  
12     format="json"  
13 )
```



ollama run llama3



Tool Calling in LLMs

❖ Blinding tools

```
● ● ●  
1 tool  
2 def search(query: str, top_k=3) → str:  
3     """Search DuckDuckGo for the given query and return the first result."""  
4     wrapper = DuckDuckGoSearchAPIWrapper(region="vn-vi", time="d", max_results=top_k)  
5     search = DuckDuckGoSearchResults(api_wrapper=wrapper, source="news")  
6     results = search.run(query)  
7  
8     return results
```

```
● ● ●  
1 @tool  
2 def multiply(a: int, b: int) → int:  
3     """Multiply a and b."""  
4     return a * b
```

```
● ● ●  
1 tools = [reverse_string, search, multiply]  
2  
3 model_with_tools = llm.bind_tools(tools)
```

Tool Calling in LLMs

❖ Tool Calling



```
1 response = model_with_tools.invoke("Tin tức mới nhất hôm nay")
2 response
```

```
AIMessage(content='', additional_kwargs={}, response_metadata={'model': 'llama3.2:3b', 'created_at': '2025-05-11T11:50:15.402212361Z', 'done': True, 'done_reason': 'stop', 'total_duration': 333856129, 'load_duration': 22755899, 'prompt_eval_count': 271, 'prompt_eval_duration': 104118940, 'eval_count': 21, 'eval_duration': 205630506, 'model_name': 'llama3.2:3b'}, id='run--e55a330c-fa03-481a-8705-2ddc5d939e5b-0', tool_calls=[{'name': 'search', 'args': {'query': 'tin tức mới nhất hôm nay'}, 'id': '47a2f459-bb3c-41a9-a7ce-340360eb0afd', 'type': 'tool_call'}], usage_metadata={'input_tokens': 271, 'output_tokens': 21, 'total_tokens': 292})
```

```
▶ ▾ response.tool_calls
[25]
...
[{'name': 'search',
 'args': {'query': 'tin tức mới nhất hôm nay'},
 'id': '47a2f459-bb3c-41a9-a7ce-340360eb0afd',
 'type': 'tool_call'}]
```

Tool Calling Instruction Fine-tuning

Tool Calling Instruction Fine-tuning

❖ 1. Setup Environment



<https://github.com/bitsandbytes-foundation/bitsandbytes>



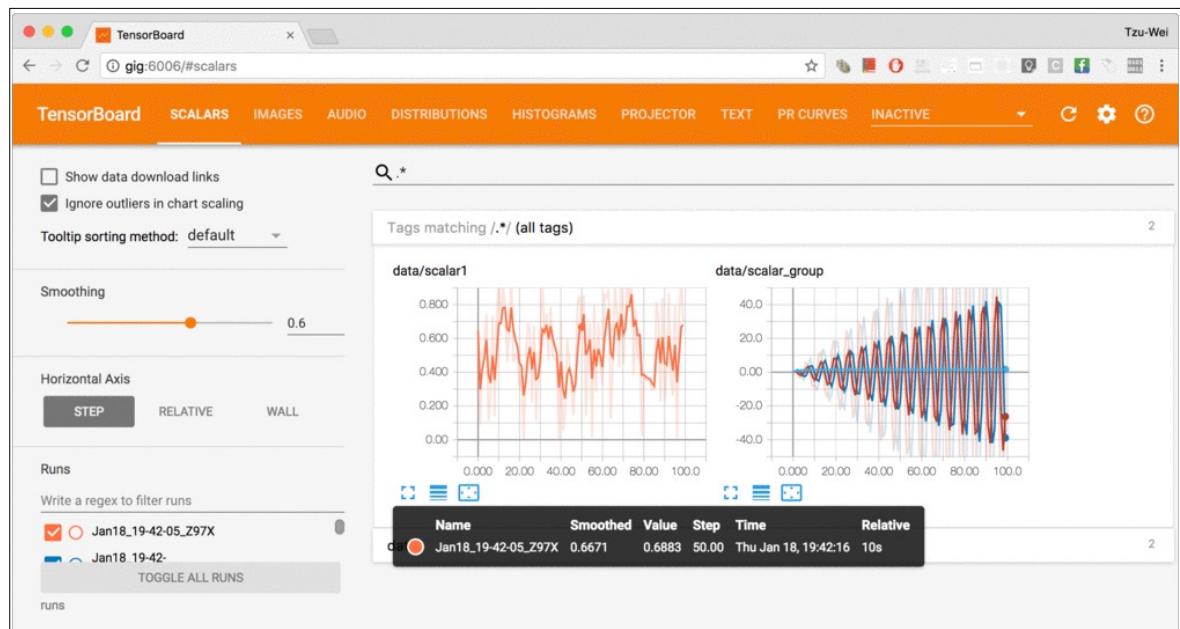
Weights & Biases

<https://github.com/wandb/wandb>

```
1 # !pip install -q -U bitsandbytes
2 # !pip install -q -U peft
3 # !pip install -q -U tensorboardX
4 # !pip install -q wandb
```



<https://github.com/huggingface/peft>



<https://github.com/lanpa/tensorboardX>

Tool Calling Instruction Fine-tuning

❖ 1. Setup Environment



<https://github.com/unslohai/unsloth>



Hugging Face

<https://huggingface.co/>

```
1 import unsloth
2 import os, torch, random, json
3 from datasets import load_dataset
4 from transformers import (
5     TrainingArguments,
6     set_seed
7 )
8 from trl import SFTTrainer, SFTConfig
9 from unsloth import FastLanguageModel
10
11
12 set_seed(59)
13 os.environ["CUDA_VISIBLE_DEVICES"] = "1"
```

Tool Calling Instruction Fine-tuning

❖ 2. Config Model & PEFT

<https://huggingface.co/meta-llama/Llama-3.2-1B>

```
1 model_id = "meta-llama/Llama-3.2-1B-Instruct"
2 max_seq_length = 2048
3 model, tokenizer = FastLanguageModel.from_pretrained(
4     model_id,
5     max_seq_length = max_seq_length,
6     load_in_4bit = True,
7 )
```

⚠ You need to agree to share your contact information to access this model

The information you provide will be collected, stored, processed and shared in accordance with the [Meta Privacy Policy](#).

LLAMA 3.2 COMMUNITY LICENSE AGREEMENT

Llama 3.2 Version Release Date: September 25, 2024

"Agreement" means the terms and conditions for use, reproduction, distribution and modification of the Llama Materials set forth herein.

"Documentation" means the specifications, manuals and documentation accompanying Llama 3.2 distributed by Meta at <https://llama.meta.com/doc/overview>.

"Licensee" or "you" means you, or your employer or any other person or entity (if you are entering into this Agreement on such person o...)

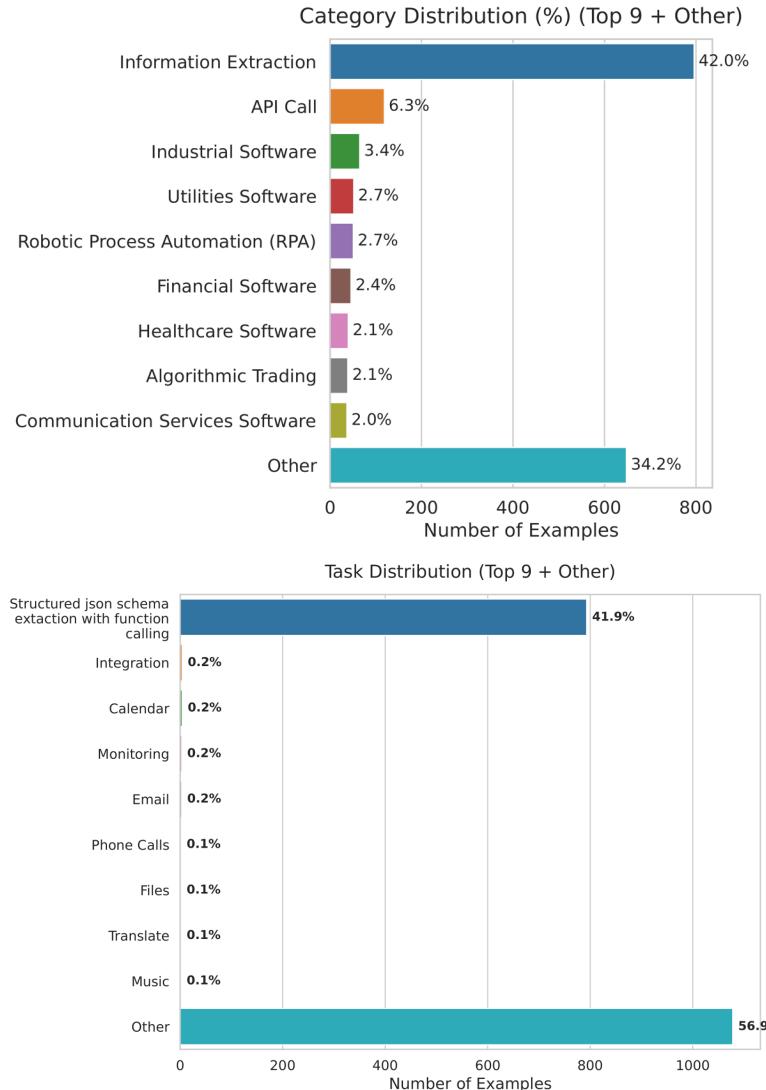
▼ Expand to review

▼ Expand to review and access

```
1 model = FastLanguageModel.get_peft_model(
2     model,
3     r = 16,
4     lora_alpha = 32,
5     lora_dropout = 0.05,
6     target_modules = ["q_proj", "k_proj", "v_proj", "o_proj",
7                       "gate_proj", "up_proj", "down_proj"],
8     bias = "none"
9 )
```

Tool Calling Instruction Fine-tuning

❖ 3. Build Function Calling Dataset



<https://huggingface.co/datasets/NousResearch/hermes-function-calling-v1>

Dataset Viewer

Subset (5) func_calling_singleturn · 1.89k rows

Split (1) train · 1.89k rows

Search this dataset

id	conversations	category	subcategory	task
36	36	3	3	63 values
85f6c398-69c7-4df2-aed1-...	[{ "from": "system", "value": "You are a function calling AI model. You are provided with function..."}	IoT and Home Automation	Security Camera Management	View and Manage Security Camera...
89ef3c87-66bd-46ee-9297-...	[{ "from": "system", "value": "You are a function calling AI model. You are provided with function..."}	IoT and Home Automation	Smart Home Setup	Set Up a Smart Home System
14657d01-d6d1-46df-8eb1-...	[{ "from": "system", "value": "You are a function calling AI model. You are provided with function..."}	IoT and Home Automation	Thermostat Control	Adjust Smart Thermostat...
c483f963-8a29-4ff0-a684-...	[{ "from": "system", "value": "You are a function calling AI model. You are provided with function..."}	IoT and Home Automation	Voice Commands for Home Tasks	Perform Home Tasks Using Voice...
81ad724a-bb74-420f-8221-...	[{ "from": "system", "value": "You are a function calling AI model. You are provided with function..."}	IoT and Home Automation	Lighting Control	Control Smart Lights in a Home
7d99abac-f27f-4ed2-a9ea-...	[{ "from": "system", "value": "You are a function calling AI model. You are provided with function..."}	IoT and Home Automation	Garage Door Automation	Open and Close Smart Garage Doors

< Previous 1 2 3 ... 19 Next >

Tool Calling Instruction Fine-tuning

❖ 3. Build Function Calling Dataset

```
1 dataset_id = "NousResearch/hermes-function-calling-v1"
2 raw_ds = load_dataset(dataset_id, split="train")
3
4 SYS_PROMPT = (
5     "Cut-off knowledge: Dec 2023\n"
6     "Today: 23 July 2024\n\n"
7     "You are a helpful assistant with tool-calling capabilities."
8 )
9
10 def to_chat(example):
11     msgs = [{"role": "system", "content": SYS_PROMPT}]
12     role = {"system": "system", "human": "user", "gpt": "assistant"}
13     for turn in example["conversations"]:
14         msgs.append({"role": role[turn["from"]], "content": turn["value"]})
15     chat_str = tokenizer.apply_chat_template(
16         msgs,
17         tokenize=False,
18         add_generation_prompt=False
19     )
20     return {"text": chat_str}
21
22 ds = raw_ds.map(to_chat).remove_columns(
23     [c for c in raw_ds.column_names if c != "text"]
24 )
```

```
[{"from": "system", "value": "You are a function calling AI model. You are provided with function signatures within <tools> </tools> XML tags. You may call one or more functions...<tool_call>\n{'arguments': <args-dict>, 'name': <function-name>}'\n</tool_call>\n"}, {"from": "human", "value": "I've recently installed a new security system at my home, and I want to ensure everything is functioning...To summarize, I request the execution of the following functions..."}, {"from": "gpt", "value": "<tool_call>\n{'arguments': {'camera_id': 'front_door', 'stream_quality': '1080p'}, 'name': 'get_camera_live_feed'}\n</tool_call>\n<tool_call>\n{'arguments': {'camera_id': 'front_door', 'duration': 30}, 'name': 'record_camera_feed'}\n</tool_call>\n<tool_call>\n{'arguments': {'camera_id': 'front_garden', 'start_time': '2023-04-22T15:00:00Z', 'end_time': '2023-04-22T17:00:00Z'}, 'name': 'get_recorded_feed'}\n</tool_call>\n"}]
```

Tool Calling Instruction Fine-tuning

❖ 4. Train Function Calling Model

```
1 training_args = TrainingArguments(  
2     output_dir = "./llama3-toolcall-lora",  
3     per_device_train_batch_size = 16,  
4     gradient_accumulation_steps = 2,  
5     num_train_epochs = 1,  
6     learning_rate = 2e-4,  
7     lr_scheduler_type = "cosine",  
8     bf16 = torch.cuda.is_bf16_supported(),  
9     logging_steps = 25,  
10    save_strategy = "epoch",  
11    push_to_hub = False,  
12 )  
13  
14 sft_cfg = SFTConfig(dataset_text_field="text",  
15                         max_seq_length=max_seq_length)  
16  
17 trainer = SFTTrainer(  
18     model          = model,  
19     tokenizer      = tokenizer,  
20     train_dataset  = ds,  
21     args           = training_args,  
22     **sft_cfg.to_dict(),  
23 )  
24 trainer.train()
```

230	1.558800
231	1.410900
232	1.106200
233	1.381900
234	1.385100
235	1.259700
236	1.432600
237	1.279800

```
1 trainer.model.save_pretrained("./llama3-toolcall-lora")  
2 tokenizer.save_pretrained("./llama3-toolcall-lora")
```

```
1 # merged = FastLanguageModel.merge_lora(trainer.model)  
2 # merged.push_to_hub("your-hf-name/llama3-1b-toolcall-merged", private=False)  
3 # tokenizer.push_to_hub("your-hf-name/llama3-1b-toolcall-merged")
```

Tool Calling Instruction Fine-tuning

❖ 5. Inference

```
1 import torch
2 from unsloth import FastLanguageModel
3 from peft import PeftModel
4
5 device = "cuda" if torch.cuda.is_available() else "cpu"
6 max_seq_length = 2048
7
8 def chat(user_msg):
9     prompt = tokenizer.apply_chat_template(
10         [{"role": "user", "content": user_msg}],
11         tokenize=False,
12         add_generation_prompt=True,
13     )
14     inputs = tokenizer(prompt, return_tensors="pt").to("cuda")
15     with torch.no_grad():
16         out = model.generate(
17             **inputs,
18             max_new_tokens=max_seq_length,
19             temperature=0.1
20         )
21     decoded_output = tokenizer.decode(out[0], skip_special_tokens=True)
22     decoded_output = decoded_output.split("assistant")[-1].strip()
23
24     return decoded_output
25
26
27 base_id = "meta-llama/Llama-3.2-1B-Instruct"
28 lora_dir = "./llama3-toolcall-lora"
29
30 model, tokenizer = FastLanguageModel.from_pretrained(
31     base_id,
32     max_seq_length = max_seq_length,
33     load_in_4bit = True
34 )
35
36 model = PeftModel.from_pretrained(model, lora_dir, is_trainable=False)
```

chat_msg = """"

You are a function calling AI model. You are provided with function signatures within <tools> </tools> XML tags.
You may call one or more functions to assist with the user query.

Don't make assumptions about what values to plug into functions.

<tools>

[

```
{'type': 'function',
'function': {
    'name': 'get_camera_live_feed',
    'description': 'Retrieves the live feed from a specified security camera.',
    'parameters': {...}},
{'type': 'function',
'function': {
    'name': 'list_all_cameras',
    'description': 'Lists all the security cameras connected to the home network.',
    'parameters': {...}},
{'type': 'function',
'function': {
    'name': 'record_camera_feed',
    'description': 'Starts recording the live feed from a specified security camera.',
    'parameters': {...}},
{'type': 'function',
'function': {
    'name': 'get_recorded_feed',
    'description': 'Retrieves a previously recorded feed from a specified security camera.',
    'parameters': {...}},
{'type': 'function',
'function': {
    'name': 'pan_tilt_camera',
    'description': 'Controls the pan and tilt functions of a PTZ (Pan-Tilt-Zoom) security camera.',
    'parameters': {...}}}
]
```

</tools>

...



continue ...

Tool Calling Instruction Fine-tuning

❖ 5. Inference

For each function call, return a JSON object with function name and arguments within <tool_call>

</tool_call> tags with the following schema:

```
<tool_call>
{'arguments': <args-dict>, 'name': <function-name>}
</tool_call>
```

```
<|eot_id|><|start_header_id|>user<|end_header_id|>
```

I've recently installed a new security system at my home, and I want to ensure everything is functioning as it should.

Specifically, I'd like to:

1. Check the live feed from the "front_door" camera in 1080p quality.
2. Start a 30-minute recording of the live feed from the same camera.
3. Retrieve the recorded feed from the "front_garden" camera between 15:00 and 17:00 on April 22, 2023.

```
'function': {
    'name': 'get_camera_live_feed',
    'description': 'Retrieves the live feed from a specified security camera.'
'function': {
    'name': 'list_all_cameras',
    'description': 'Lists all the security cameras connected to the home network.'
'function': {
    'name': 'record_camera_feed',
    'description': 'Starts recording the live feed from a specified security camera.'
'function': {
    'name': 'get_recorded_feed',
    'description': 'Retrieves a previously recorded feed from a specified security camera.'
'function': {
    'name': 'pan_tilt_camera',
    'description': 'Controls the pan and tilt functions of a PTZ (Pan-Tilt-Zoom) security camera.'
```

```
41 response = chat(chat_msg)
42 print(response)

<tool_call>
{'arguments': {'camera_id': 'front_door', 'stream_quality': '1080p'}, 'name': 'get_camera_live_feed'}
</tool_call>
<tool_call>
{'arguments': {'camera_id': 'front_door', 'duration': 30}, 'name': 'record_camera_feed'}
</tool_call>
<tool_call>
{'arguments': {'camera_id': 'front_garden', 'start_time': '2023-04-22T15:00:00', 'end_time': '2023-04-22T17:00:00'}, 'name': 'get_recorded_feed'}
</tool_call>
```



AI

AI VIET NAM
@aivietnam.edu.vn

QUIZ

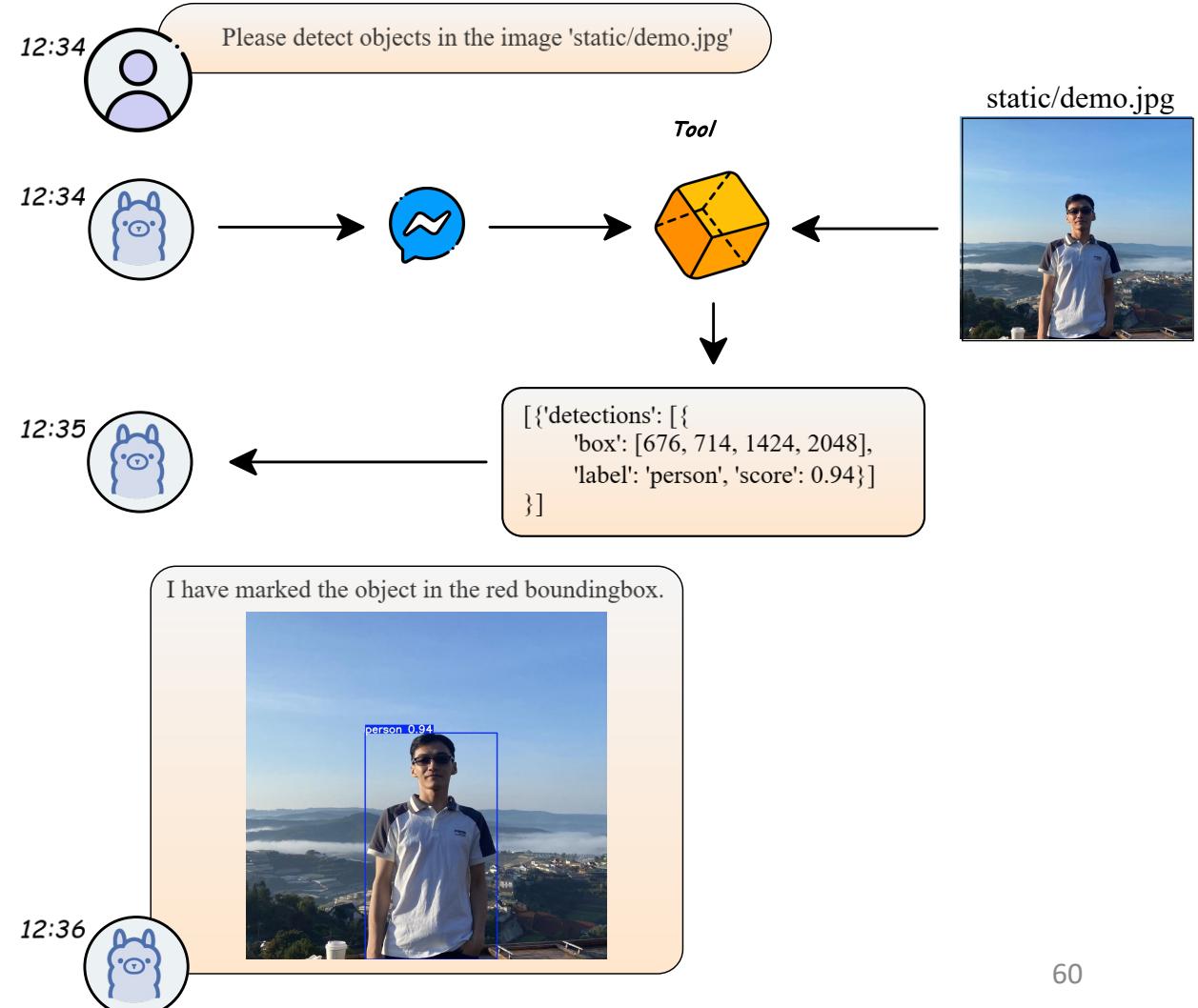
Tool Calling Vision Task

Tool Calling in LLMs

❖ Solving Computer Vision Tasks with Tool Calling

Problem Description: Build a program capable of performing Object Detection and Object Segmentation on a given image.

- **Input:** An image file and a text instruction (e.g., “Detect objects in this image”).
- **Output:** Detected objects with bounding boxes or segmentation masks, returned as structured data.





Tool Calling in LLMs

❖ Step 1: Install and import necessary libraries

```
1 !pip install -q --upgrade ultralytics opencv-python transformers torch vllm unsloth
```

```
1 import os
2 import torch
3 import json
4 import ast
5 import sys
6 from transformers import SamModel, SamProcessor
7 from transformers import AutoModelForCausalLM, AutoTokenizer
8 from transformers import GenerationConfig
9 from PIL import Image
10 from ultralytics import YOLO
11
12
13 os.environ["CUDA_VISIBLE_DEVICES"] = "1"
14 device = torch.device("cuda" if torch.cuda.is_available() else "cpu")
```





Tool Calling in LLMs

❖ Step 2: Define Tools

```
1 detection_model_id = "yolo11n.pt"
2 detection_model = YOLO(detection_model_id)
3 def run_detection(image_path: str, is_visualize: bool = False):
4     """YOLOv11: return list of {box, label, score} for a single image."""
5     results = detection_model(image_path)
6     r = results[0]
7
8     detections = []
9     for box in r.boxes:
10         # box.xyxy is a 1x4 tensor, box.conf is a 1-element tensor, box.cls likewise
11         coords = box.xyxy.cpu().numpy().flatten().tolist()
12         score = float(box.conf.cpu().numpy().item())
13         cls_id = int(box.cls.cpu().numpy().item())
14         detections.append({
15             "box": coords,
16             "label": r.names[cls_id],
17             "score": score,
18         })
19
20     if is_visualize:
21         r.save()
22         r.show()
23
24     return {"detections": detections}
```

The detection function wraps the full YOLOv11 pipeline—input, inference, JSON output, and visualization—into a single callable, letting the LLM use it without handling internal model logic.

Tool Calling in LLMs

❖ Step 2: Define Tools

```
26 segmentation_model_id = "facebook/sam-vit-base"
27 sam_processor = SamProcessor.from_pretrained(segmentation_model_id)
28 sam_model = SamModel.from_pretrained(segmentation_model_id)
29 def run_segmentation(image_path: str):
30     """SAM: return binary masks as nested lists"""
31     img = Image.open(image_path).convert("RGB")
32     inputs = sam_processor(images=img, return_tensors="pt")
33     outputs = sam_model(**inputs)
34
35     masks = outputs.pred_masks.squeeze(0).cpu().detach().numpy().tolist()
36     return {"masks": masks}
```

Like detection, the segmentation function extracts binary masks from SAM and returns pure Python data, enabling seamless data flow between components.

To JSON 

```
def run_detection(image_path: str, is_visualize: bool = False):
    """YOLOv11: return list of {box, label, score} for a
    single image."""

    detections = detection_model(image_path)[0]

    return {"detections": detections}
```

```
{
  "type": "function",
  "function": {
    "name": "run_detection",
    "description": "Detect objects in an image and return
bounding boxes and labels.",
    "parameters": {
      "type": "object",
      "properties": {
        "image_path": {"type": "string", "description":
"Local path to the image file."},
        "is_visualize": {"type": "bool", "description":
"Boolean flag to visualize the detection results. Default is
False."}
      },
      "required": ["image_path"]
    }
  }
}
```



Tool Calling in LLMs

❖ Step 2: Define Tools

```
1 tool_functions = {  
2     "run_detection": run_detection,  
3     "run_segmentation": run_segmentation
```

```
1 functions_lst = [  
2     {  
3         "type": "function",  
4         "function": {  
5             "name": "run_detection",  
6             "description": "Detect objects in an image and return bounding boxes and labels.",  
7             "parameters": {  
8                 "type": "object",  
9                 "properties": {  
10                     "image_path": {"type": "string", "description": "Local path to the image file."},  
11                     "is_visualize": {"type": "bool", "description": "Boolean flag to visualize the detection results. Default is False."}  
12                 },  
13                 "required": ["image_path"]  
14             }  
15         }  
16     },  
17     {  
18         "type": "function",  
19         "function": {  
20             "name": "run_segmentation",  
21             "description": "Segment objects in an image and return binary masks.",  
22             "parameters": {  
23                 "type": "object",  
24                 "properties": {  
25                     "image_path": {"type": "string", "description": "Local path to the image file."}  
26                 },  
27                 "required": ["image_path"]  
28             }  
29         }  
30     }  
31 ]  
32 
```



Tool Calling in LLMs

❖ Step 3: Define system prompt

```
1 SYSTEM_PROMPT = """
2 You are an expert in composing functions. You are given a question and a set of possible functions.
3 Based on the question, you will need to make one or more function/tool calls to achieve the purpose.
4 If none of the function can be used, point it out. If the given question lacks the parameters required by the function,
5 also point it out. You should only return the function call in tools call sections.
6
7 If you decide to invoke any of the function(s), you MUST put it in the format of [func_name1(params_name1=params_value1, params_name2=params_value2...), func_name2(params)]\n
8 You SHOULD NOT include any other text in the response.
9
10 Here is a list of functions in JSON format that you can invoke.\n\n{functions}\n""".format(functions=functions_lst)
```

Once having the function list in JSON format, we define a system prompt to help the model understand the task and available tools.

Tool Calling in LLMs

❖ Step 4: Load Pre-trained model

```
1 model_id = "meta-llama/Llama-3.2-1B-Instruct"
2 llm = AutoModelForCausalLM.from_pretrained(model_id,
3                                              device_map=device,
4                                              torch_dtype=torch.bfloat16)
5 tokenizer = AutoTokenizer.from_pretrained(model_id)
```

[meta-llama/Llama-3.2-1B-Instruct](#) like 924

Text Generation Transformers Safetensors PyTorch

text-generation-inference arxiv:2204.05149 arxiv:2405.16406

Inference Providers SambaNova

Text Generation Reset Chat Examples ▾

Input a message to start chatting with meta-llama/Llama-3.2-1B-Instruct.

What can you do? Make it briefly in one sentence

I can provide information and entertainment on a wide range of topics, from general knowledge and trivia to creative writing and language translation.

Your sentence here... Send

Tool Calling in LLMs

❖ Step 5: Initialize generation config & conversation

```
1 messages = [  
2     {  
3         "role": "system",  
4         "content": SYSTEM_PROMPT  
5     },  
6     {  
7         "role": "user",  
8         "content": "Please segment objects in the image 'static/cat.jpg'."  
9     }  
10 ]
```

```
11 generation_config = GenerationConfig(  
12     max_new_tokens=512,  
13     temperature=0.5,  
14     top_p=0.98,  
15     top_k=50,  
16     repetition_penalty=1.2,  
17     do_sample=True,  
18     num_return_sequences=1,  
19     bos_token_id=tokenizer.bos_token_id,  
20     eos_token_id=tokenizer.eos_token_id,  
21     pad_token_id=tokenizer.eos_token_id,  
22 )  
23
```

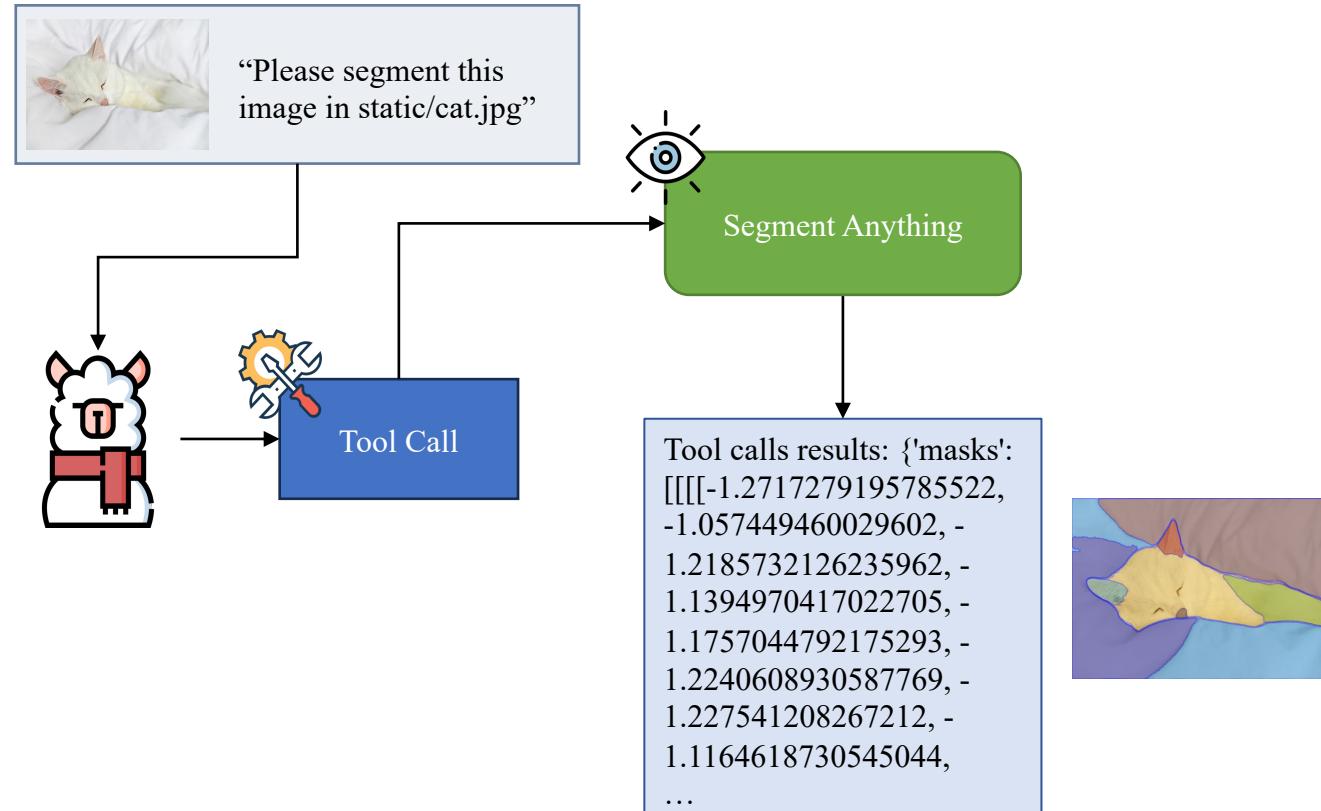
Tool Calling in LLMs

❖ Step 6: Tokenization

Tool Calling in LLMs

❖ Step 7: Run Tool Call

```
43 tool_calls_result = []
44 for call in call_nodes:
45     function_name = call.func.id
46     parameters = {
47         kw.arg: ast.literal_eval(kw.value)
48         for kw in call.keywords
49     }
50
51     result = tool_functions[function_name](**parameters)
52     tool_calls_result.append(result)
53
54 print("Tool calls results:", tool_calls_result)
```



Summarization and Q&A

