

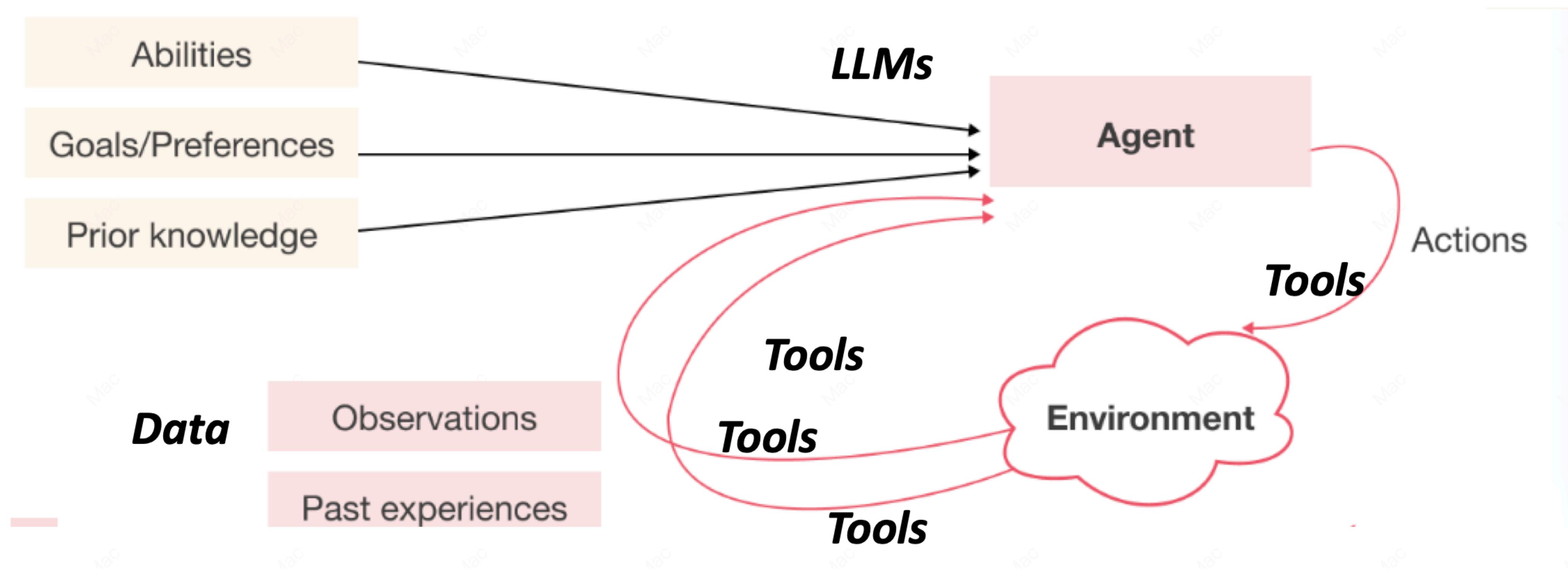
# Language Agents and Tools

Junxian He

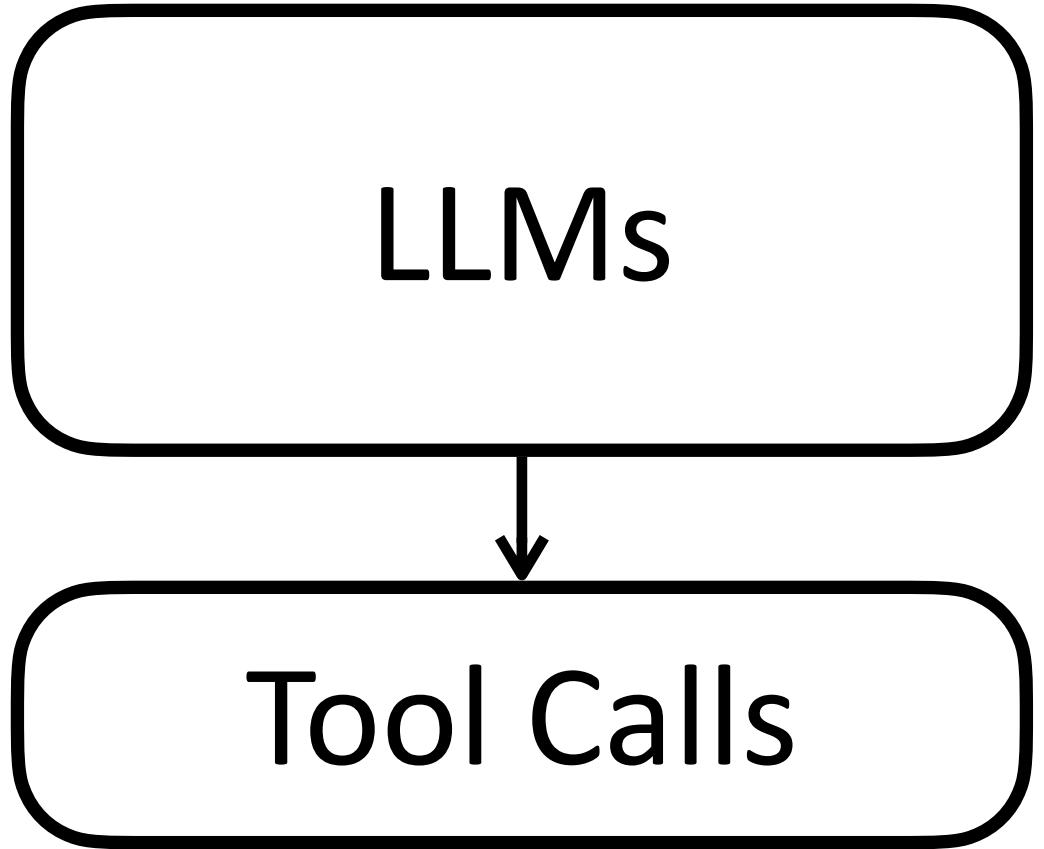
Nov 12, 2025

# Recap: What are Agents

Anything that can be viewed as **perceiving** its environment through sensors and **acting** upon that environment through actuators.

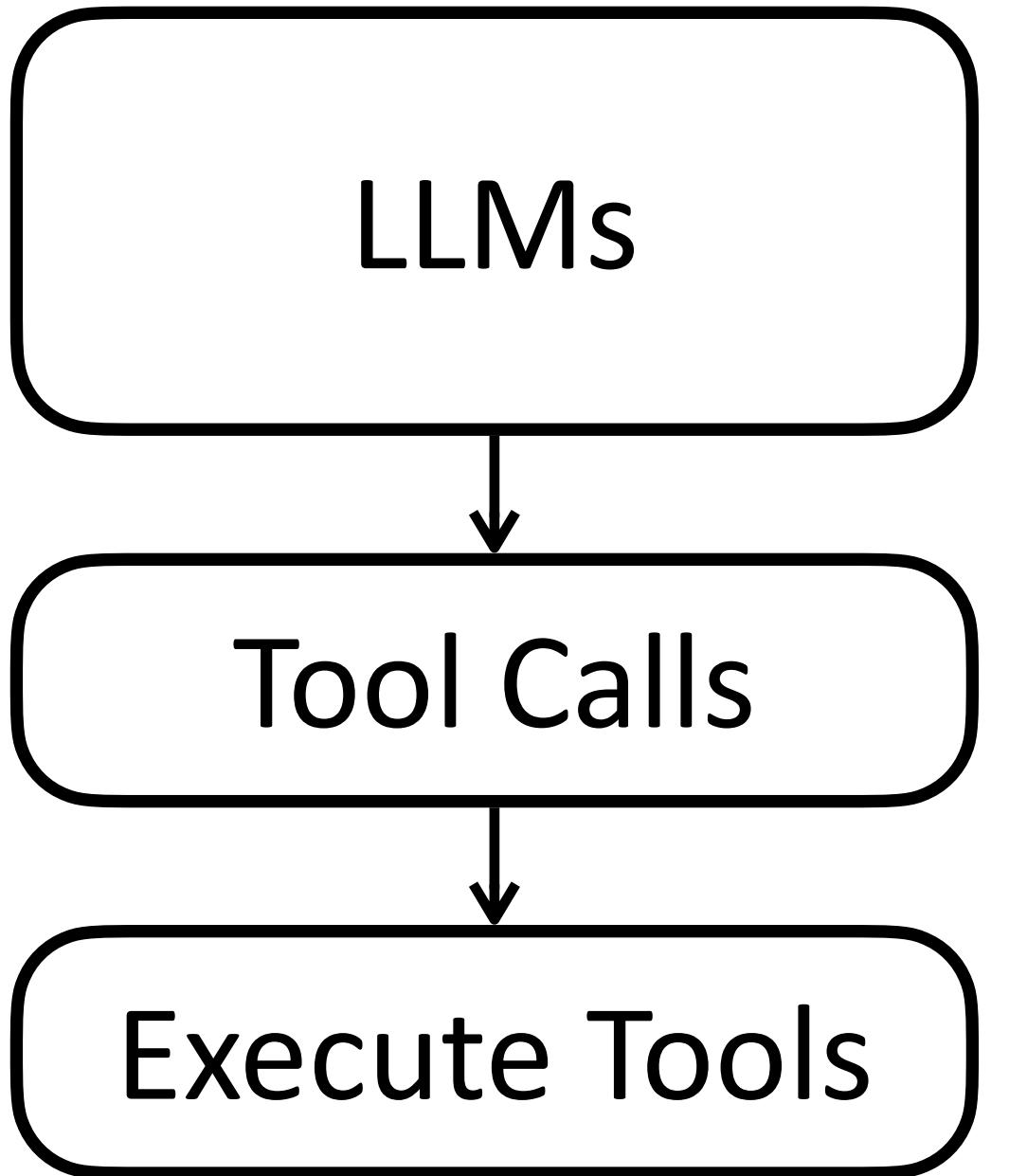


# One-Step Tool Call -> Agents

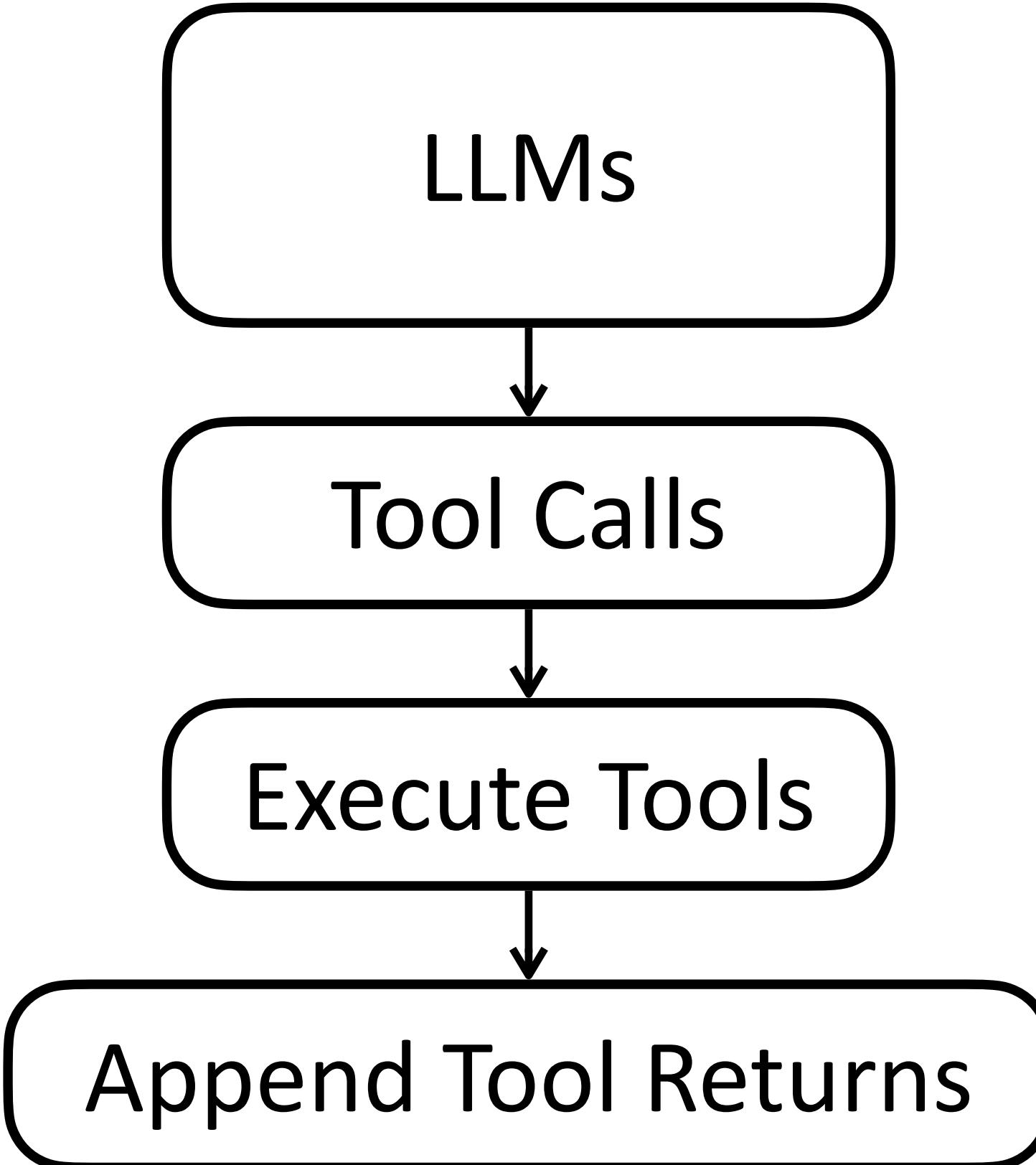


```
{  
  "response": "Sure, I'll check the current weather for you.",  
  "reasoning": "I need real-time conditions so the user's route  
recommendation is accurate.",  
  "tool_calls": [  
    {  
      "name": "get_weather",  
      "arguments": {  
        "location": "San Jose, CA, US",  
        "date": "2025-11-07"  
      }  
    }  
  ],  
}
```

# One-Step Tool Call -> Agents

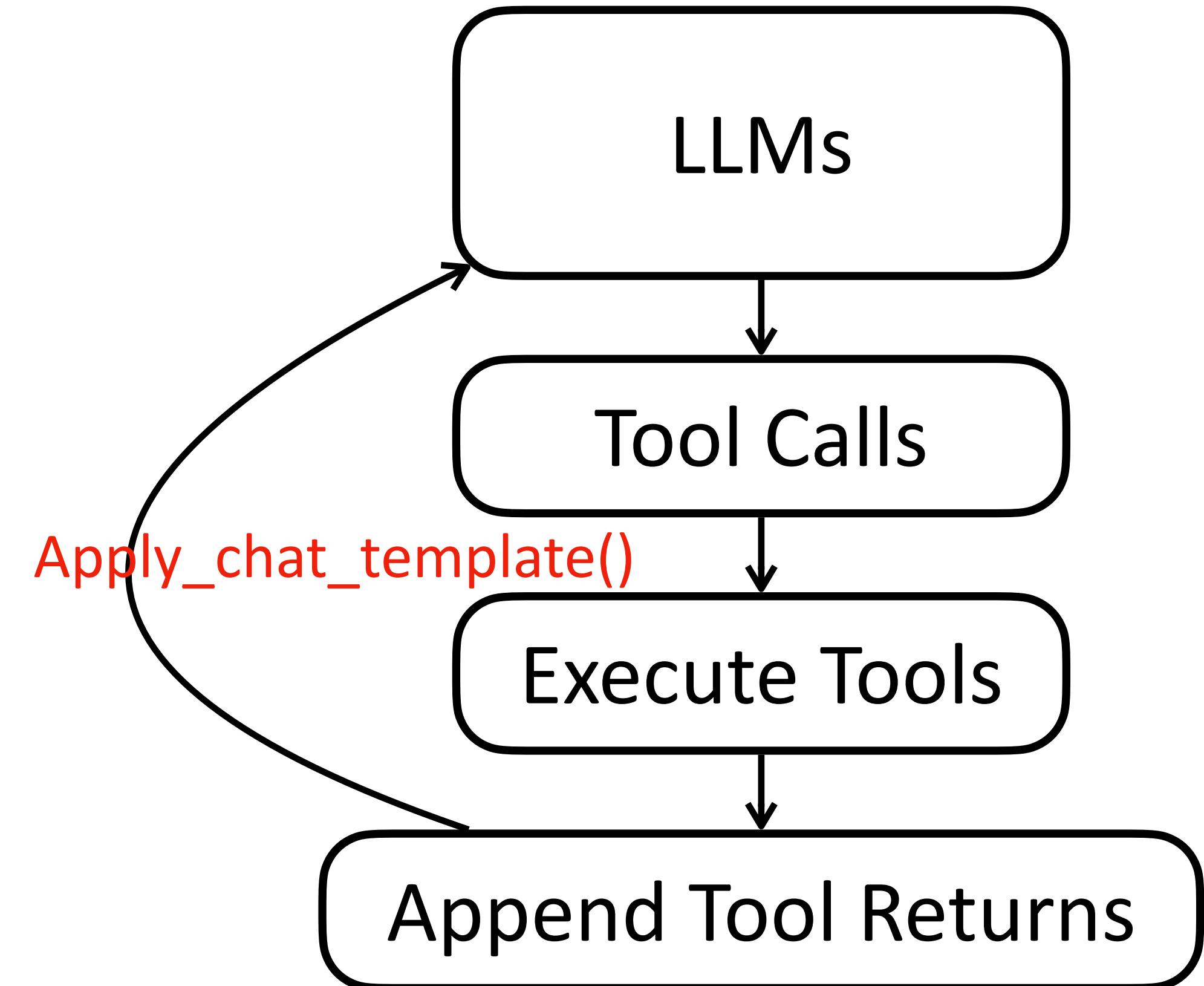


# One-Step Tool Call -> Agents



```
{  
    "response": "Sure, I'll check the current weather for you.",  
    "reasoning": "I need real-time conditions so the user's route recommendation is accurate.",  
    "tool_calls": [  
        {  
            "name": "get_weather",  
            "arguments": {  
                "location": "San Jose, CA, US",  
                "date": "2025-11-07"  
            }  
        }  
    ],  
    "tool_return": {  
        "temperature": 21.5,  
        "condition": "clear",  
        "humidity": 60,  
        "wind_speed": 10,  
        "location": "San Jose, CA, US",  
        "date": "2025-11-07"  
    }  
}
```

# One-Step Tool Call -> Agents



```
{  
    "response": "Sure, I'll check the current weather for you.",  
    "reasoning": "I need real-time conditions so the user's route recommendation is accurate.",  
    "tool_calls": [  
        {  
            "name": "get_weather",  
            "arguments": {  
                "location": "San Jose, CA, US",  
                "date": "2025-11-07"  
            }  
        }  
    ],  
    "tool_return": {  
        "temperature": 21.5,  
        "condition": "clear",  
        "humidity": 60,  
        "wind_speed": 10,  
        "location": "San Jose, CA, US",  
        "date": "2025-11-07"  
    }  
}
```

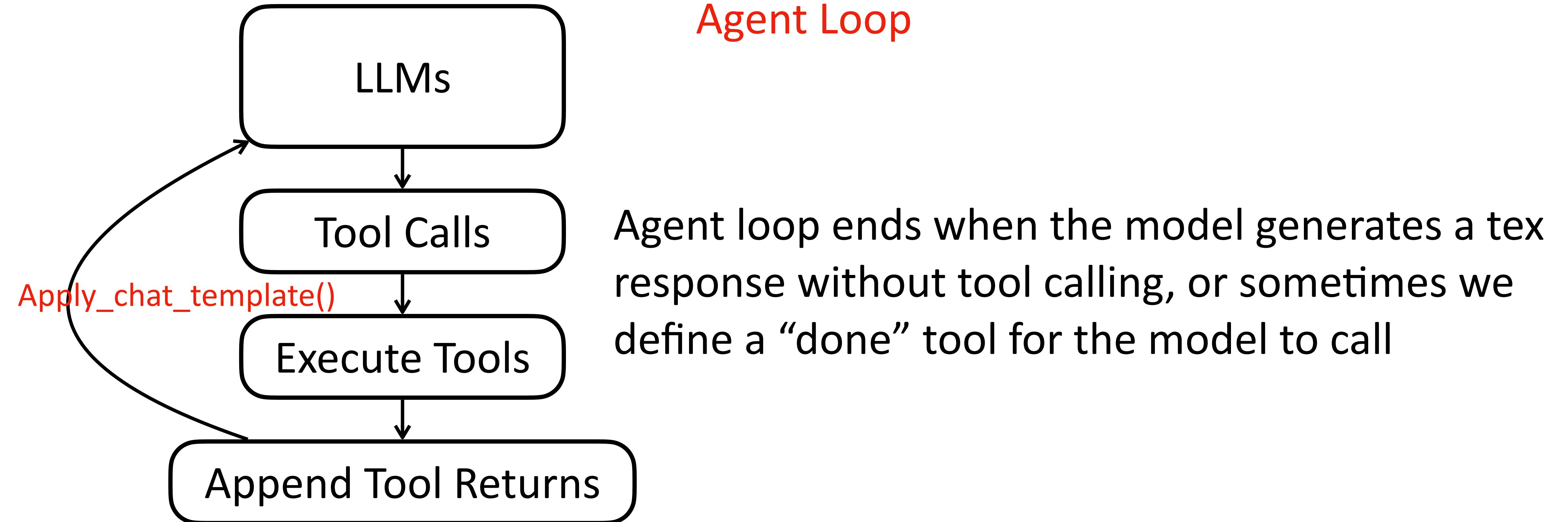
Sure, I'll check the current weather for you.

[thinking] I need real-time conditions so the user's route recommendation is accurate. [/thinking]

```
<tool_call>
{"name": "get_weather", "arguments": {"location": "San Jose, CA, US", "date": "2025-11-07"}}
</tool_call>
<tool_return>
{
    "temperature": 21.5,
    "condition": "clear",
    "humidity": 60,
    "wind_speed": 10,
    "location": "San Jose, CA, US",
    "date": "2025-11-07"
}
</tool_return>
```

This is the context fed back to the model to continue generation

# One-Step Tool Call -> Agents



```

def send_messages(messages):
    response = client.chat.completions.create(
        model="deepseek-chat",
        messages=messages,
        tools=tools
    )
    return response.choices[0].message

client = OpenAI(
    api_key=<your api key>,
    base_url="https://api.deepseek.com",
)
tools = [
    {
        "type": "function",
        "function": {
            "name": "get_weather",
            "description": "Get weather of a location, the user should supply a location first.",
            "parameters": {
                "type": "object",
                "properties": {
                    "location": {
                        "type": "string",
                        "description": "The city and state, e.g. San Francisco, CA",
                    }
                },
                "required": ["location"]
            },
        },
    },
]
messages = [{"role": "user", "content": "How's the weather in Hangzhou, Zhejiang?"}]
message = send_messages(messages)
print(f"User>\t {messages[0]['content']}")

tool = message.tool_calls[0]
messages.append(message)

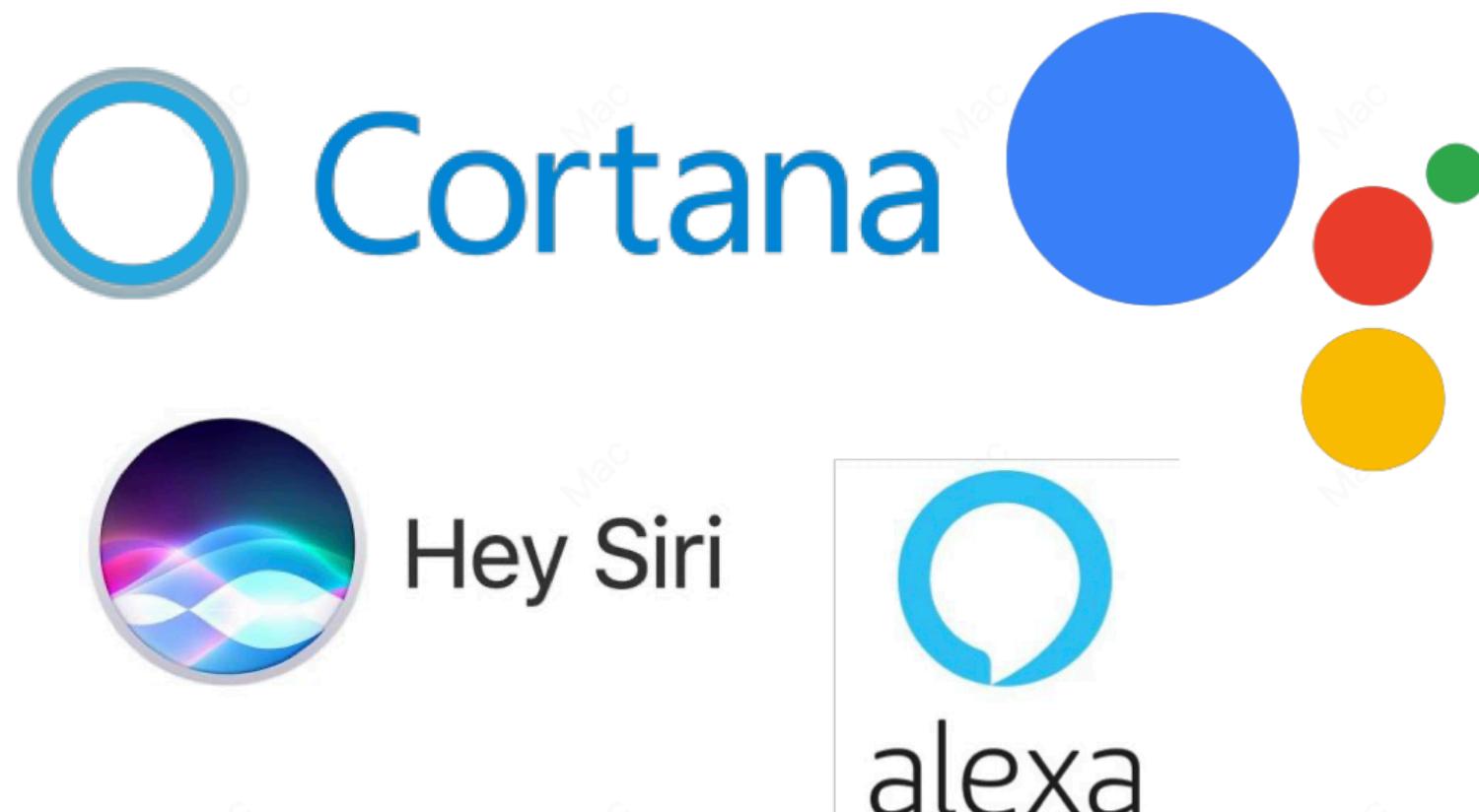
messages.append({"role": "tool", "tool_call_id": tool.id, "content": "24"})
message = send_messages(messages)
print(f"Model>\t {message.content}")

```

## One-step Example

# Why Do We Want Agents

Imagine if things get done by just talking...



```
Untitled-1
```

```
my_list = [3, 5, 1]
sort in descending order →
sorted(my_list, reverse=True)
```

The screenshot shows a Python code editor window titled "Untitled-1". The code defines a list `my\_list` with values [3, 5, 1], then uses the `sorted` function to sort it in descending order. The third line of code, `sorted(my\_list, reverse=True)`, is highlighted in green. At the bottom of the editor, it says "master\*" and "Python 3.6.5 64-bit".

## Virtual Assistants

- Set an alarm at 7 AM
- Remind me for the meeting at 5pm
- Play Jay Chou's latest album

## Natural Language Programming

- Sort my\_list in descending order
- Copy my\_file to home folder
- Dump my\_dict as a csv file output.csv

# Training-free Methods for Building Agents

# How to Let LLM Become an Agent

Sure, I'll check the current weather for you.

[thinking] I need real-time conditions so the user's route recommendation

```
<tool_call>
{"name": "get_weather", "arguments": {"location": "San Jose, CA, US", "dat
</tool_call>
```

The New England Journal of Medicine is a registered trademark of [QA("Who is the publisher of The New England Journal of Medicine?") → Massachusetts Medical Society] the MMS.

Out of 1400 participants, 400 (or [Calculator(400 / 1400) → 0.29] 29%) passed the test.

The name derives from "la tortuga", the Spanish word for [MT("tortuga") → turtle] turtle.

The Brown Act is California's law [WikiSearch("Brown Act") → The Ralph M. Brown Act is an act of the California State Legislature that guarantees the public's right to attend and participate in meetings of local legislative bodies.] that requires legislative bodies, like city councils, to hold their meetings open to the public.

We just need the LLMs to output certain formats of tool calls that we can parse

# Just Prompt

Supposing you have the following APIs, answer the given question.

Available APIs:

1. weather(location): for getting the weather information given a location
2. location(): for getting the location of the user.
3. bus\_route(start, end): for finding the current bus route from the start position to the end position.
4. count\_characters(s): for counting the number of characters in a string s.

Question: Is it okay to go hiking today?

To know the answer, we first need to know the current location of the user.

API call: location() = Seattle.

Next, we need to know the current weather in Seattle.

API call: weather(Seattle) = Cloudy with a chance of rain.

Based on this information, it is not recommended to go hiking today in Seattle.

We just need the LLMs to output certain formats of tool calls that we can parse

Available APIs:

1. weather(location): for getting the weather information given a location
2. location(): for getting the location of the user.
3. bus\_route(start, end): for finding the current bus route from the start position to the end position.
4. count\_characters(s): for counting the number of characters in a string s.

Question: how long does it take to take a bus to Vancouver?

To answer this question, we first need to know the current location of the user.

API call: location() = Seattle.

Next, we need to find the bus route from Seattle to Vancouver.

API call: bus\_route(Seattle, Vancouver) = 4 hours.

Therefore, it takes 4 hours to take a bus from Seattle to Vancouver.

# Evaluating Language Agents

# Evaluation of LLM Agents

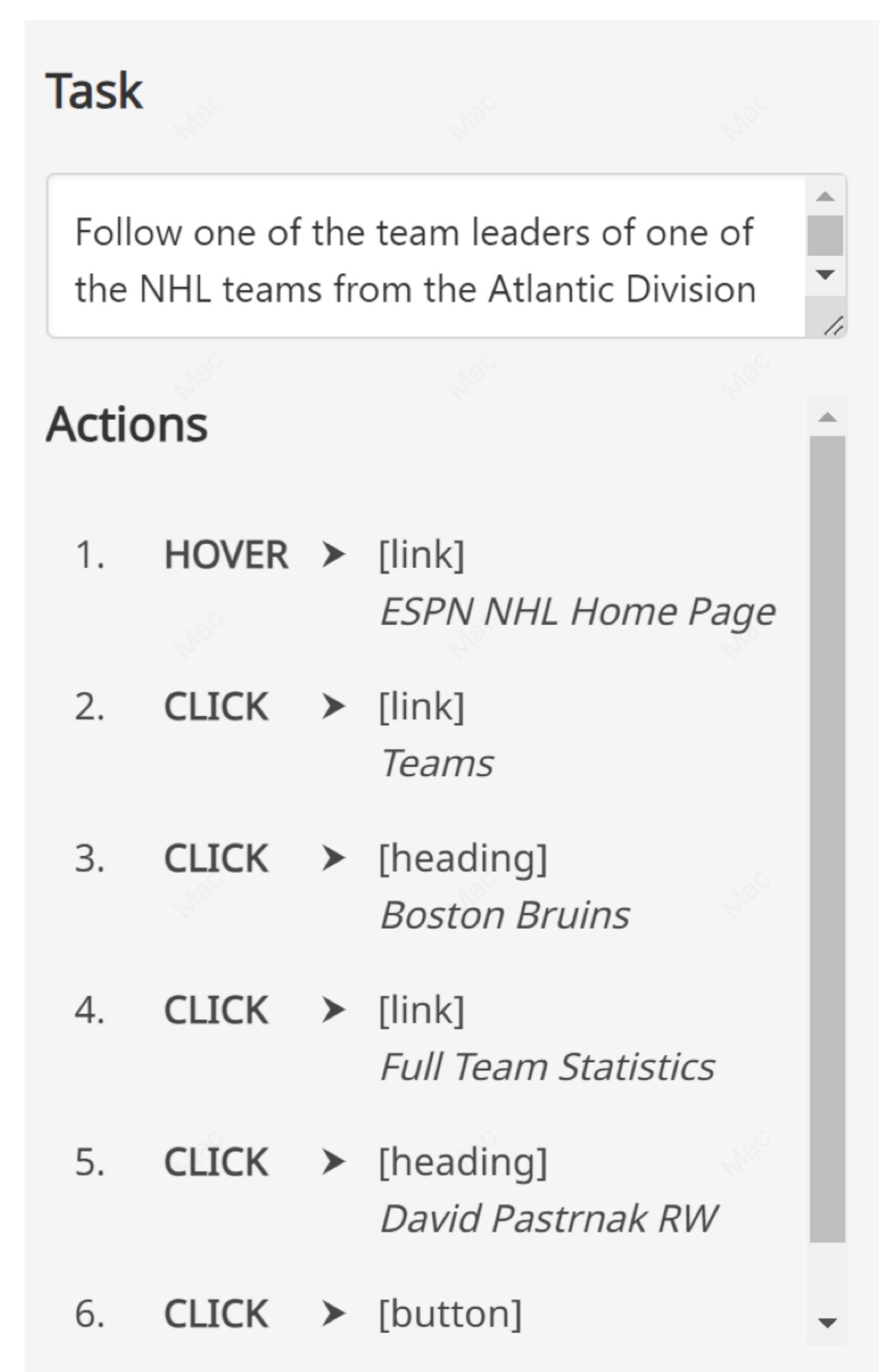
- Simplified environments and basic tasks
- Performance is saturating.

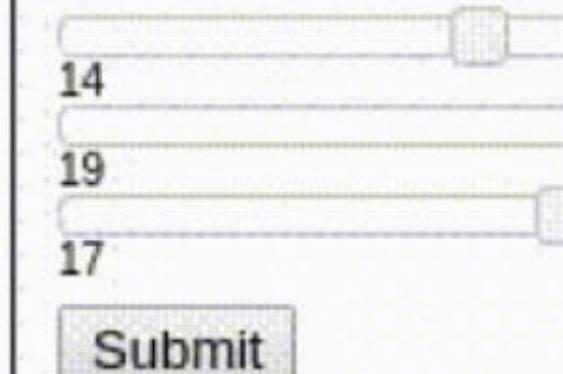
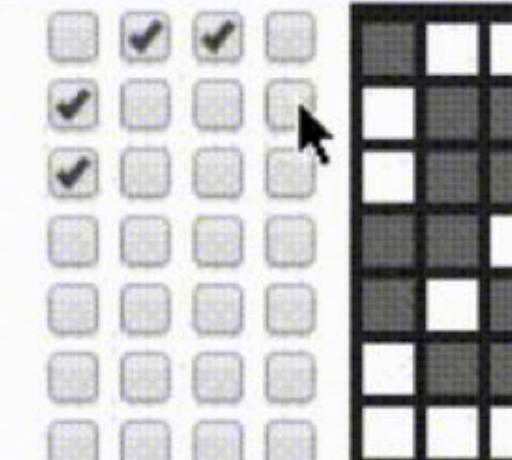
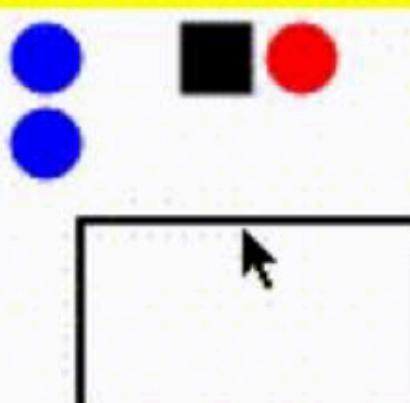
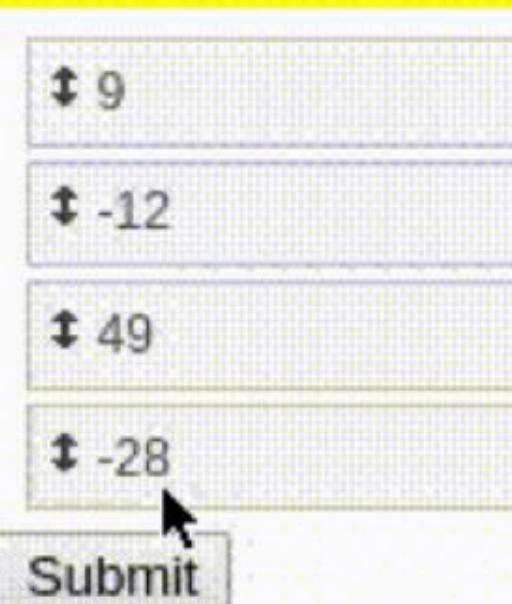
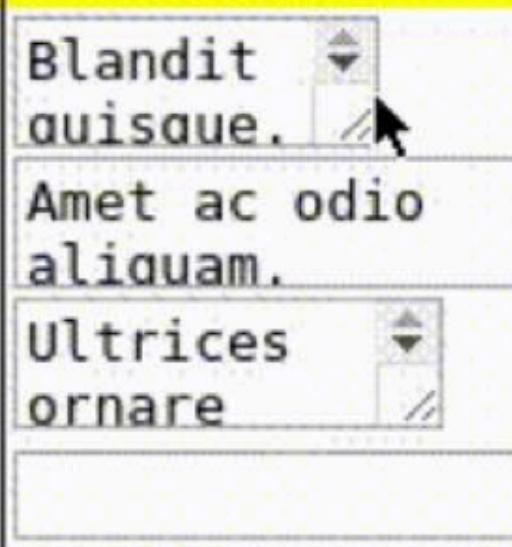
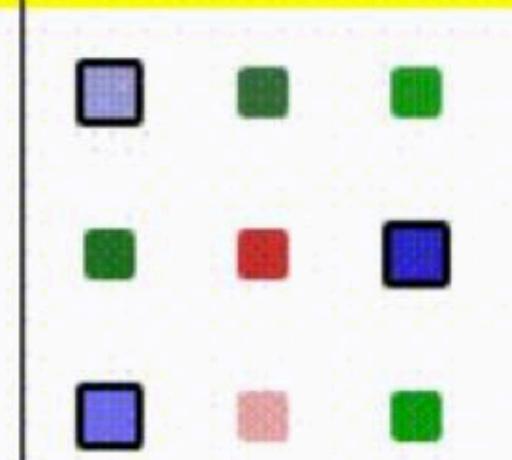
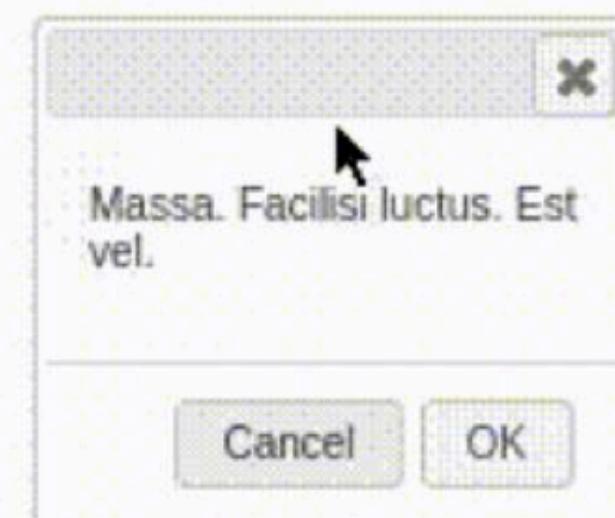
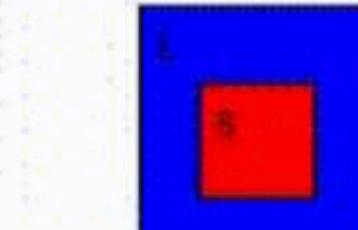
1. Stateless, non interactive environment, e.g.

*Mind2Web (Deng et al. 2023)* has only dumped pages.

2. Checking action sequence accuracy (step-wise, surface form only)

3. Simple interactive environment, short horizon, e.g.  
*WebShop (Yao et al. 2023), MiniWoB++ (Humphreys et al. 2022)*



<p>Move the cube around so that "5" is the active side facing the user.</p>  <p><b>Submit</b></p>	<p>Set the sliders to the combination [13,20,13] and submit.</p>  <p><b>Submit</b></p>	<p>Draw the number "2" in the checkboxes using the example on the right and press Submit when finished.</p>  <p><b>Submit</b></p>	<p>Drag Ree to the 4th position.</p> 	<p>Keep your mouse inside the circle as it moves around.</p> 	<p>Enter the value of Country into the text field and press Submit.</p> <table border="1"> <tr> <td>Gender</td> <td>Male</td> </tr> <tr> <td>First name</td> <td>Annecorinne</td> </tr> <tr> <td>Country</td> <td>Guam</td> </tr> <tr> <td>Year of Birth</td> <td>1934</td> </tr> <tr> <td>Religion</td> <td>Hinduism</td> </tr> </table> <p><b>Submit</b></p>	Gender	Male	First name	Annecorinne	Country	Guam	Year of Birth	1934	Religion	Hinduism
Gender	Male														
First name	Annecorinne														
Country	Guam														
Year of Birth	1934														
Religion	Hinduism														
<p>Drag all triangles into the black box.</p>  <p><b>Submit</b></p>	<p>Select 09/23/2016 as the date and hit submit.</p>  <p><b>Submit</b></p>	<p>Sort the numbers in increasing order, starting with the lowest number at the top of the list.</p>  <p><b>Submit</b></p>	<p>Copy the text from the 1st text area below and paste it into the text input.</p>  <p><b>Submit</b></p>	<p>Select all the shades of blue and press Submit.</p>  <p><b>Submit</b></p>	<p>Find the 4th word in the paragraph, type that into the textbox and press "Submit".</p> <p>Non arcu ut <b>ultricies</b> est. Gravida gravida. Porta erat nulla eget condimentum posuere a</p> <p><b>Submit</b></p>										
<p>Click the button in the dialog box labeled "Cancel".</p>  <p><b>Submit</b></p>	<p>Highlight the text in the paragraph below and click submit.</p> <p>Ultricies. Sagittis in.</p> <p><b>Submit</b></p>	<p>Highlight the text in the paragraph below and click submit.</p> <p>Tempor posuere nibh. Vel nisl, faucibus. Feugiat condimentum</p> <p><b>Submit</b></p>	<p>Find the 11th word in the paragraph, type that into the textbox and press "Submit".</p> <p>Ullamcorper aliquet amet ullamcorper. Elit. Mattis luctus diam. Lobortis nulla fermentum ornare faucibus</p> <p><b>Submit</b></p>	<p>Move the cube around so that "2" is the active side facing the user.</p>  <p><b>Submit</b></p>	<p>Drag the smaller box so that it is completely inside the larger box.</p>  <p><b>Submit</b></p>										

Instruction: i am looking for x-large, red color women faux fur lined winter warm jacket coat, and price lower than 70.00 dollars

Current Query: women fur jacket coat

## Results

Page 1 (1-10) of 50 total results

[Back to Search](#)

[Next >](#)



[B09KP78G37](#)

Women Faux Fur Lined Jacket Coat  
Winter Warm Thick Fleece Outwear  
Trench Zipper Plus Size Long  
Sleeve Plush Overcoat



[B07ZXBGDXF](#)

Women's Coat, FORUU Winter Faux  
Fur Fleece Outwear Warm Lapel  
Biker Motor Aviator Jacket



[B098XT346Y](#)

Fjackets Real Lambskin Sherpa  
Jacket - Mens Leather Jacket

 4.7

**Current Action: click [Fjackets Real Lambskin...]**

# Key to Agent Benchmarks

## Environment:

- Diverse functionality.
- Rich and realistic content.
- Interactive
- Easily Extendable
- Reproducible

## Tasks:

- Long horizon tasks
- Enough difficulty
- Involves multiple websites

## Evaluation:

- Reliable metrics
- Encourage final goal rather than partial satisfaction.

# WebArena

## Environment:

- Diverse functionality.
- Rich and realistic content.
- Interactive
- Easily Extendable
- Reproducible

## A sandbox Internet:

- Open source, production-ready implementation of the websites
- Data populated from real-world websites
- Easily distributable – Dockers, AWS images, etc.

Self-hosted fully functional web applications

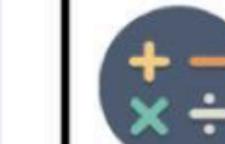


OneStopShop

CMS

reddit

GitLab



Toolbox



Knowledge resources

WebArena

# Example Tasks in WebArena



“Create a plan to visit Pittsburgh’s art museums with minimal driving distance starting from Schenley Park. Log the order in my “awesome-northeast-us-travel” repository”

webarena.wikipedia.com

List of museums in Pittsburgh

This list of museums in Pittsburgh, Pennsylvania encompasses museums defined for this context as institutions (including nonprofit organizations, government entities, and private businesses) that collect and care for objects of cultural, artistic, scientific, or historical interest and make their collections or related exhibits available for public viewing. Also included are university and non-profit art galleries. Museums that exist only in cyberspace (i.e., virtual museums) are not included.

Wikimedia Commons has media related to [Museums in Pittsburgh](#).

See also: [List of museums in Pennsylvania](#)

Museums

Search for museums in Pittsburgh

webarena.openstreetmap.com

OpenStreetMap Edit History Export

Schenley Park, Pittsburgh, Allegheny County

The Andy Warhol Museum, 117, Sandusky St

Car (OSRM)

Reverse Directions

Directions

Distance: 7.1km. Time: 0:10.

1. Start on Panther Hollow Road 300m  
2. Slight right onto unnamed road 160m

webarena.gitlab.com

Update README.md

README.md 158 B

Edit Replace

Travel in Northeast US

Pittsburgh

- + Miller Gallery at Carnegie Mellon University
- + American Jewish Museum
- + Carnegie Museum of Art

Record the optimized results to the repo

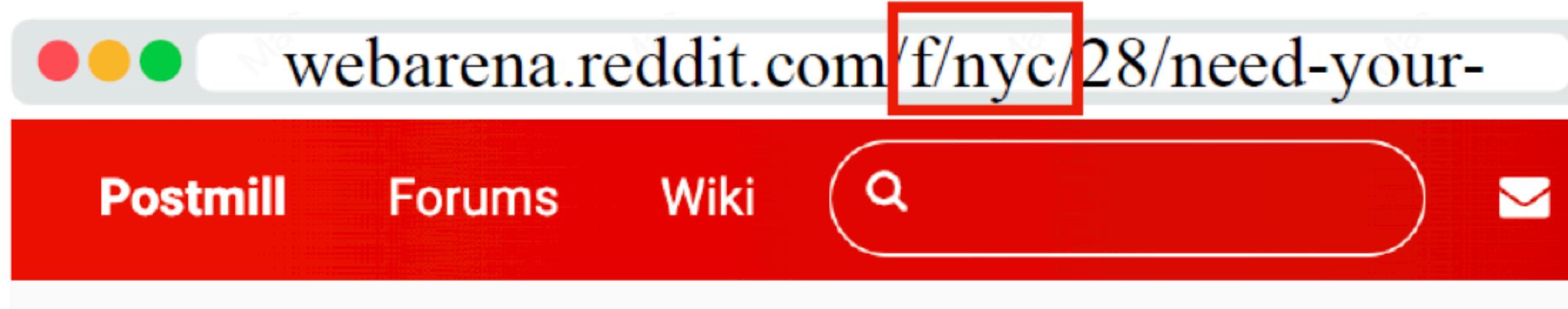
# Outcome/Execution-based Evaluation

Goal: directly validate the correctness of the execution

- “When was the last time I bought shampoo?”
- **Directly compare with the annotated answer:** Answer is “Dec 15th, 2022”

# Outcome/Execution-based Evaluation

*Post my question, “is car necessary in NYC”, in a subreddit where I’m likely to get an answer*



← /f/nyc

▲ Need your answer

1

Submitted by convexitgtaaaa 0 seconds ago in nyc

is car necessary in NYC?

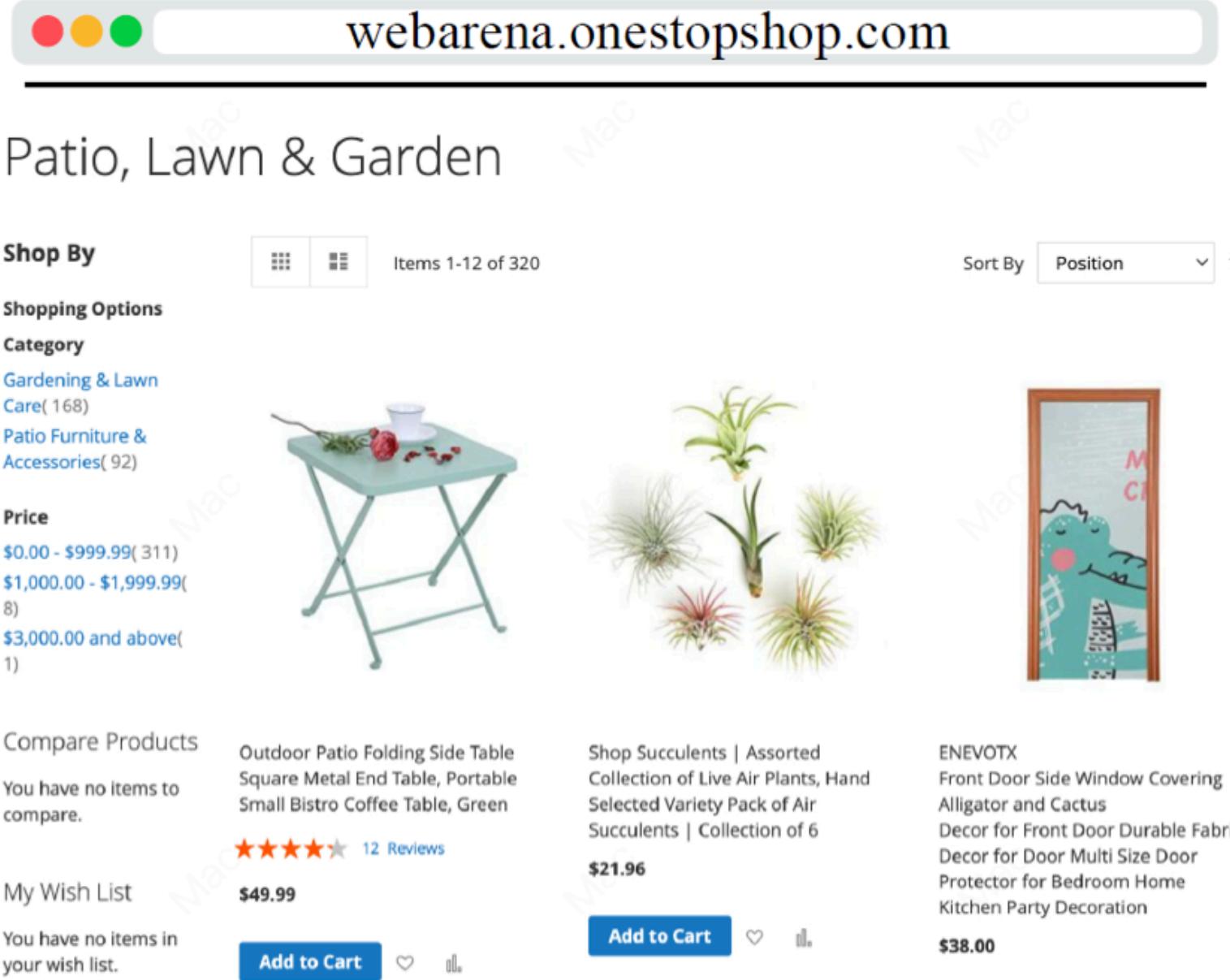
No comments Edit Delete

```
▼<div class="submission__row"> flex
  ▼<div class="submission__inner"> == $0
    ►<header class="submission__header">...</header>
    ▼<div class="submission__content flow-slim">
      ▼<div class="submission__body break-text text-flow">
        <p lang="en" dir="ltr">is car necessary in NYC?</p>
      </div>
    ...
```

“f/nyc” in page.url

“Is car necessary in NYC?” in  
document.querySelector(“.su  
bmission\_\_inner”).outText

# Observation & Action Space



Screenshot

Keyboard: type

```
<li>
<div>
<a href="..."></a>
<div class>
<a href="...">>Outdoor Patio ...
</a>
<div>
<span>Rating:</span>
<div>
<span>82%</span>
</div>
<a href="#reviews">12
<span>Reviews</span></a>
```

Text

Mouse: click, hover, scroll

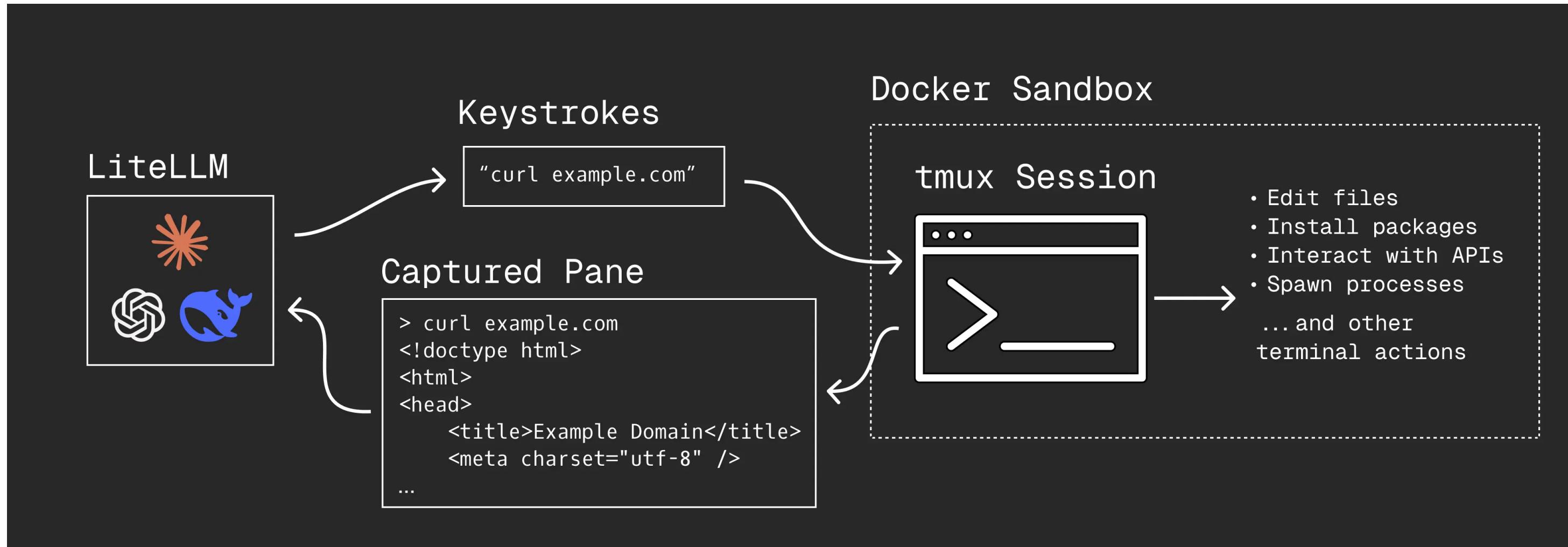
```
RootWebArea 'Patio, Lawn ..'
link 'Image'
img 'Image'
link 'Outdoor Patio..'
LayoutTable "
  StaticText 'Rating:'
    generic '82%'
  link '12 Reviews'
  StaticText '$49.99'
  button 'Add to Cart' focusable: True
  button 'Wish List' focusable: ...
  button 'Compare' focusable: ...
```

Accessibility tree

Browser: New tab, go back

Another type of web agents, GUI agents, directly takes image as input observations

# TerminalBench



## Terminal-Bench: Possible Agent Actions

- Run shell commands (`ls`, `cd`, `make`, `python`, etc.)
- Manage tmux sessions/panes (`new-session`, `split-window`, `select-pane`)
- Edit files (`vim`, `nano`, `echo > file`, etc.)
- Install/build software (`apt install`, `gcc`, `make`)
- Read & analyze outputs/logs (`cat`, `less`, `grep`)
- Navigate directories and view help (`cd`, `ls`, `--help`)
- Verify or fix results (re-run scripts, check outputs)

# Training Methods for Improving Agents

# Learning of LLM Agents

- In-Context Learning – Learning from few-shot exemplars
- Supervised Finetuning – Learning From *Experts*
- Reinforcement Learning – Learning from *Environment*

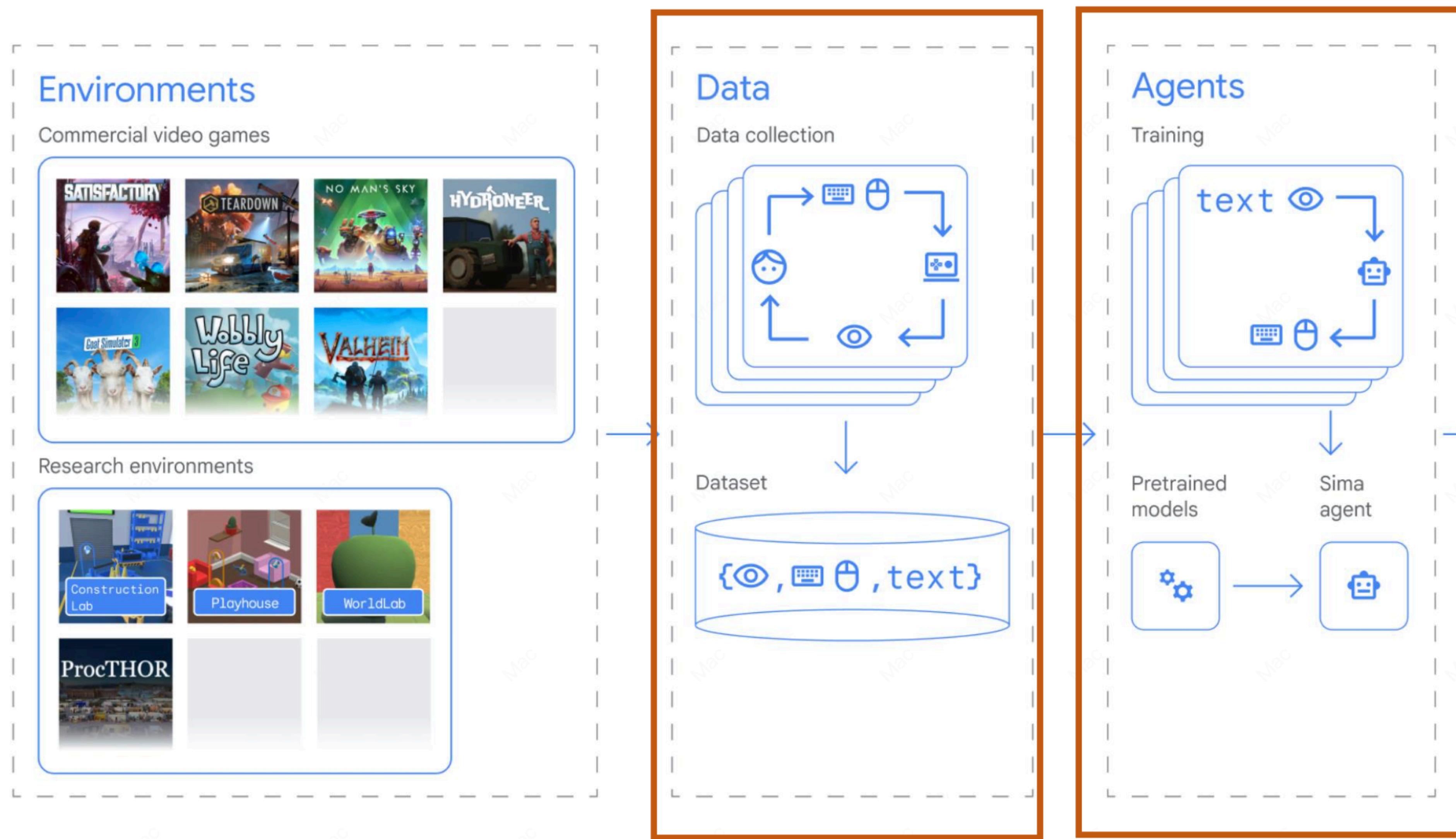
# Supervised Finetuning

- Collect large amount of expert trajectories (e.g. from human annotation)

task\_intent, [(obs\_1, action\_1), ..., (obs\_N, action\_N)]

- Finetune the LLM with standard cross-entropy loss.

# Supervised Finetuning



# Supervised Finetuning

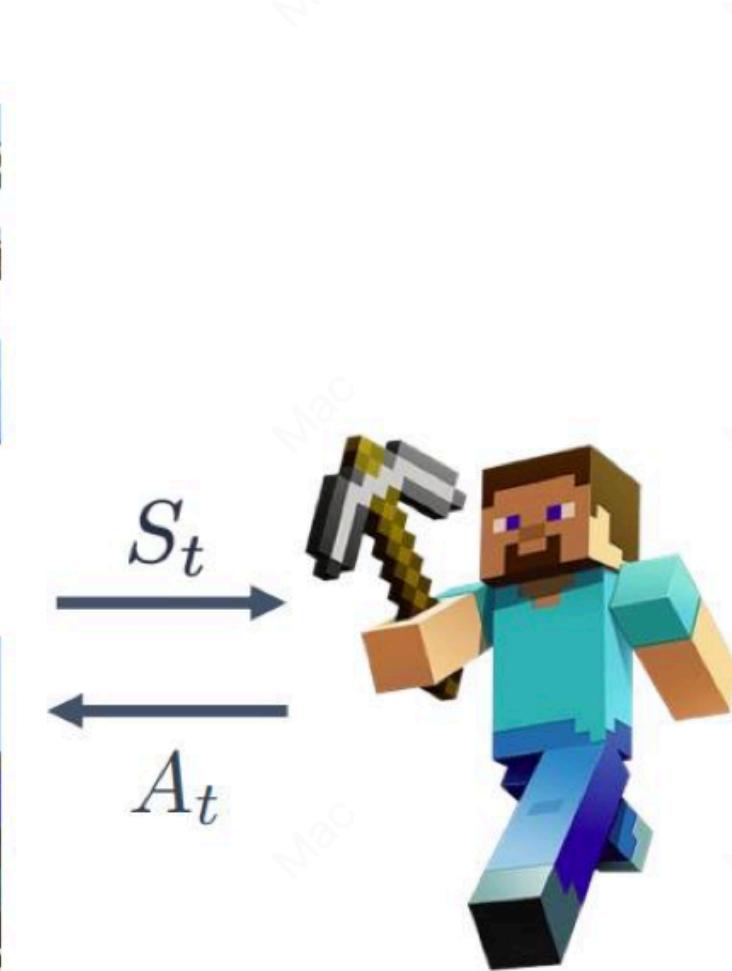
- Data hungry
- Cannot learn much from failed trajectories
  - $a_1, a_2, a_3, \dots, a_{10}$  - Success
  - $a_1, a_2, a_3, \dots, a_{10}$  - Fail (Wasted)
- Need human trajectory?
  - Data augmentation techniques

# Create More Training Data

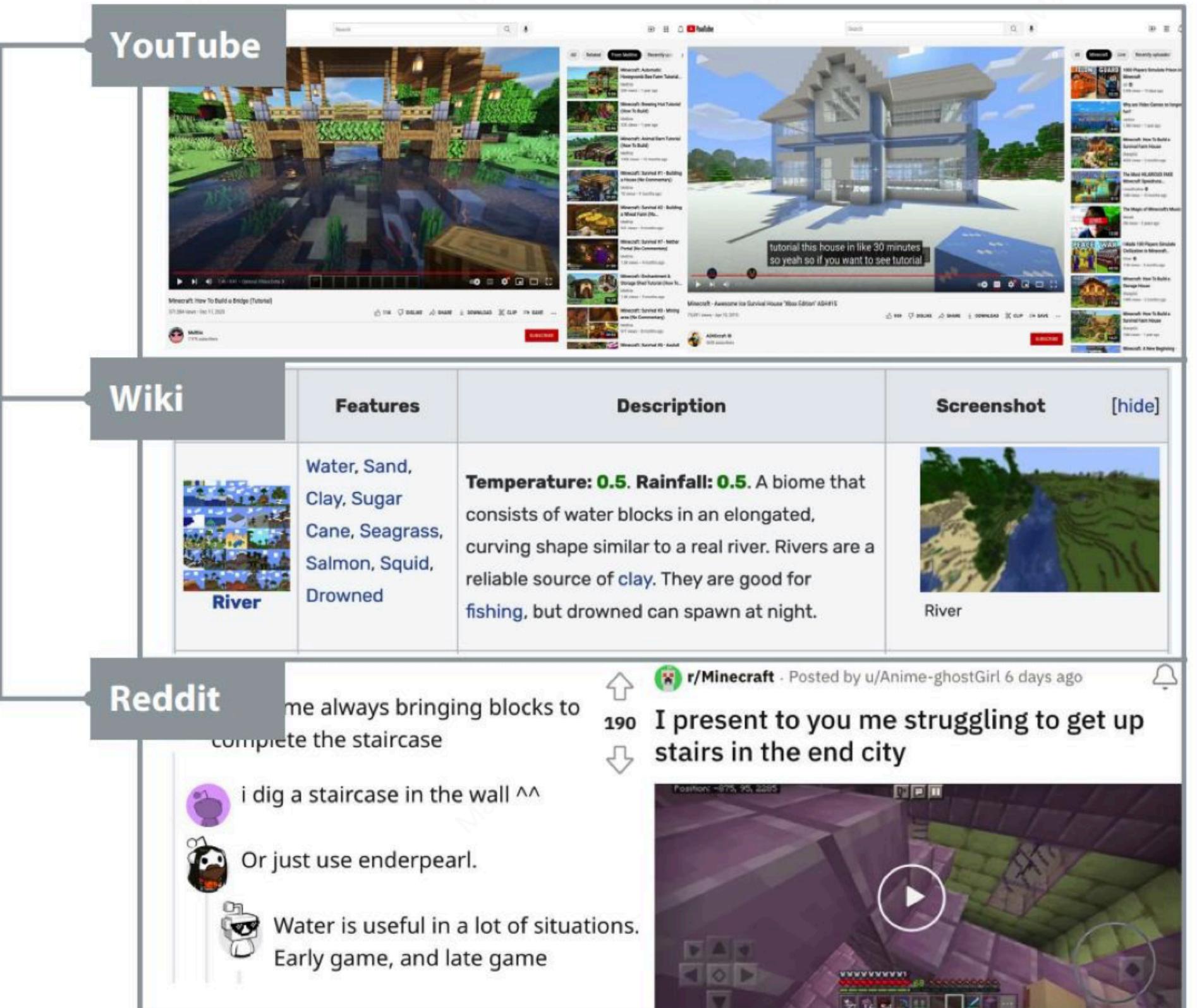
## Open-ended Environments



## Generalist Agent

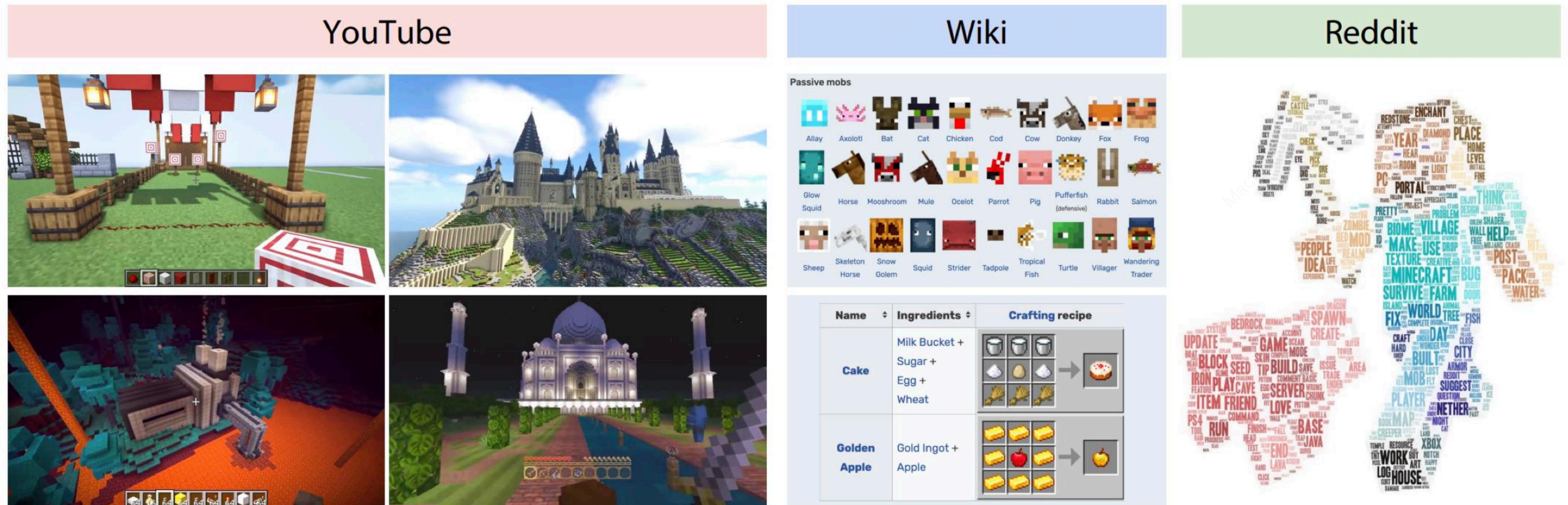


## Internet-scale Knowledge Base



# Data Augmentation

- Continue pre-train on large amount of data automatically mined
  - Even noisy, not clear trajectories, provide domain adaptation.

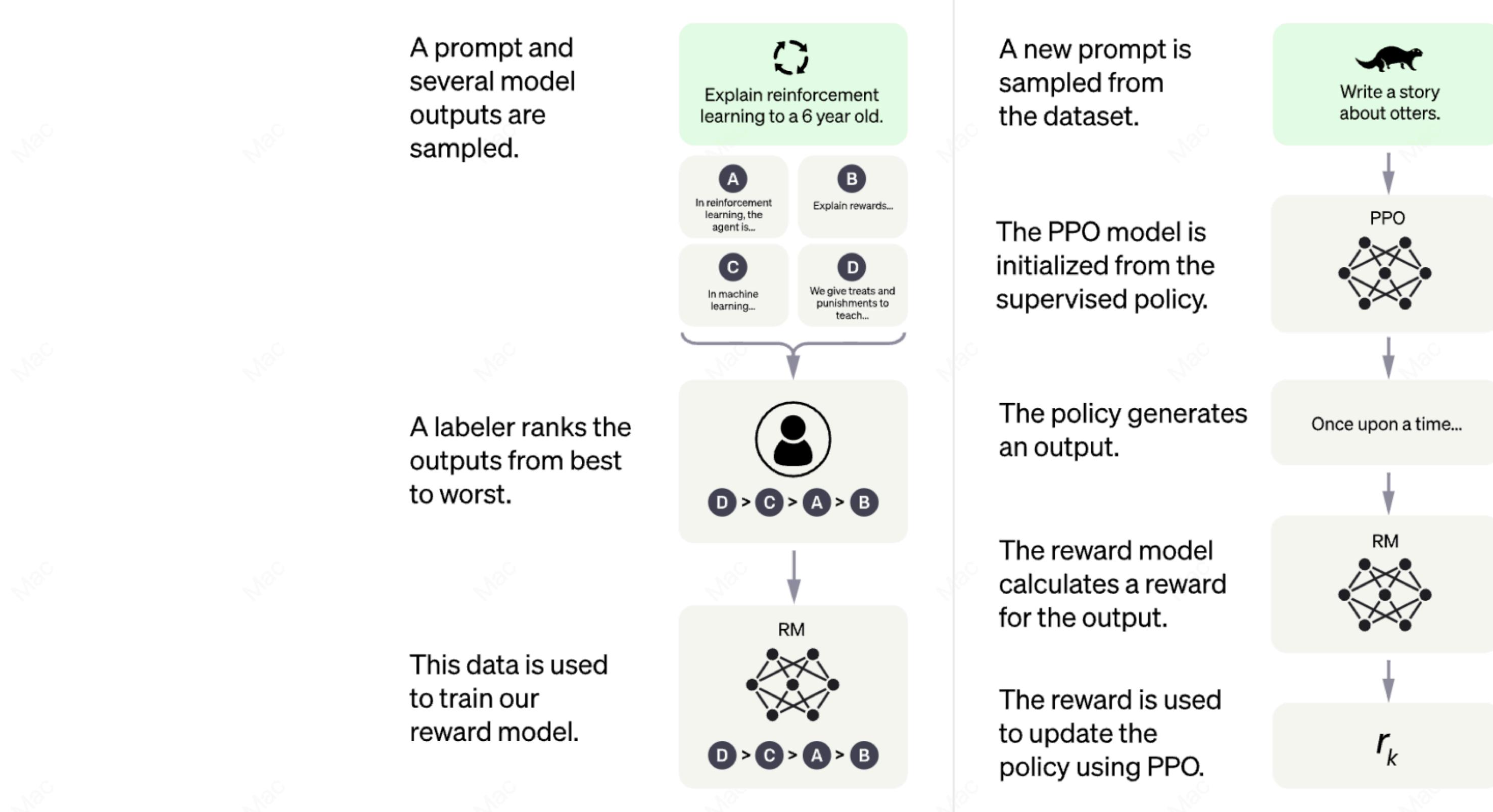


MineDojo, Fan et al. 22'  
Don't Stop Pretraining, Gururangan et al., 20'

# Reinforcement Learning

Lots of on-going research in this area!

Recall RLHF: Reinforcement Learning from Human Feedback:

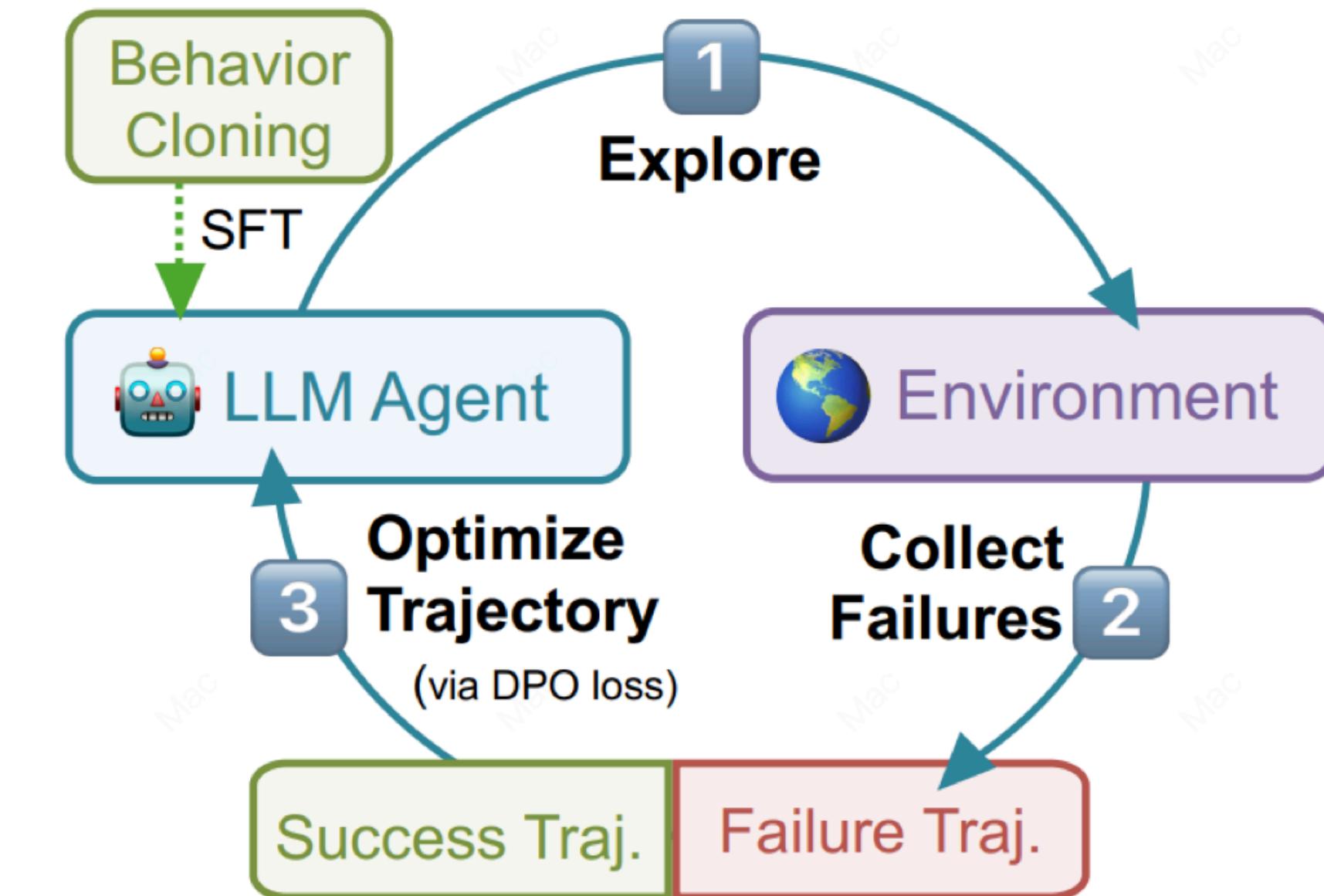
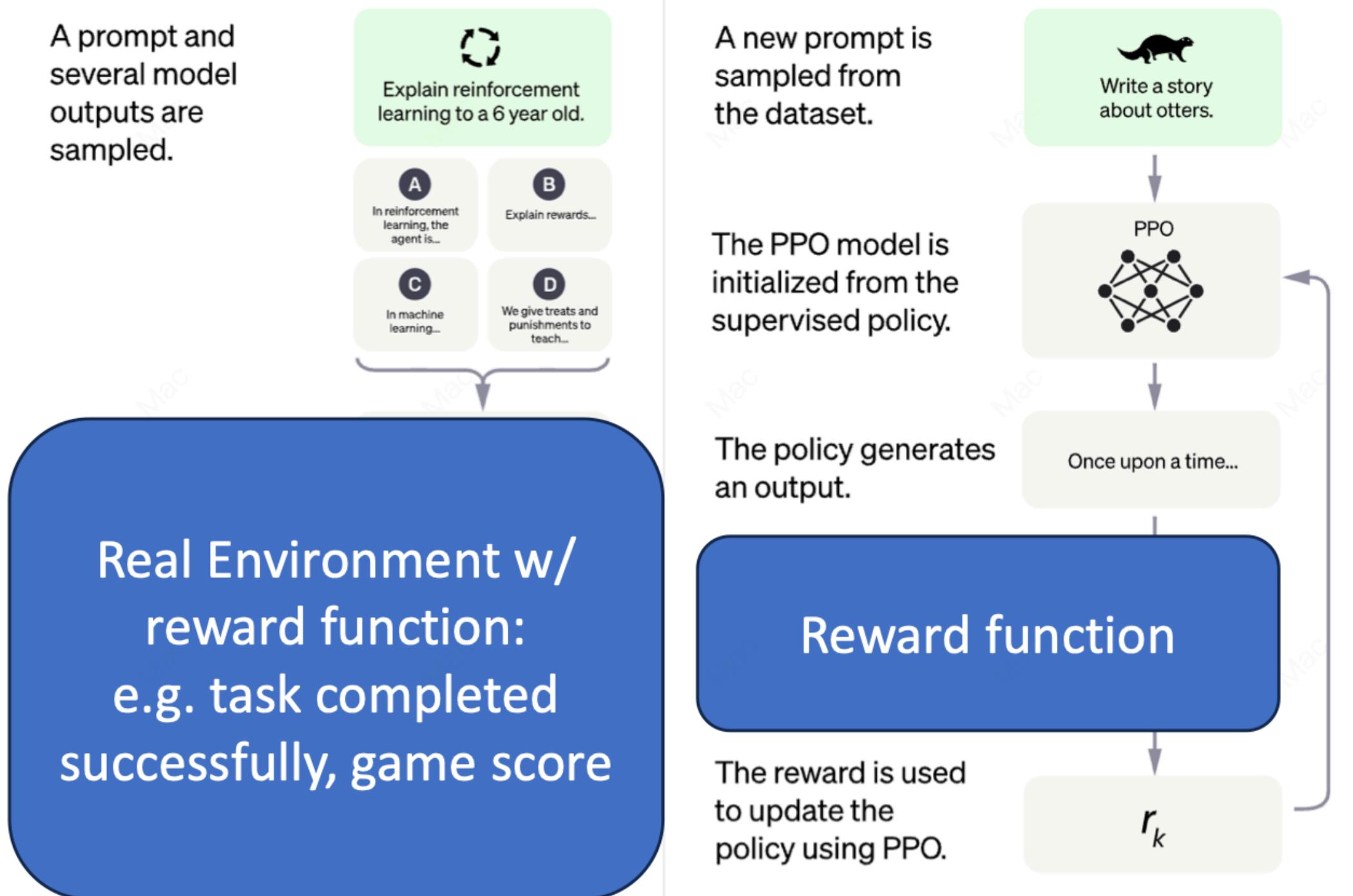


RLHF, Ouyang, et al. 22'

# Reinforcement Learning

Compared to RLHF:

Given environment, *reward function*  
(trajectory, reward) pairs without human



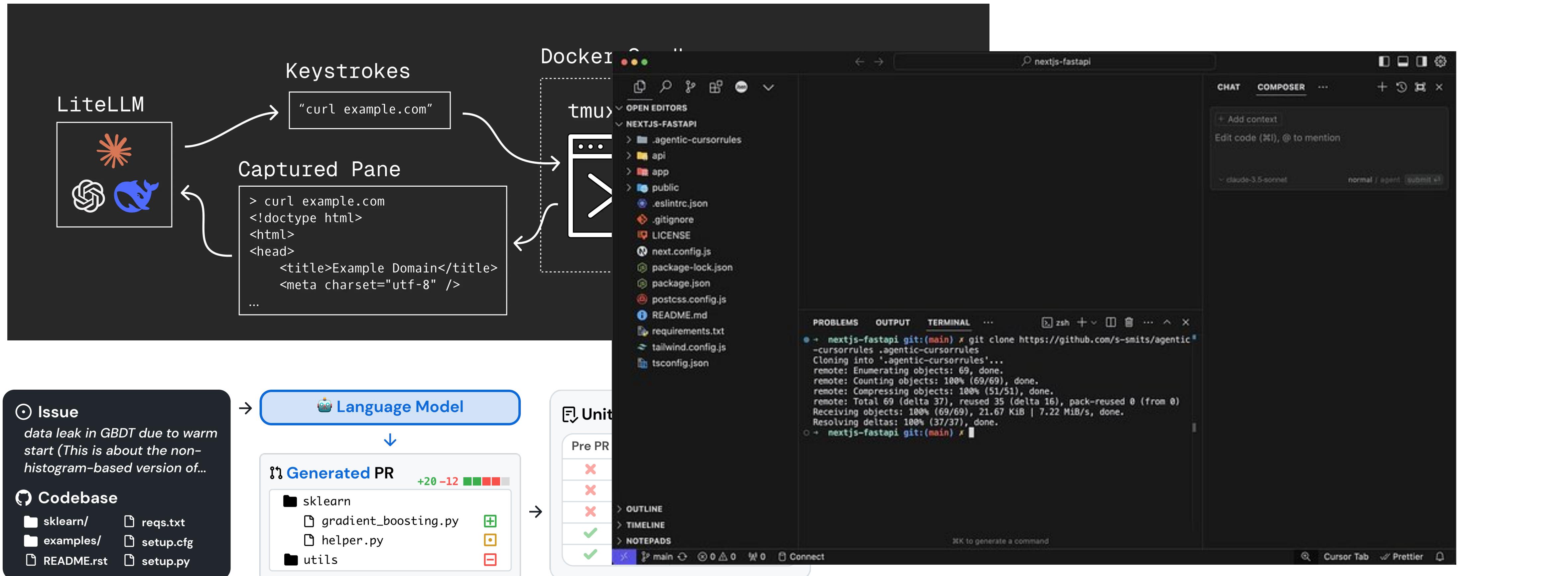
Trial and Error (Song et al. 24')

# Reinforcement Learning

- Closed loop, interactive environment
- Need good reward functions
  - What if the task success/fail is not easy to automatically assess?
- Need good initial models
  - Has decent basic knowledge ability, sparse rewards
- Scalability
  - The environment takes 10 seconds to env.step()
  - The reward function takes 100 seconds to get a scalar reward

# RL Environments

Environments and benchmarks typically come together



Research and Products are really close nowadays, and we can directly RL in real, product-level environments

**Thank You!**