

COMP7607 Final Project:

Foundation Agents for General Digital Control

Zhiyong Wu

2023/09/09



1. Final Project Announcement
2. Background & Research Sharing
3. Exemplary Projects



Self-Introduction: Zhiyong Wu

Education

- 2017-2021, PhD, University of Hong Kong
- 2013-2017, BEng, Wuhan University

Experience

- 2021.12- now, Shanghai AI Lab, Research Scientist

Research Interest

- Large Language Models
- Language Agent

Jarvis From the Iron Man



COMP7607 Final Project: Building Your Own JARVIS

Teams:

5-8 students

Topics:

Create a fully functional OS agent (JARVIS-like digital copilot)



JARVIS: Render complete in 5 minutes.

COMP7607: Building Jarvis-like Digital Copilot

Embodied AI,

but in the virtual world!

Value:

Improve the productivity of the
whole society

Enhance technology accessibility
for the elderly or individuals with lower
education levels.



JARVIS: Render complete in 5 minutes.

WHY

WHY: An 'underestimated' track where the research attention (**red curve**) is not proportional to the capital investment.

=» The opportunity to become a field leader ↑

=> Lots of job opp in industry

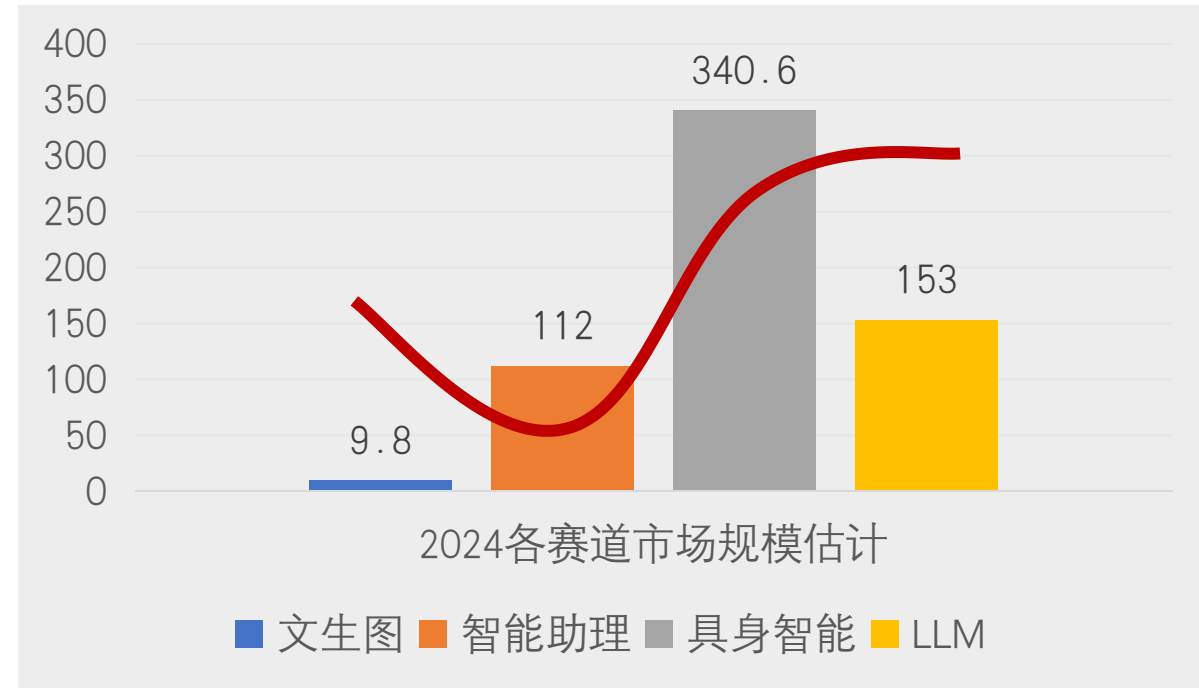
Text-to-Image: Stable Diffusion, Dall-E

Embodied AI: Stanford、Berkeley

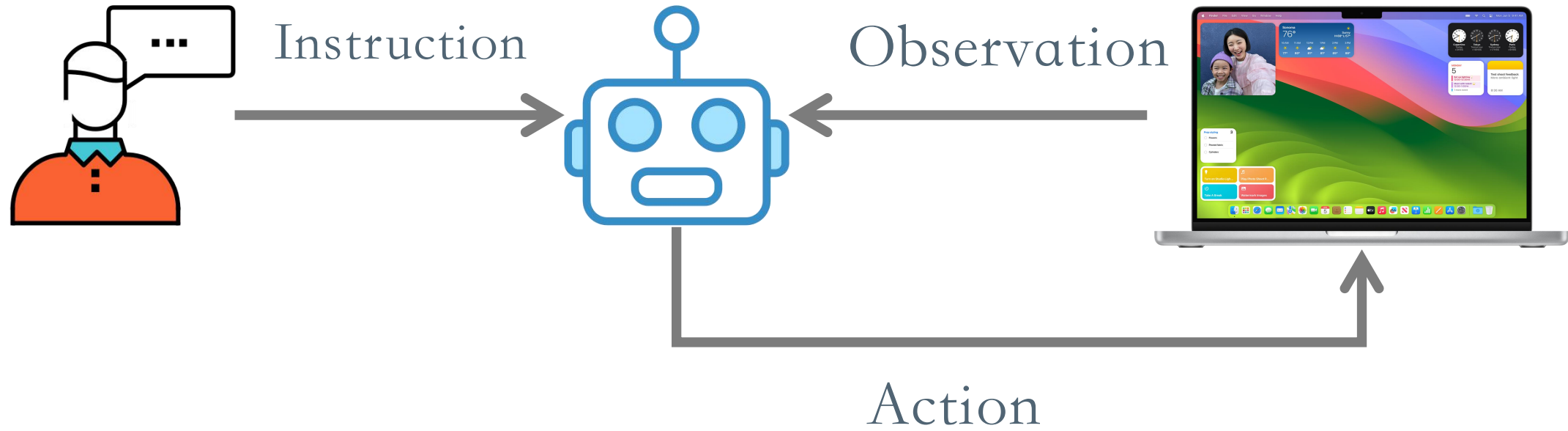
LLM: OpenAI, Anthropic

Digital Assistant: ?

Siri, Cortana, Alexa



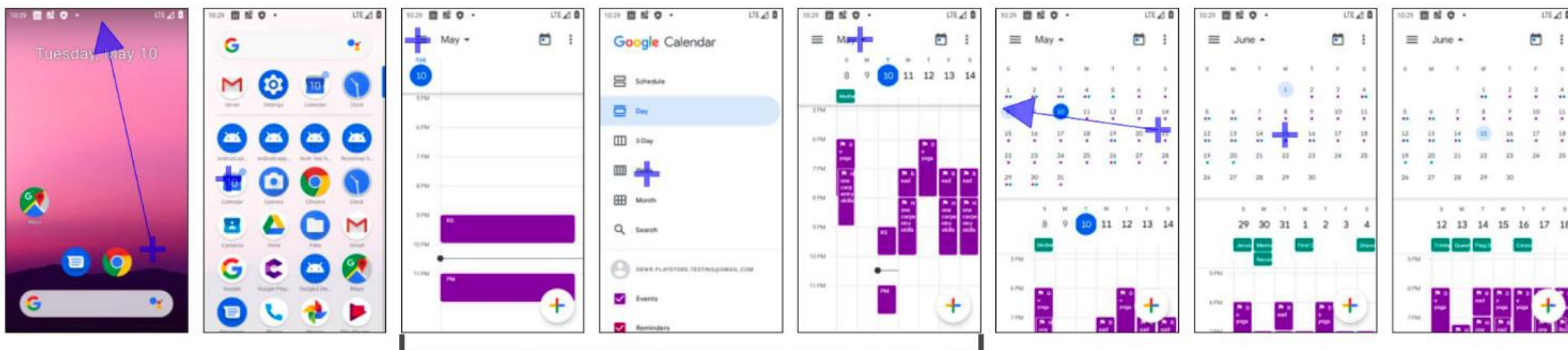
OS Agent System Overview



Challenge 1: Data Scarcity

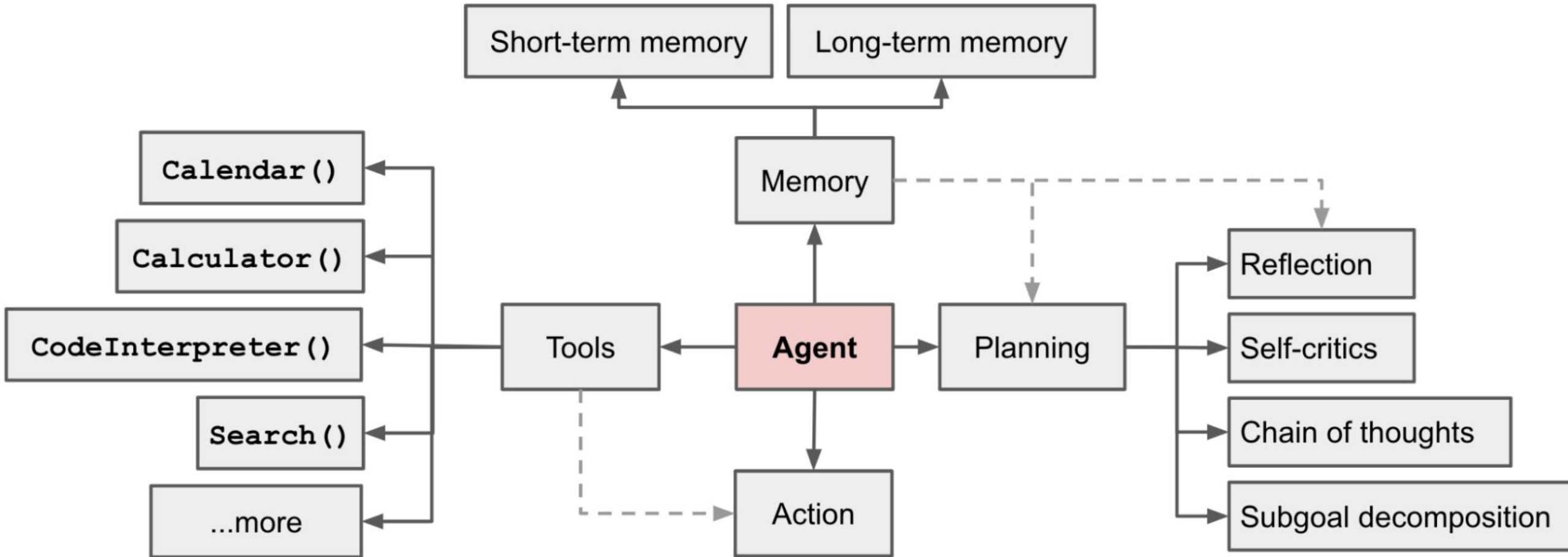
Trajectory annotations are scarce and expensive

e.g., DeepMind/Android/15k traj/833 app/20 ppl/4 mo



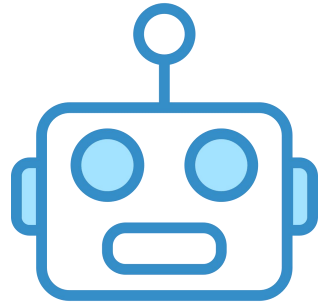
Challenge 2: from foundation model to ACTION model

Foundation models can read and write, but they can't act



Challenge 3: Safety and Interpretability

Help me clean my desktop



```
cd Desktop  
rm -rf ./
```

Web agent for shopping



Research Challenge

1. Scarcity of data

Trajectory annotations are scarce and expensive

e.g., DeepMind/Android/15k traj/833 app/20 ppl/4 mo

2. Foundation models can read and write, but they can't act

Special training is needed for digital grounding

3. Safety and Interpretability

Lack of transparency require more user trust

Research Challenge

1. Scarcity of data

Trajectory annotations are scarce and expensive

e.g., DeepMind/Android/15k traj/833 app/20 ppl/4 mo

2. Foundation models can read and write, but can't act

Specific training is needed for digital grounding

=> **Executable Language Grounding**

3. Safety and Interpretability

Lack of transparency require more user trust

Executable Language Grounding

Symbol-LLM (ACL'24) OS-Copilot (ICLR'24 WS)

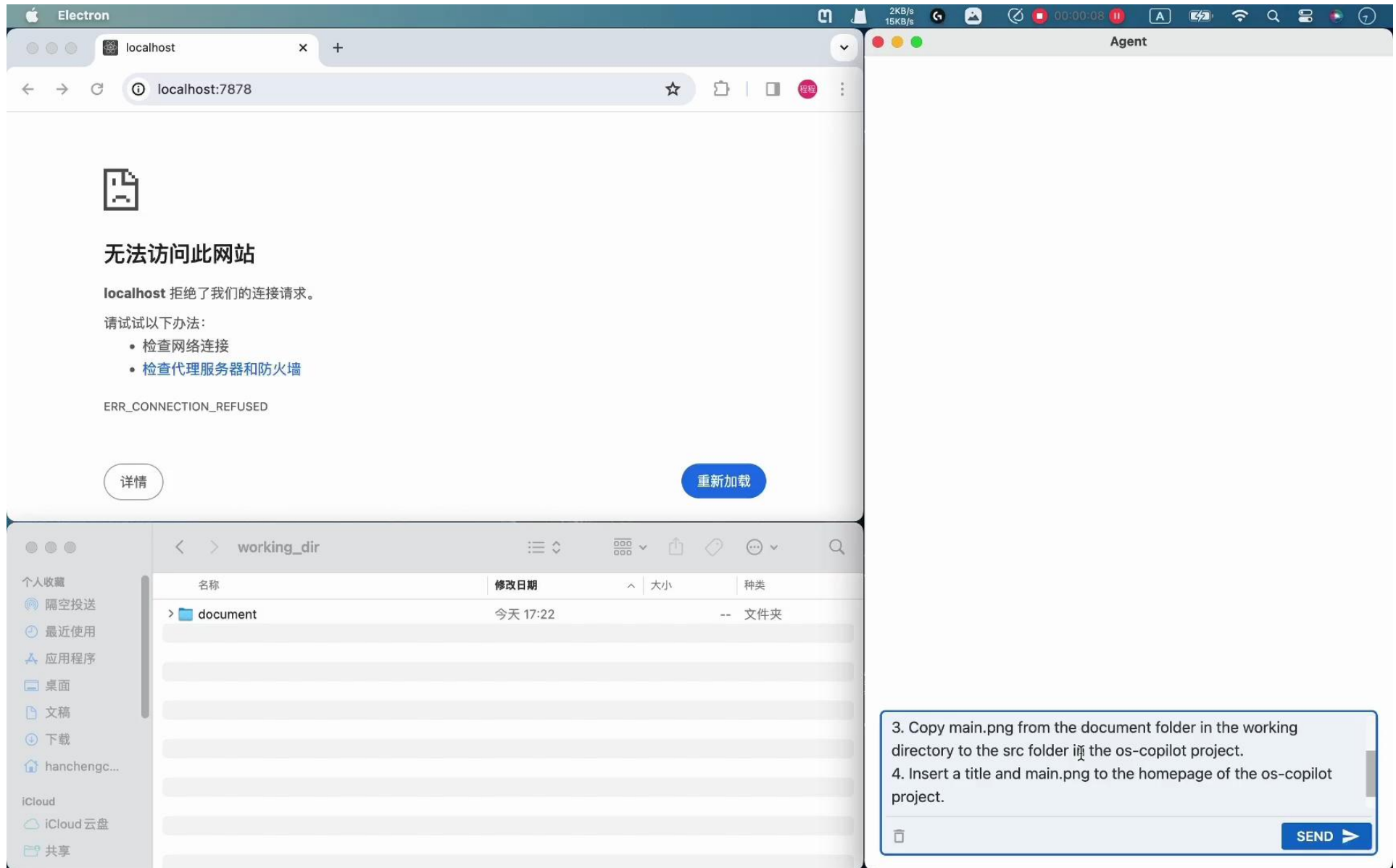


1. Grounding via code: turn natural language instructions into executable programs

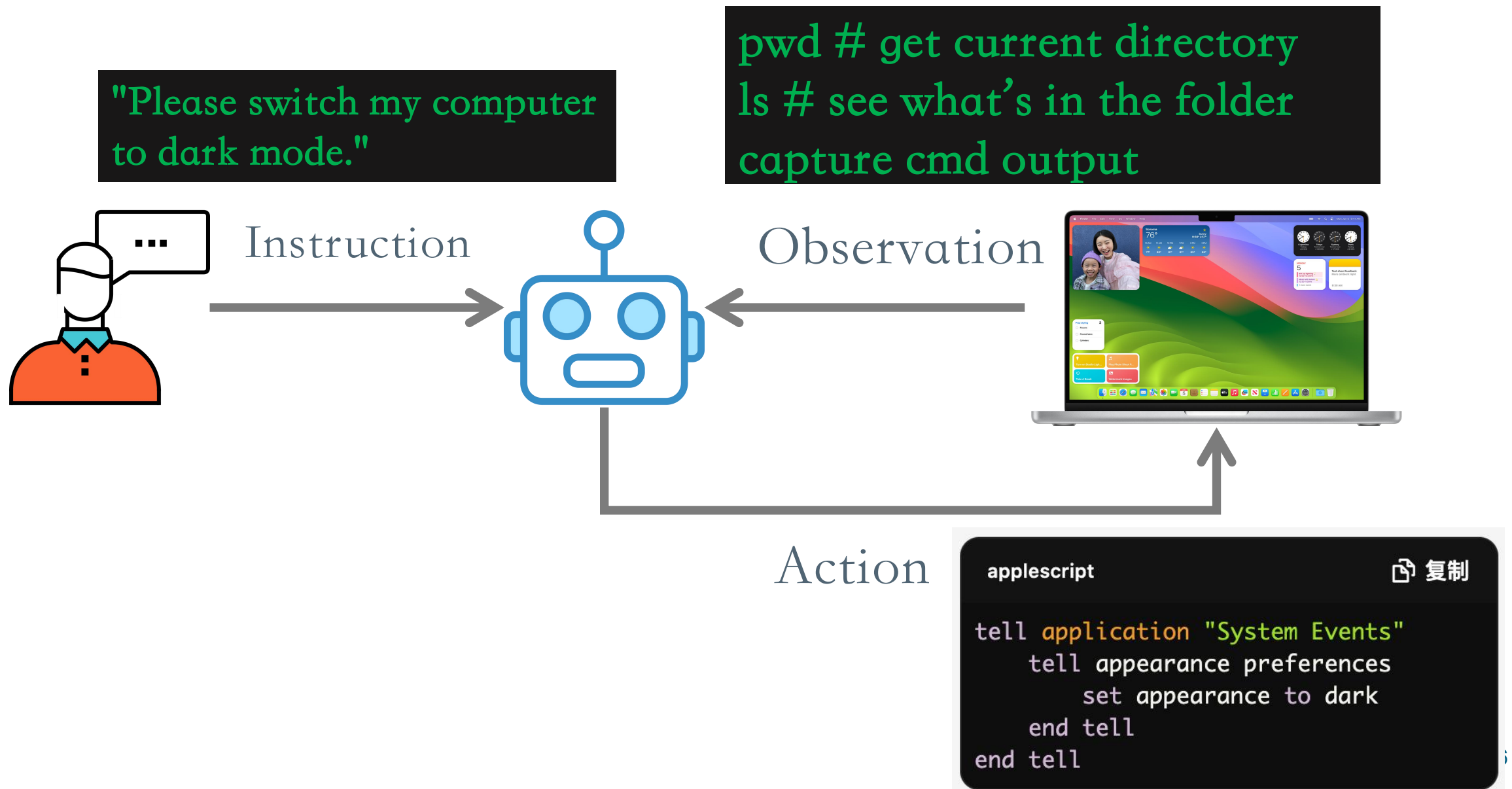
Example: "Please switch my computer to dark mode."

```
applescript 复制  
  
tell application "System Events"  
  tell appearance preferences  
    set appearance to dark  
  end tell  
end tell
```

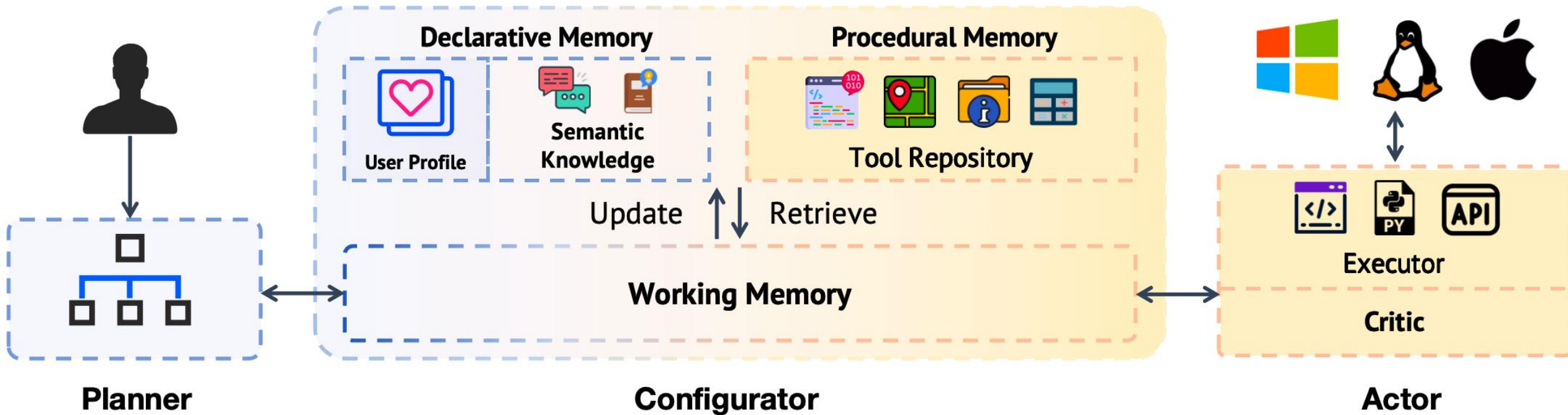
Demo: Building a React Website



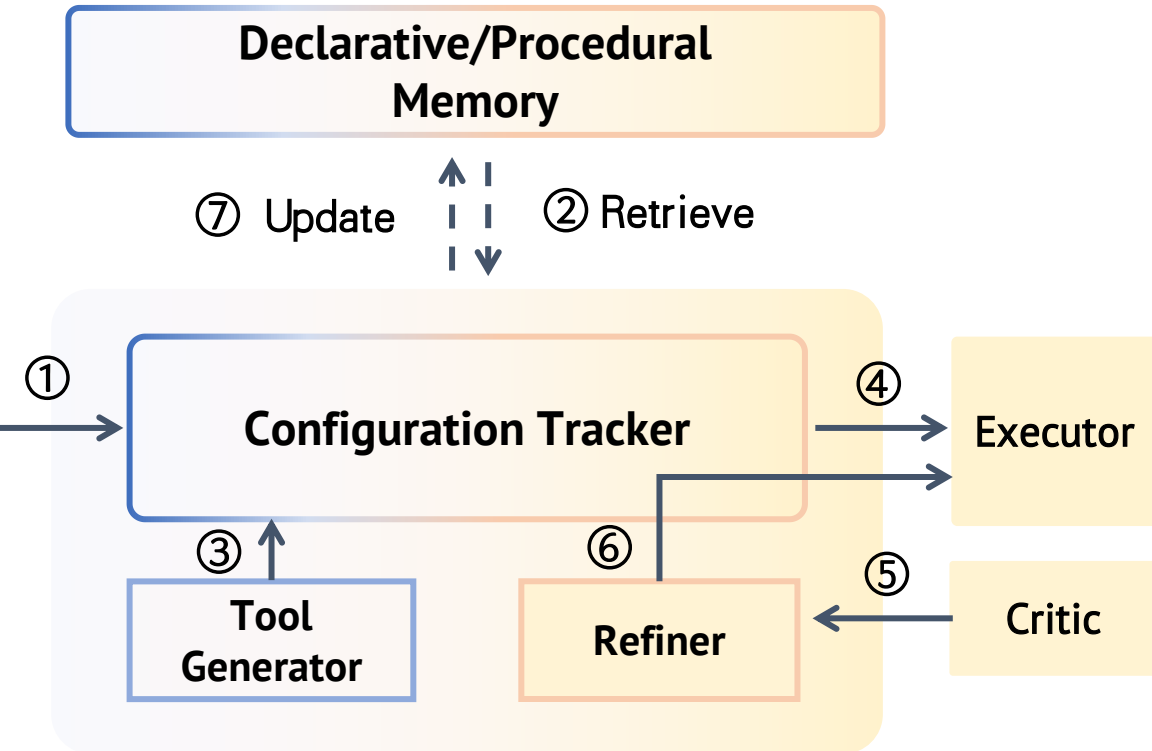
OS-Copilot System Overview



OS-Copilot System Overview



OS-Copilot: Configurator



Subtask: Change the system into the Dark mode

Tool Generator:

```
class change_system_appearance(BaseAction):  
    ...  
    script = 'tell app "System Events" to tell appearance  
preferences to set dark mode to true'  
    ...
```

Executor:

- (i) Save the tool to `change_system_appearance.py`
- (ii) Execute the tool

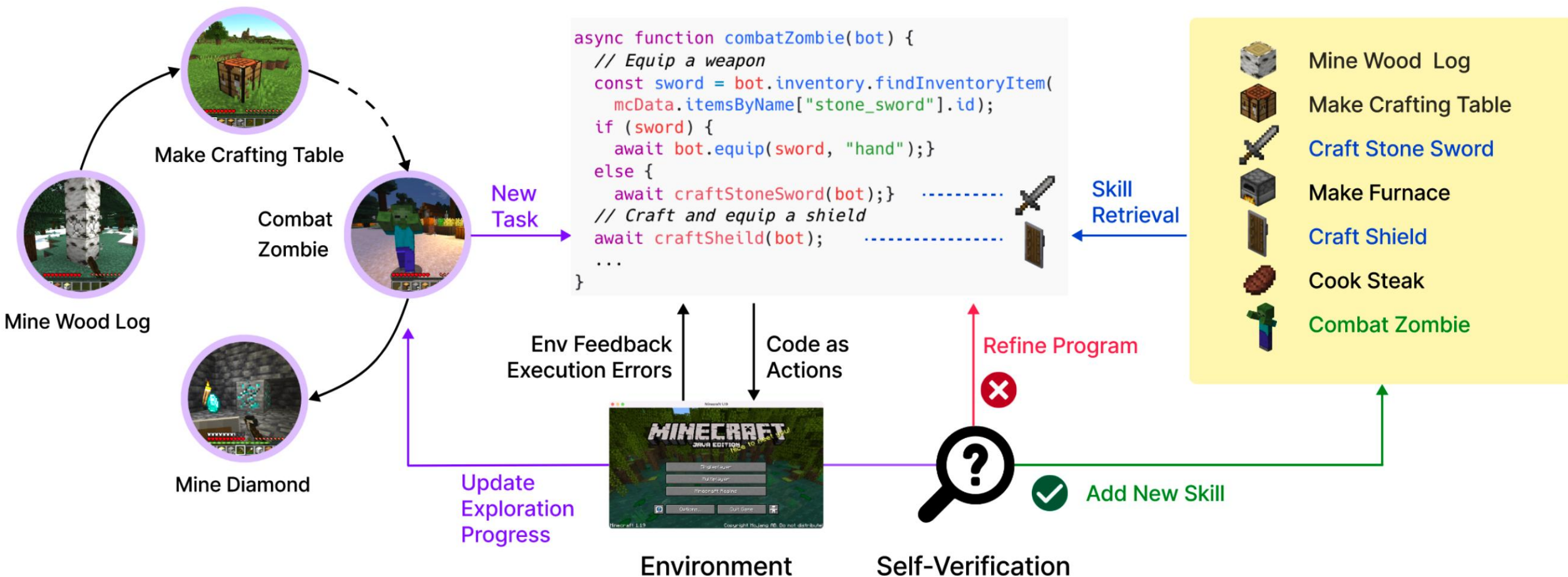
```
> _ python change_system_appearance.py dark
```

Populate the memory through self-directed learning

Automatic Curriculum

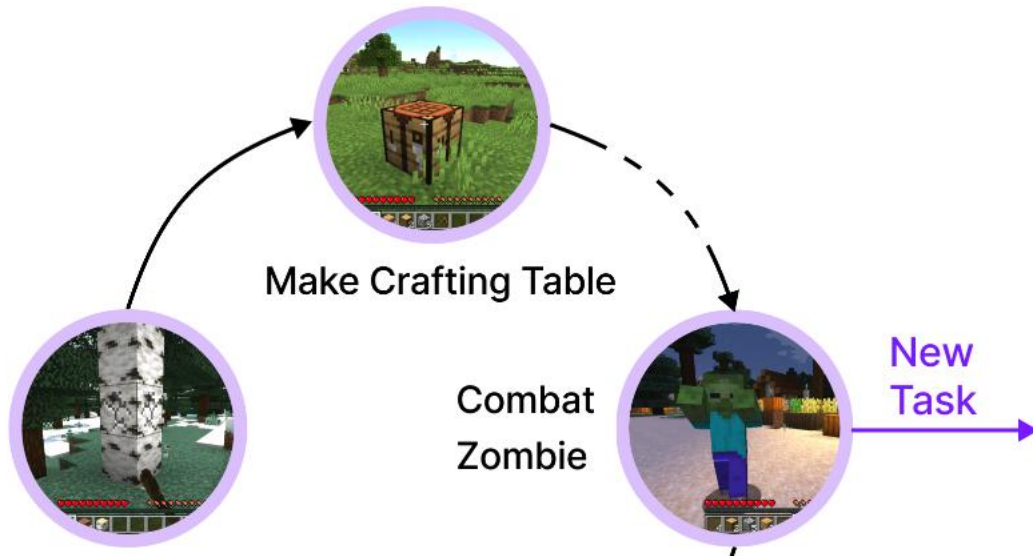
Iterative Prompting Mechanism

Skill Library



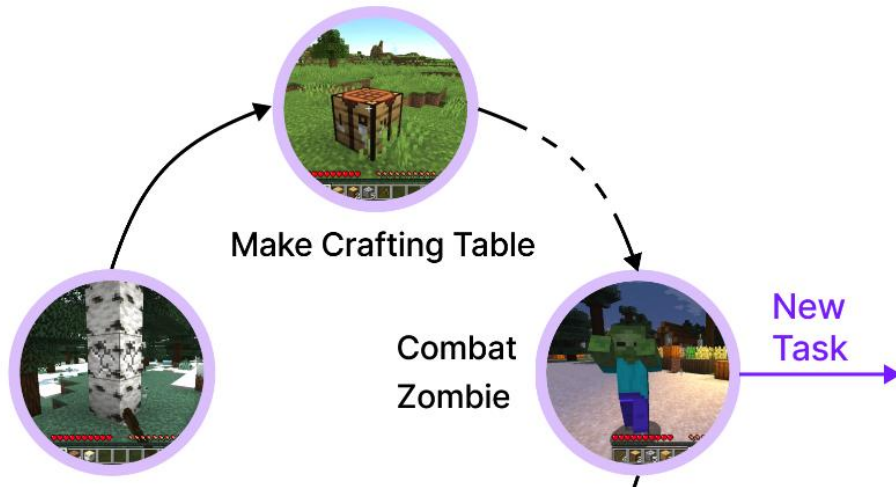
Populate the memory through self-directed learning

Automatic Curriculum



Populate the memory through self-directed learning

Automatic Curriculum



Iterative Prompting Mechanism

```
async function combatZombie(bot) {  
  // Equip a weapon  
  const sword = bot.inventory.findInventoryItem(  
    mcData.itemsByName["stone_sword"].id);  
  if (sword) {  
    await bot.equip(sword, "hand");  
  }  
  else {  
    await craftStoneSword(bot); .....  
  }  
  // Craft and equip a shield  
  await craftShield(bot); .....  
  ...  
}
```



Skill
Retrieval

Skill Library

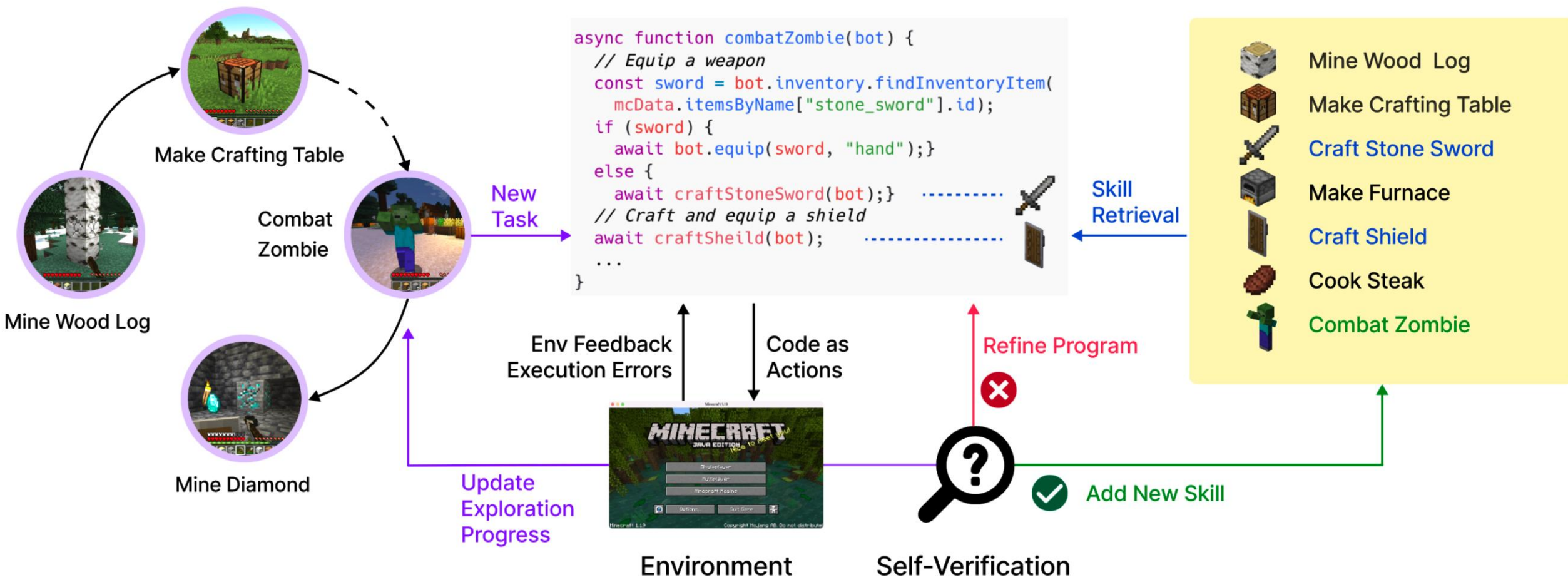
-  Mine Wood Log
-  Make Crafting Table
-  Craft Stone Sword
-  Make Furnace
-  Craft Shield
-  Cook Steak

Populate the memory through self-directed learning

Automatic Curriculum

Iterative Prompting Mechanism

Skill Library



1. Grounding via Code

Symbol-LLM (ACL'24) OS-Copilot (ICLR'24 WS)



1. Grounding via code:

Example: "Please switch my computer to dark mode."

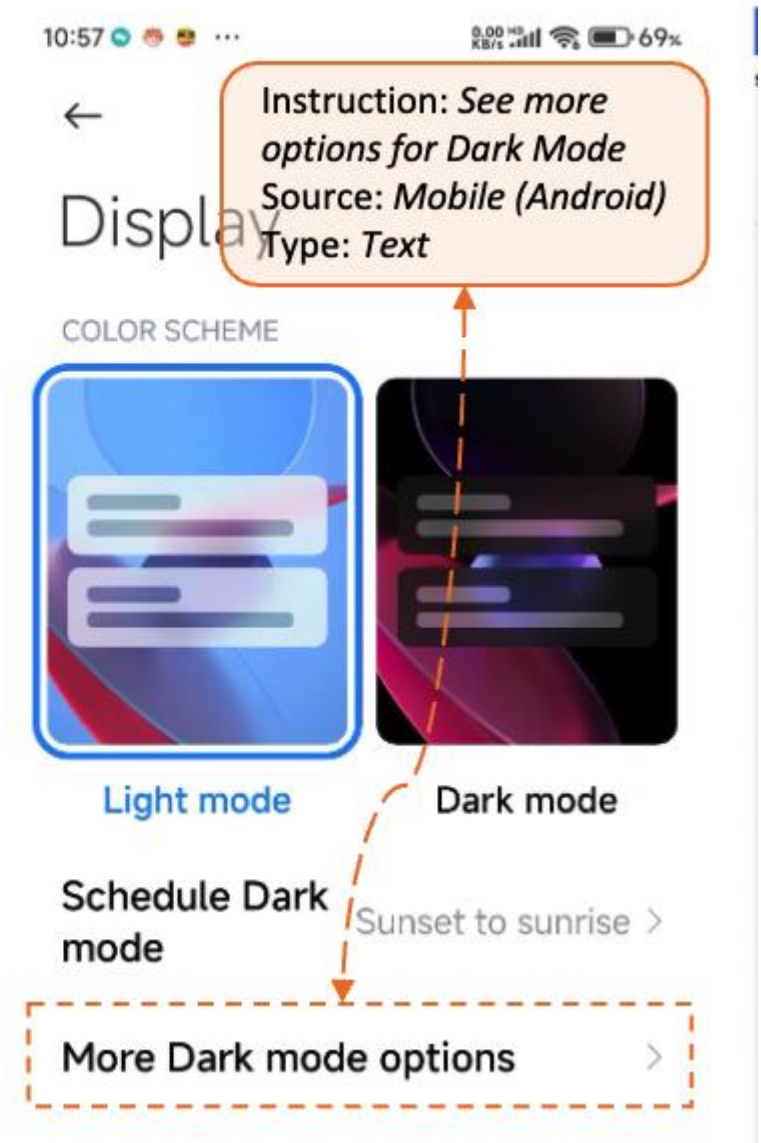
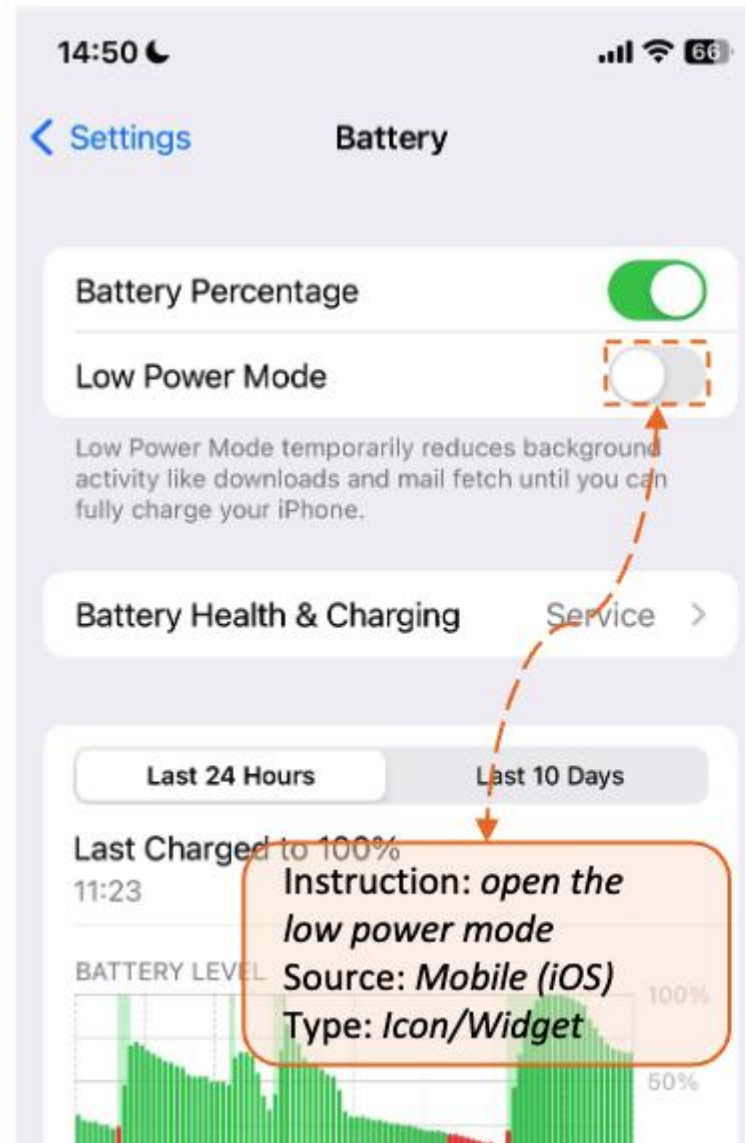
```
applescript 复制  
  
tell application "System Events"  
    tell appearance preferences  
        set appearance to dark  
    end tell  
end tell
```

CON: **Most software is not open sourced!**

2. Grounding via GUI Interaction

2. Grounding via GUI interaction: Turn instruction in human-like GUI interaction (e.g., Click)

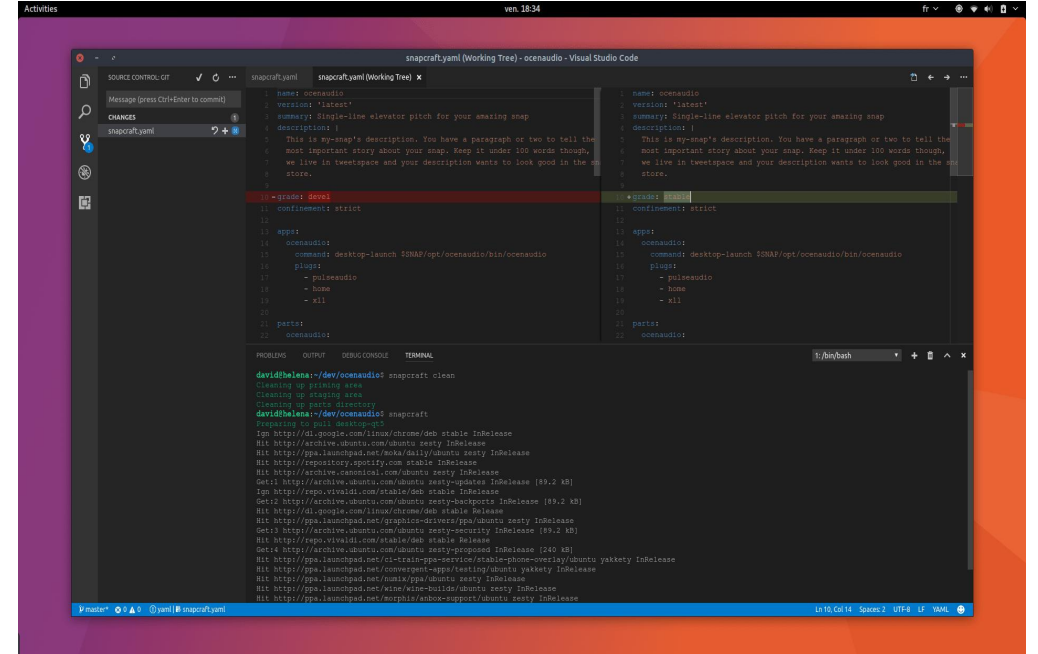
Challenge: understand and locate GUI element



2. Grounding via GUI Interaction



VS

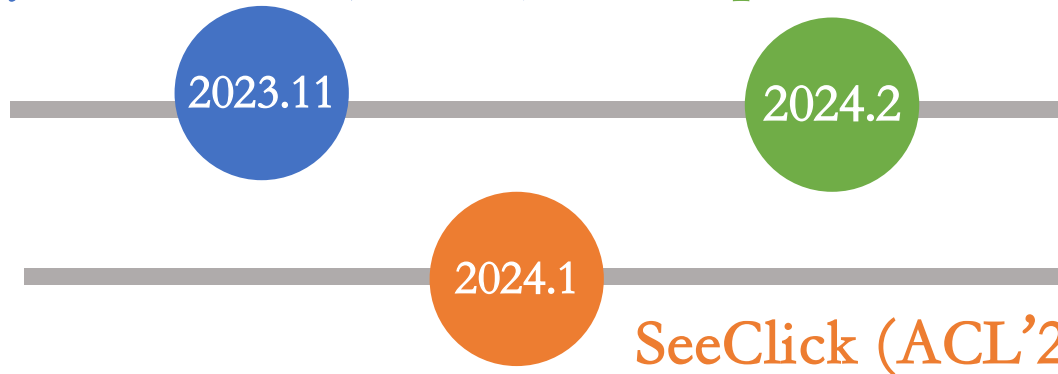


Current VLMs are mostly trained with natural images

OS Agent often deals with “unnatural” images

Executable Language Grounding

Symbol-LLM (ACL'24) OS-Copilot (ICLR'24 WS)



CON: **low-efficiency,**
do we really need human-like agents?

	Pros	Cons
GUI	Highly generalizable	Low-efficiency, Upperbounded by current VLM
CLI	High-efficiency	Poor generalizability

COMP7607 Final Project: Building Your Own JARVIS

Teams:

5-8 students

Topics:

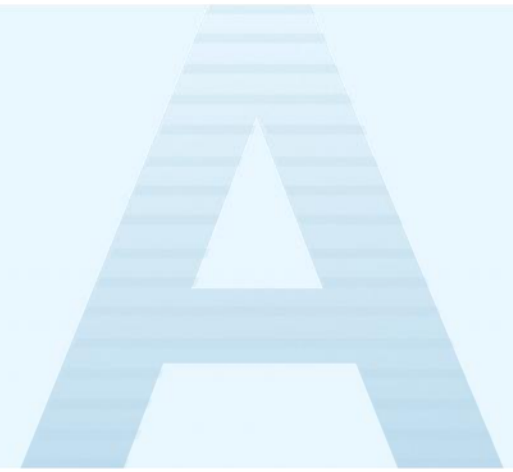
Using OS-Copilot as the base,
creating a fully functional OS agent
(JARVIS-like digital copilot)

Contributing to OS-Copilot by
making a PR/MR!



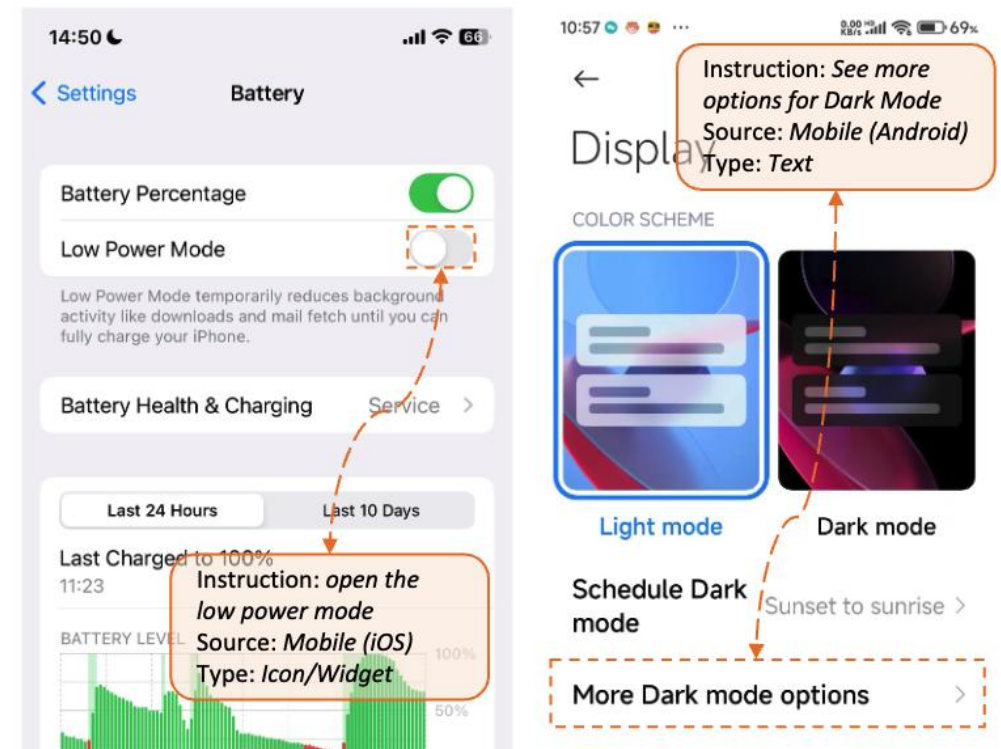
JARVIS: Render complete in 5 minutes.

| Exemplary Projects



Exemplary Project 1: Mobile Agent

Extending the current OS-Copilot framework to support mobile control

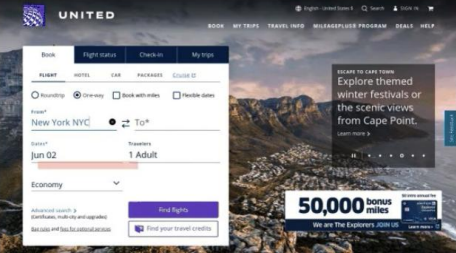


1. SeeClick: Harnessing GUI Grounding for Advanced Visual GUI Agents. Cheng et al., 2024
2. Mobile-Agent: The Powerful Mobile Device Operation Assistant Family. Wang et al., 2024
3. On the Effects of Data Scale on Computer Control Agents. Li et al., 2024
4. <https://www.youtube.com/watch?v=EMblpzqJld0>

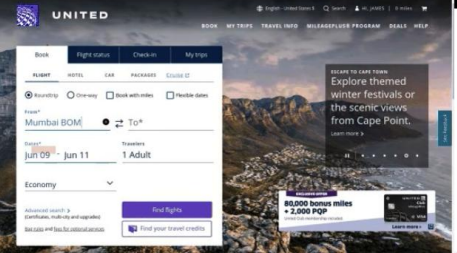
Exemplary Project 2: Web Agent

Improving the current framework to focus on website navigation and control

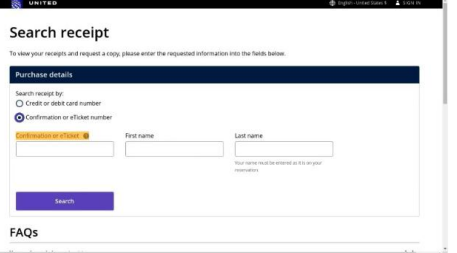
(a) Find one-way flights from New York to Toronto.



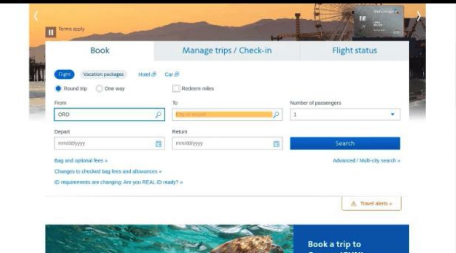
(b) Book a roundtrip on July 1 from Mumbai to London and vice versa on July 5 for two adults...




(c) Search receipt with the eTicket 12345678 for the trip reserved by Jason Two



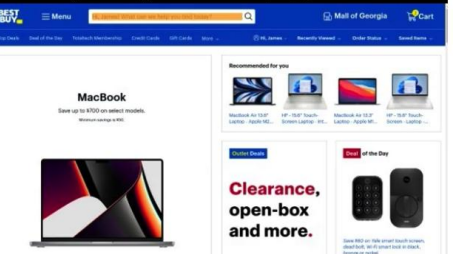
(d) Find a flight from Chicago to London on 20 April and return on 23 April.




(e) Search for the interactions between ibuprofen and aspirin.



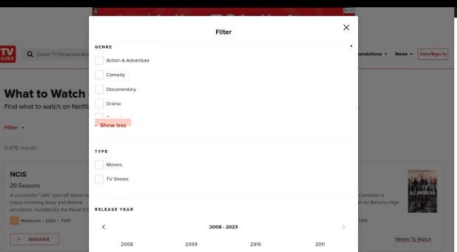
(f) As a Verizon user, finance a blue iPhone 13 with 256gb along with monthly apple care.



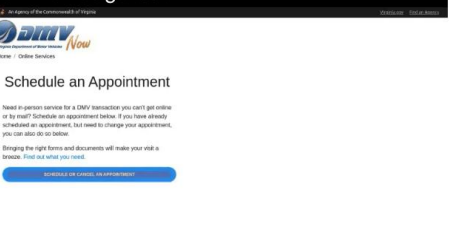
(g) Find Elon Musk's profile and start following, start notifications and like the latest tweet.



(h) Browse comedy films streaming on Netflix that was released from 1992 to 2007.



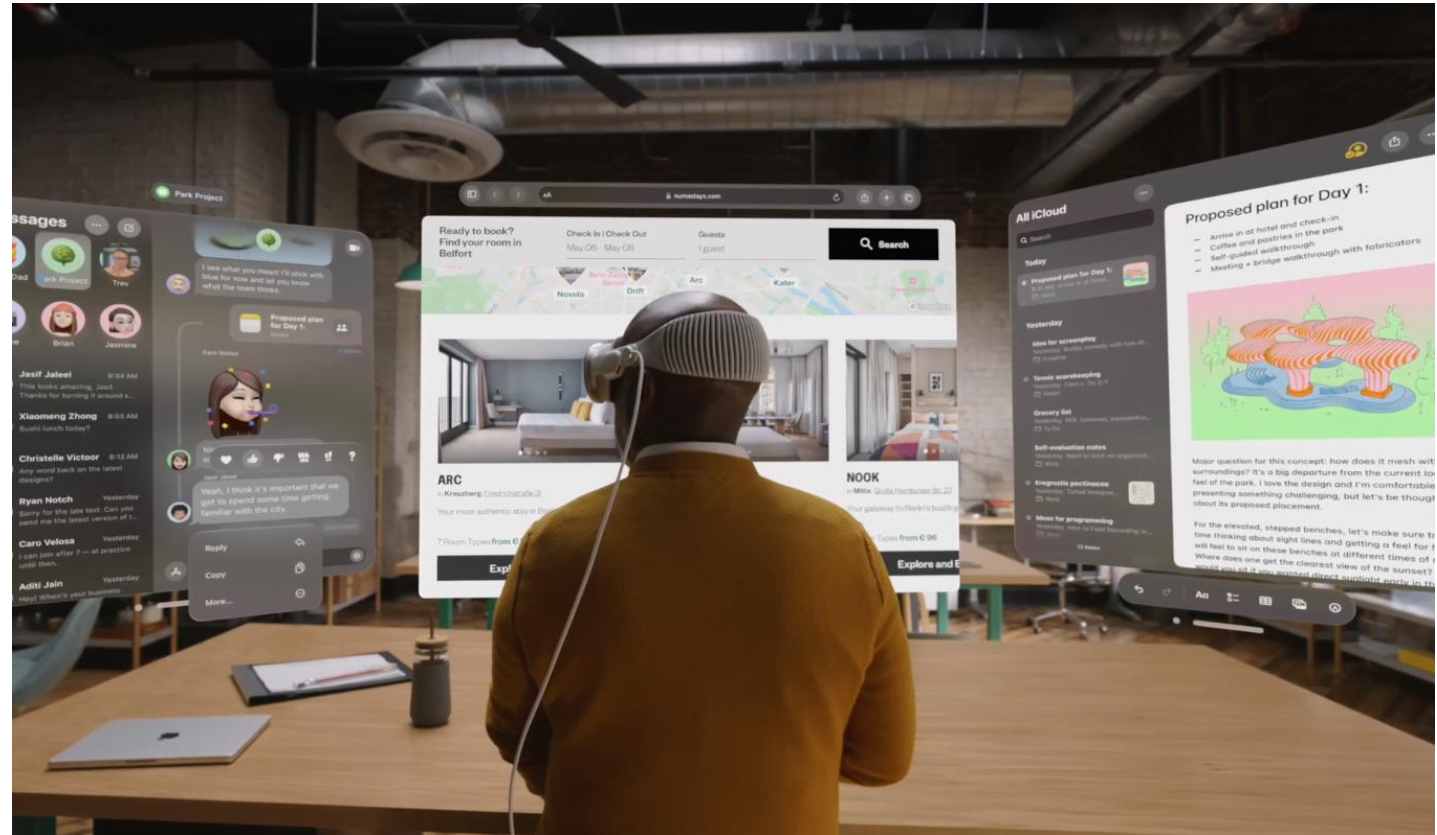
(i) Open page to schedule an appointment for car knowledge test.



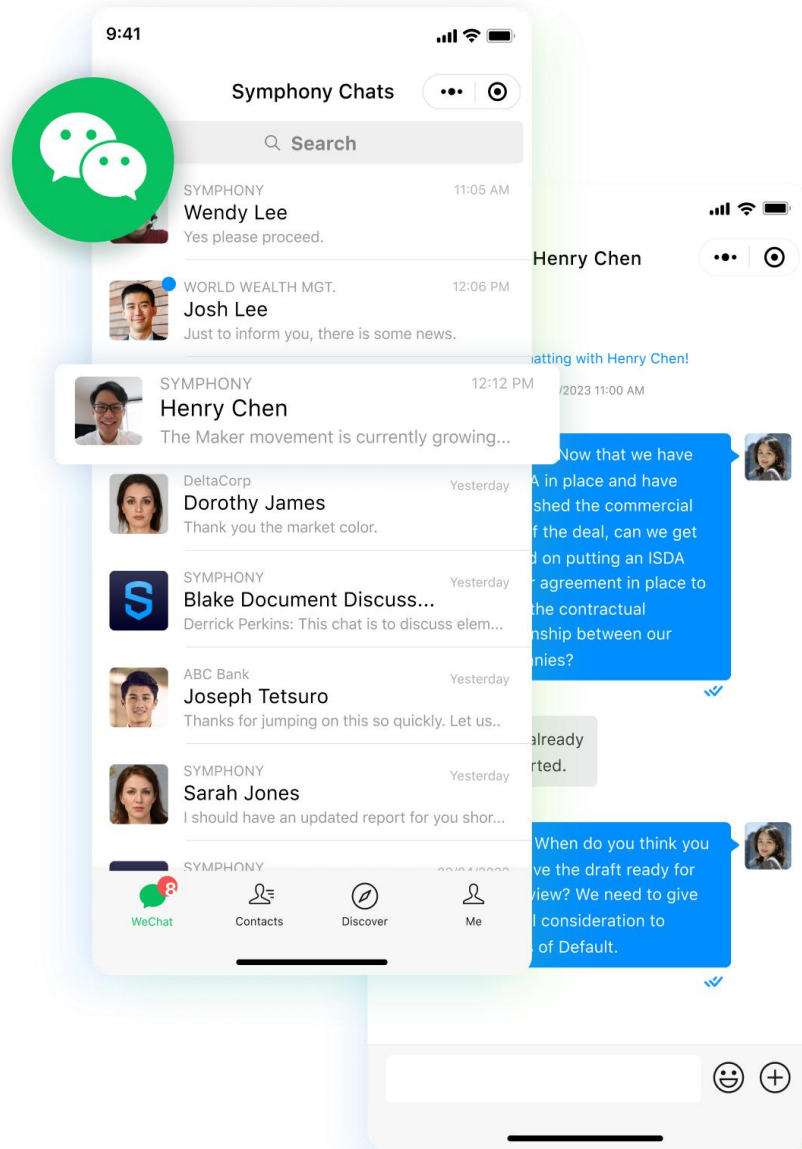
1. <https://osu-nlp-group.github.io/Mind2Web/>
2. <https://www.multion.ai/>

Exemplary Project 3: VR-based Agent

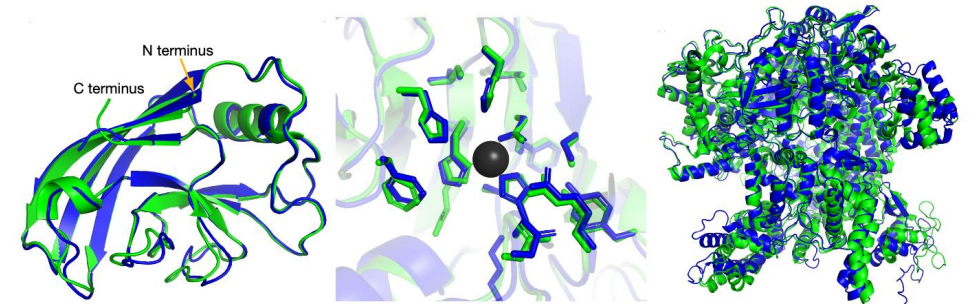
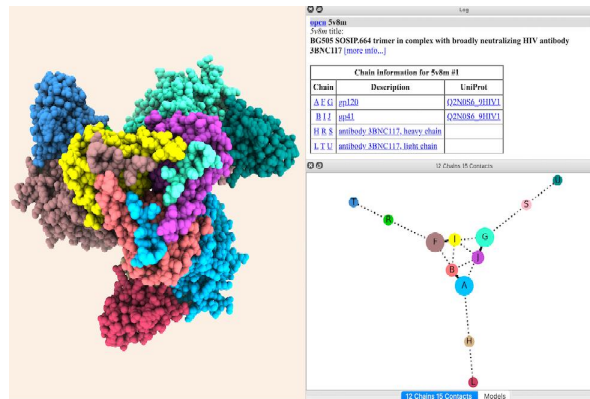
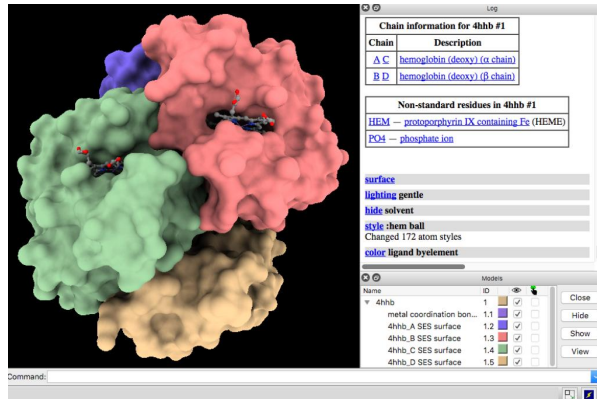
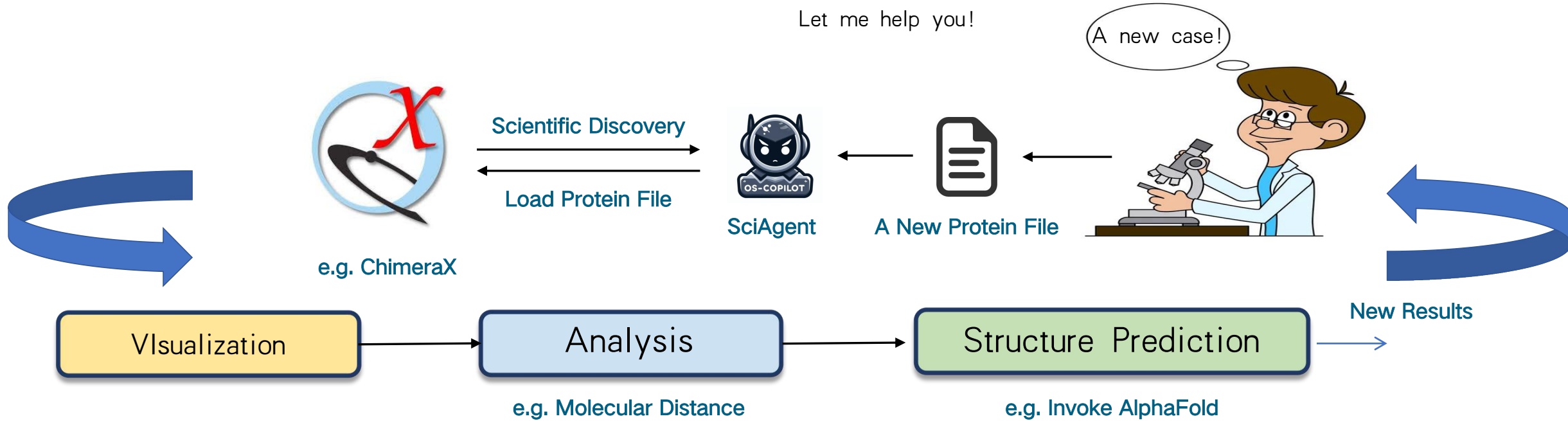
Extending the current framework to support VR platform (with voice control)



Exemplary Project 4: App-specific Agent

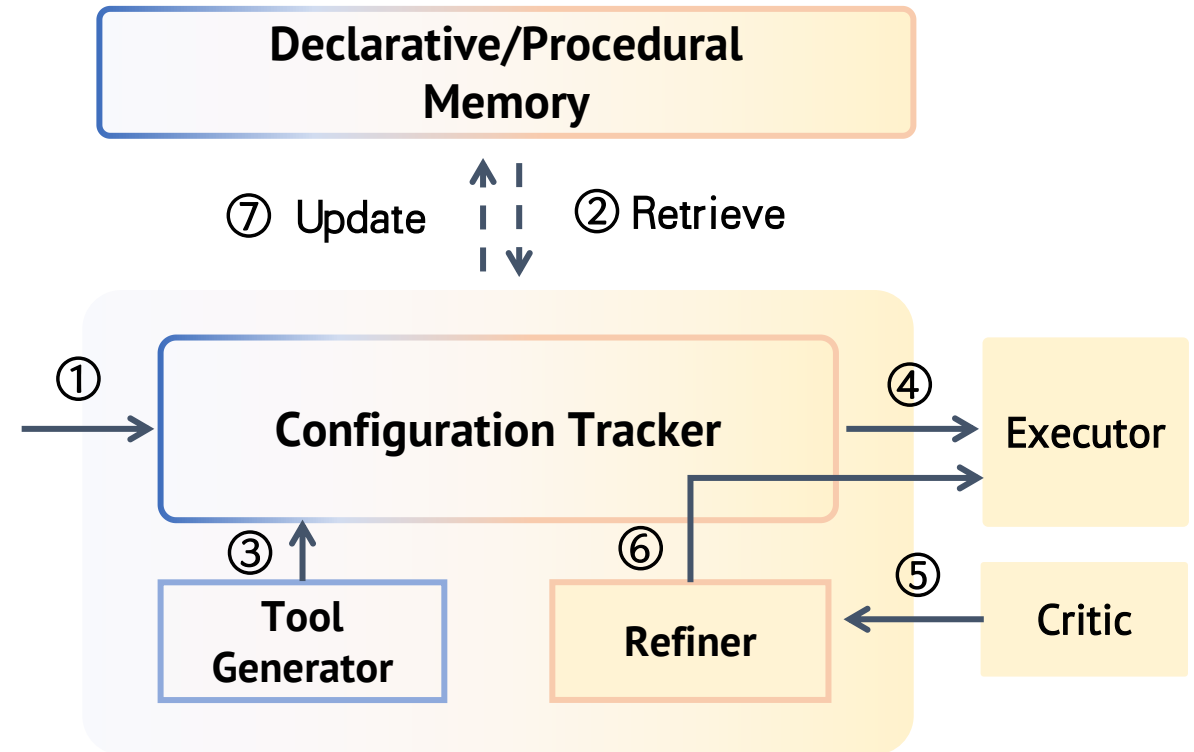


Exemplary Project 5: Agent for Scientific Discovery



Exemplary Project 6: New Methodology

Support GUI interaction
Improve tool learning
Support more long-term memory
Personalization
Better planner
Better self-refine
....
Safety...
Interpretability...



Looking forward to exciting demonstrations from all of you!



Deliverables

1. Project proposal (0%):

Up to one page proposal containing team information and brief introduction of the project

2. Presentation and demonstration (35%):

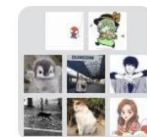
Present your project in the class and showcase your exciting demos!

Evaluated based on excitement(15%), soundness(15%), and presentation clearness(5%).

3. Project Report (15%):

4-8 pages report concludes the project.

- OS-Copilot
 - Paper: <https://arxiv.org/pdf/2402.07456>
 - Code: <https://github.com/OS-Copilot/OS-Copilot>
 - Make PR/MR and become a contributor!
- Reading List:
 - released with the handout



群聊: COMP 7607 OS-Copilot
答疑群



该二维码7天内(9月15日前)有效, 重新进入将更新