# Windows Agent Arena:
## Evaluating Multi-Modal OS Agents at Scale

https://microsoft.github.io/WindowsAgentArena/

**Rogerio Bonatti, Dan Zhao,**
**Francesco Bonacci, Dillon Dupont, Sara Abdali, Yinheng Li, Yadong Lu,**
**Justin Wagle, Kazuhito Koishida, Arthur Bucker, Lawrence Jang, Zack Hui**

# What are agents? What can they do?

Reason

Prompt / objective

Your favorite LLM/VLM

Agent

Output

1. Lorem ipsum dolor sit amet, consectetur adipiscing elit, sed do eiusmod tempor incididunt ut labore et dolore magna aliqua.
2. Ut enim ad minim veniam, quis nostrud exercitation ullamco laboris nisi ut aliquip ex ea commodo consequat.
3. Duis aute irure dolor in reprehenderit in voluptate velit esse cillum dolore eu fugiat nulla pariatur.
4. Excepteur sint occaecat cupidatat non proident, sunt in culpa qui officia deserunt mollit anim id est laborum.

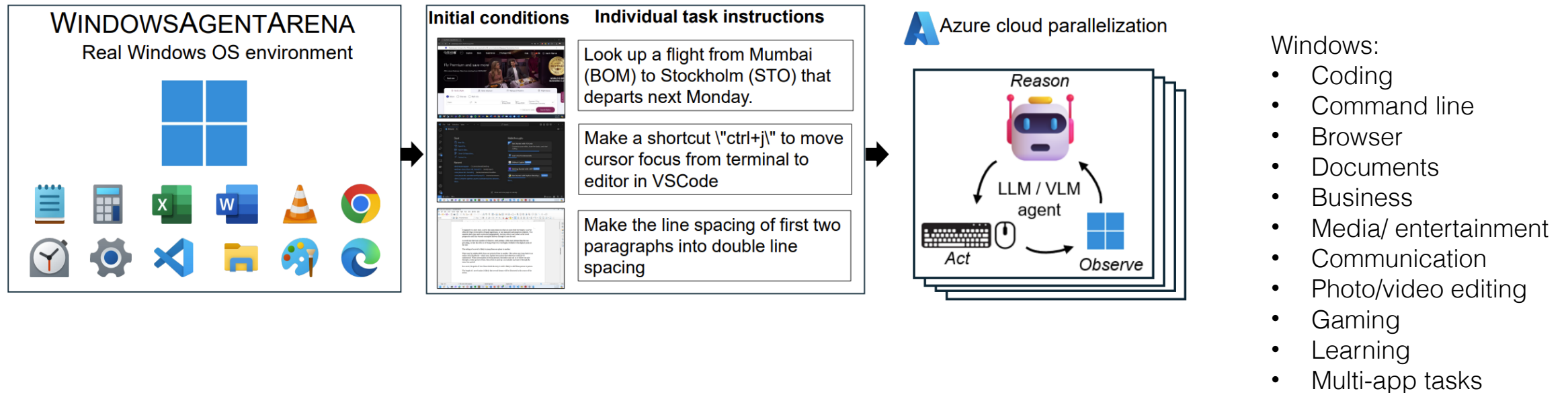Observe

Act

Environment:
- Codebase
- Browser
- Computer

**Challenge 1:** Generalist OS agents are a superset of browser and domain-specific agents.

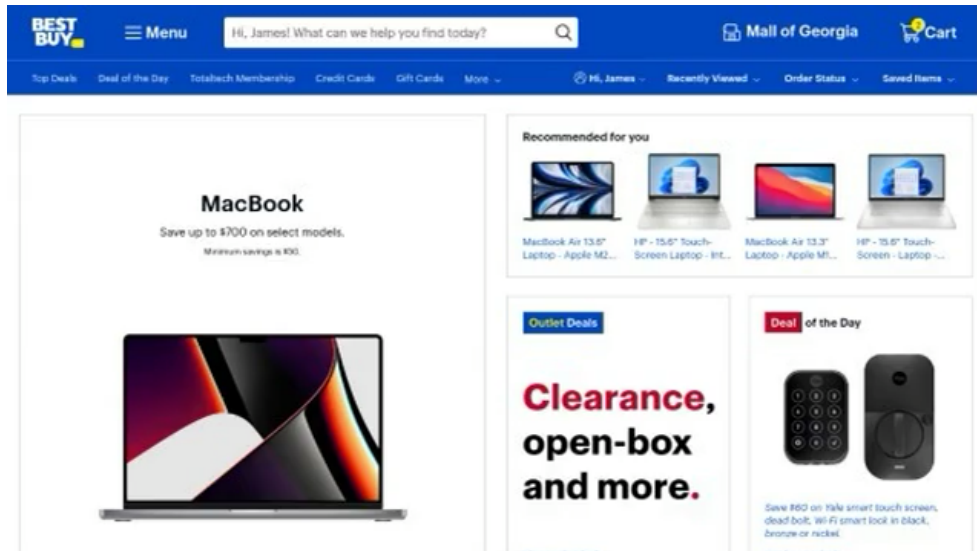But where's Windows OS?

First agent benchmark for Windows OS



WINDOWSAGENTARENA
Real Windows OS environment

Initial conditions | Individual task instructions

Look up a flight from Mumbai (BOM) to Stockholm (STO) that departs next Monday.

Make a shortcut \"ctrl+j\" to move cursor focus from terminal to editor in VSCode

Make the line spacing of first two paragraphs into double line spacing

Azure cloud parallelization

Reason
LLM / VLM agent
Act          Observe

Windows:
- Coding
- Command line
- Browser
- Documents
- Business
- Media/ entertainment
- Communication
- Photo/video editing
- Gaming
- Learning
- Multi-app tasks

# Challenge 2: You can't improve what you can't measure. But how do we measure a computer agent's performance?

- ## Static trajectory matching

*Task: finance a blue iPhone 13 with 256gb*

Human demonstration (Mind2Web benchmark):



- ## Execution-based evaluation

*Task: finance a blue iPhone 13 with 256gb*



Black box agent execution

Trajectory is irrelevant. Success iff:
- iPhone added to cart
- iPhone memory = 256gb
- Payment method is financing

# Challenge 3: execution-based evaluation is (very) slow

- Total eval time = # tasks * # steps * # seconds per step
  - (20)          (30 s)

|  | Series | Parallel |
|---|---|---|
| 100 tasks | 17 hours | 10 min |
| 1,000 tasks | 1 week | 10 min |
| 10,000 tasks | 2.5 months | 10 min |
| 1M trajectories | 19 years | 10 min |

Benchmark sizes today

Ideal benchmark size

Ideal dataset size

# Windows Agent Arena <span style="color:red">contributions</span>

- First agent benchmark for Windows OS
- Execution-based evaluation
- 150+ tasks across 11 domains (browser, 1P/3P apps)
- Parallel evaluation in Azure in under 20min
- Open-source benchmark. Check it out!

# Problem definition

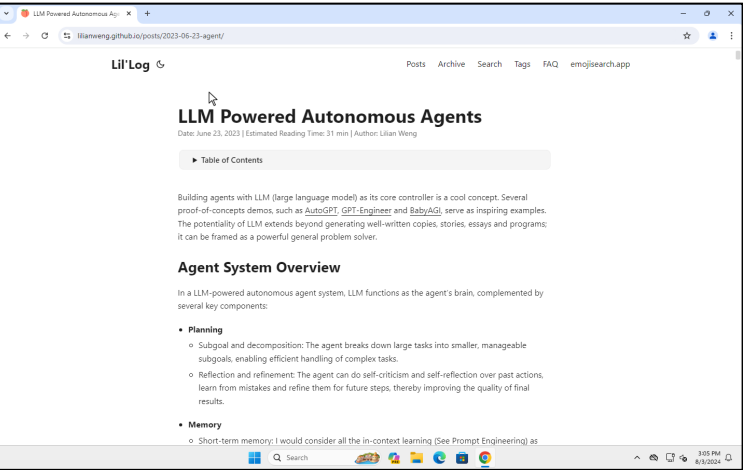- POMDP: $S, \mathcal{O}, \mathcal{A}, T, \mathcal{R}$
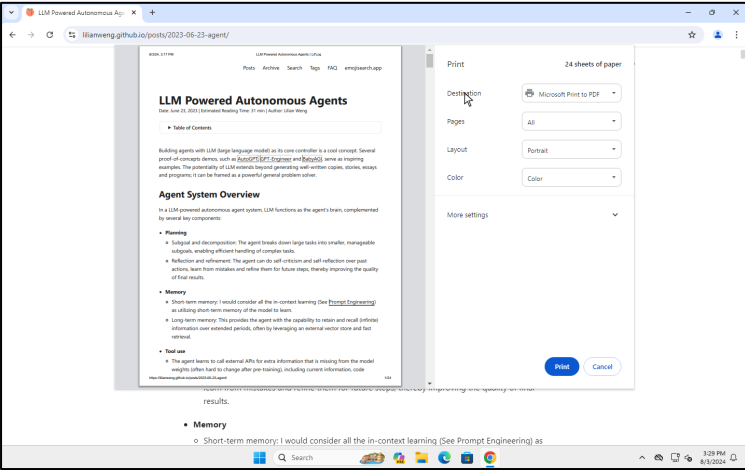
- Observation:



Pixels +UIA

- Action: 

- Transition function $T : S \times \mathcal{A} \rightarrow S$

- Reward: $\mathcal{R} : S \times \mathcal{A} \rightarrow \mathbb{R}$

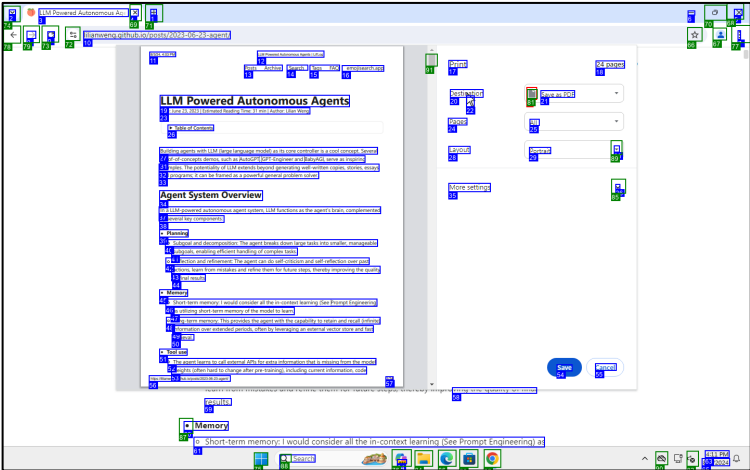| Group | Function |
|-------|----------|
| computer.mouse | move_id(id)<br>move_abs(x,y)<br>single_click()<br>double_click()<br>right_click()<br>scroll(direction) |
| computer.keyboard | write(text)<br>press(key) |
| computer.clipboard | copy_text(text)<br>copy_image(image)<br>paste() |
| computer.os | open_program(program) |
| computer.window_manager | switch_to_application(window) |

# Task: *Computer, can you turn the webpage I'm looking at into a PDF file and put it on my main screen, you know, the Desktop?*
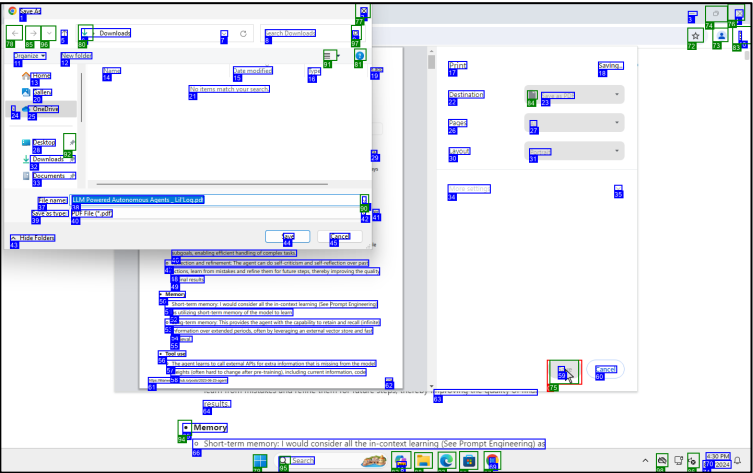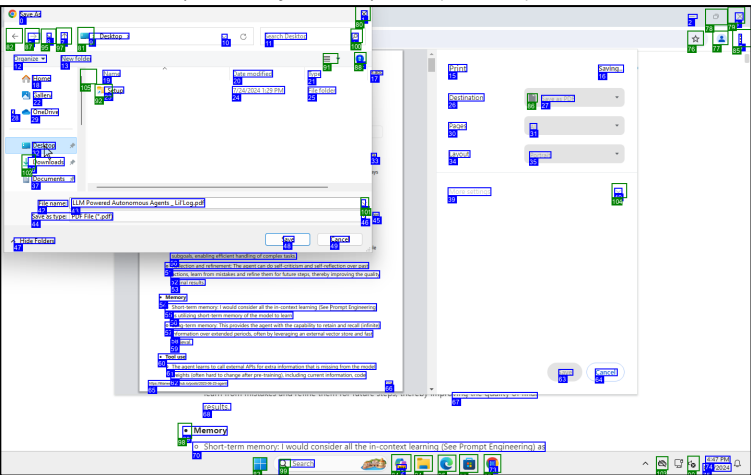


Step 1: computer.keyboard.press("ctrl+p")



Step 2:
```
# Move the focus to the 'Destination' dropdown
computer.keyboard.press("tab")
# Attempt to navigate the dropdown with the keyboard
computer.keyboard.press("down")
```
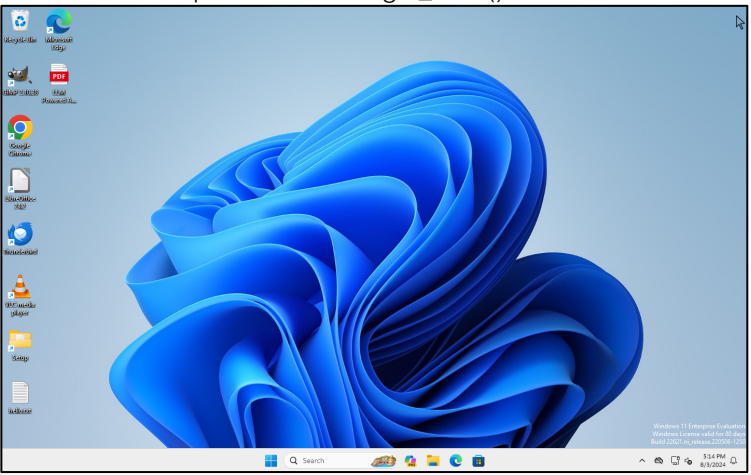


Step 3:
```
# Move to the "Save" button
computer.mouse.move_id(id=54)
# Click to save the webpage as a PDF
computer.mouse.single_click()
```



Step 4:
```
# Move to "Desktop" (pinned) to change the
location computer.mouse.move_id(id=28)
computer.mouse.single_click()
```
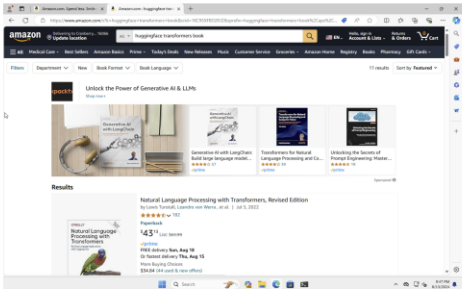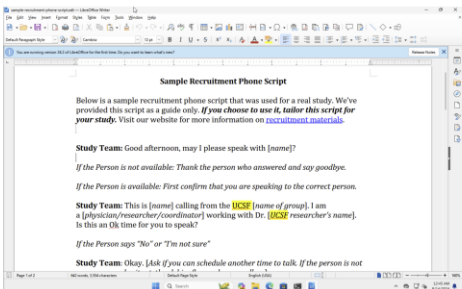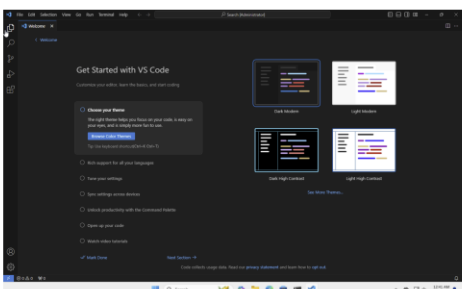


Step 5:
```
# Move cursor to the "Save" button
computer.mouse.move_id(id=48)
# Click to save the webpage as a PDF to the
Desktop
```



Result: Task successful, reward 1.0
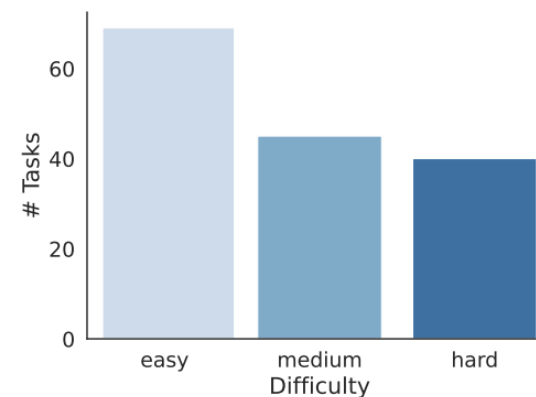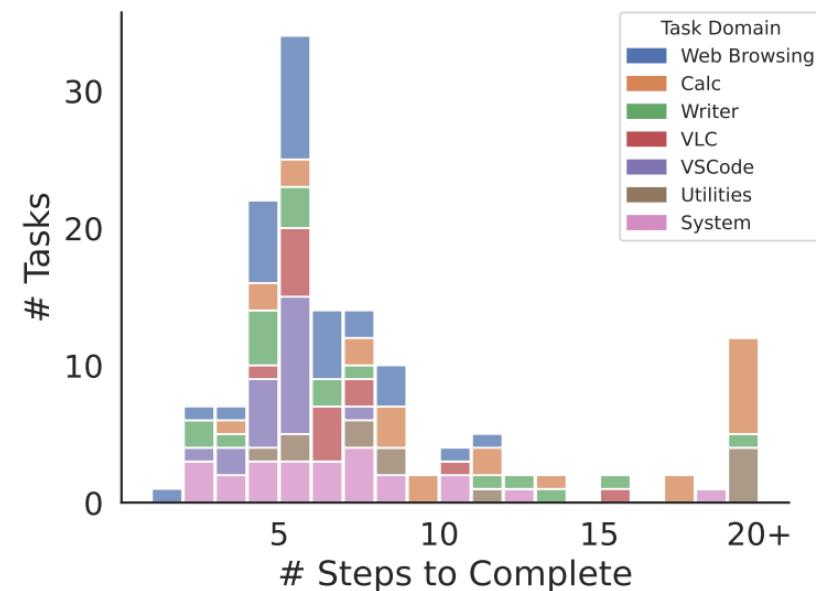
# ✔ Execution-based evaluation

| Initial State | Task Instruction | Evaluation Script (Simplified) |
|---|---|---|
|  | "Can you help me clean up my computer by getting rid of all the tracking things that Amazon might have saved? I want to make sure my browsing is private and those sites don't remember me." | ```
cookie_data = get_cookie_data(env)
rule = {"type":"rule",
        "domains": ["amazon.com"]}
is_cookie_deleted(cookie_data, rule)
``` |
|  | "I have been editing my document and some words that needed to be rewritten are highlighted in yellow. As I fix those words, please help me remove all highlight. I want to make sure that there is no highlighted word." | ```
vm_file = get_file(env, "file.doc")
golden_file = get_file(drive,
                       "file.doc")
check_highlighted_words(vm_file,
                        golden_file)
``` |
|  | "Please help me modify the setting of VS Code to keep my cursor focused on the debug console when debugging in VS Code, instead of automatically focusing back on the Editor." | ```
vscode_data = get_vscode_config(env)
rule = {"expected":
  {"debug.focusEditorOnBreak": false}
        }
check_json_settings(vscode_data, rule)
``` |

# ⊕ 154 tasks across 12 domains 📝🧮📊📘🔺🔴

*Unassisted human performance = 74% success rate*

| Domain | Activity Examples | # Tasks | % |
|---|---|---|---|
| Office | Libreoffice Writer (Editing, writing) and LibreOffice Calc (spreadsheets, plotting, data manipulation) | 43 | 28% |
| Web Browsing | Online shopping & search, customizing app settings with Edge and Chrome | 30 | 19% |
| Windows System | File Explorer (navigate, customize files), Windows Settings, customization | 24 | 16% |
| Coding | Editing code, customizing IDE, installing extensions with VSCode | 24 | 16% |
| Media & Video | Watching videos in VLC, listening to music, settings | 21 | 14% |
| Windows Utilities | Miscellaneous tasks with Notepad, Clock, and Paint | 12 | 8% |
| **Total** | | **154** | **100%** |

# Secure cloud parallelization in Azure ML



- Golden snapshot copied for each VM
- Full separation between instances
- Security: no external VM commands

# Results

| UIA tree | OCR | Icon det. | Image det. | Model | Office | Web Browser | Windows System | Coding | Media & Video | Windows Utils | Total |
|---|---|---|---|---|---|---|---|---|---|---|---|
| ✗ | Pytesseract + DOM | | Grounding DINO | Phi3.5-V | 0.0% | 4.8% | 0.0% | 0.0% | 0.0% | 0.0% | 0.9% |
| | | | | Phi3-V | 0.0% | 0.0% | 4.2% | 4.3% | 0.0% | 0.0% | 1.3% |
| | | | | GPT-4o-mini | **2.3%** | 6.7% | 12.5% | 8.3% | 14.6% | 0.0% | 7.2% |
| | | | | GPT-4o | 0.0% | 0.0% | 29.2% | 0.0% | 5.0% | 0.0% | 5.2% |
| | | | | GPT-4V-1106 | 0.0% | 10.3% | 21.3% | 12.5% | 9.8% | 0.0% | 8.6% |
| ✓ | Pytesseract+ DOM | | Grounding DINO | Phi3.5-V | 0.0% | 3.6% | 4.2% | 0.0% | 0.4% | 0.0% | 1.4% |
| | | | | Phi3-V | 0.0% | 0.0% | 4.2% | 4.3% | 5.0% | 0.0% | 2.0% |
| | | | | GPT-4o-mini | 0.0% | 3.3% | 8.3% | 0.0% | 5.0% | 0.0% | 2.6% |
| | | | | GPT-4o | 0.0% | 10.0% | 29.2% | 8.3% | 14.6% | 0.0% | 9.8% |
| | | | | GPT-4V-1106 | 0.0% | 13.3% | 25.0% | 13.0% | 28.9% | **8.3%** | 13.1% |
| ✗ | OneOCR | | Proprietary models | Phi3.5-V | 0.0% | 4.3% | 29.6% | 0.0% | 0.0% | 0.0% | 7.0% |
| | | | | Phi3-V | 0.0% | 6.7% | 8.3% | 0.0% | 4.8% | 0.0% | 3.2% |
| | | | | GPT-4o-mini | 0.0% | 3.3% | 12.5% | 4.5% | 14.6% | 0.0% | 5.3% |
| | | | | GPT-4o | **2.3%** | 17.3% | 20.8% | 4.5% | 9.8% | 0.0% | 9.3% |
| | | | | GPT-4V-1106 | **2.3%** | 13.7% | 16.7% | 13.6% | 19.3% | **8.3%** | 11.3% |
| ✓ | OneOCR | | Proprietary models | Phi3.5-V | 0.0% | 8.0% | 0.0% | 0.0% | 0.7% | 0.0% | 1.7% |
| | | | | Phi3-V | 0.0% | 6.7% | 8.3% | 4.5% | 5.0% | 0.0% | 4.0% |
| | | | | GPT-4o-mini | 0.0% | 7.3% | 20.8% | 8.3% | 9.8% | 0.0% | 7.3% |
| | | | | GPT-4o | 0.0% | 20.0% | 29.2% | 9.1% | 25.3% | 0.0% | 13.3% |
| | | | | GPT-4V-1106 | 0.0% | 26.3% | 16.7% | 17.4% | 19.3% | 0.0% | 13.1% |
| ✗ | | Omniparser | | Phi3.5-V | 0.0% | 6.9% | 8.8% | 0.0% | 0.5% | 0.0% | 2.8% |
| | | | | Phi3-V | 0.0% | 0.0% | 8.6% | 0.0% | 5.0% | 0.0% | 2.0% |
| | | | | o1 | 0.0% | 13.8% | 29.2% | 25.0% | **33.3%** | 8.3% | 16.3% |
| | | | | GPT-4o-mini | 0.0% | 0.0% | 12.5% | 0.0% | 5.3% | 0.0% | 2.7% |
| | | | | GPT-4o | 0.0% | 6.7% | 30.3% | 4.3% | 15.3% | **8.3%** | 9.4% |
| | | | | GPT-4V-1106 | **2.3%** | 23.6% | 20.8% | 8.3% | 20.0% | 0.0% | 12.5% |
| ✓ | | Omniparser | | Phi3.5-V | 0.0% | 6.7% | 8.3% | 0.0% | 0.3% | 0.0% | 2.6% |
| | | | | Phi3-V | 0.0% | 6.9% | 8.3% | 0.0% | 6.2% | 0.0% | 3.5% |
| | | | | o1 | 0.0% | 20.2% | 20.8% | **37.5%** | 19.0% | 0.0% | 15.6% |
| | | | | GPT-4o-mini | 0.0% | 14.9% | 8.3% | 0.0% | 0.0% | 0.0% | 4.2% |
| | | | | GPT-4o | 0.0% | 13.7% | 29.2% | 0.0% | 10.3% | 0.0% | 8.6% |
| | | | | GPT-4V-1106 | 0.0% | **27.3%** | **33.3%** | 27.3% | 30.3% | **8.3%** | **19.5%** |
| | | | | **Human performance** | 75.8% | 76.7% | 83.3% | 68.4% | 42.8% | 91.7% | 74.5% |

# A tale in 3 plots:



Benchmark score %

2.2x

Omniparser + UIA

GPT-4o   o1   GPT4-V   Human

GPT4-V >> GPT4-o
(for agents)

Benchmark score %

1.5x

Pixel-only   Pixel + UIA

GPT4-V

UIA / DOM boosts
pixel models

Benchmark Compute Time vs Benchmark Score (with o1)

21x

Benchmark Score

Screen Parsing Model
Open-source
Proprietary
Omniparser
o1

Benchmark Compute Time

Model speed matters
(even longer without
our parallelization)

# Where do agents fail?



**Task:** *Convert all uppercase text to lowercase*

**Step 3:** computer.mouse.move_id(id=5) # Move to the 'Format' menu option
computer.mouse.single_click() # Open the 'Format' menu

Screen parsing errors



**Task:** *Open Paint and draw a red circle*

**Step 7:** computer.mouse.move_id(id=59) # Move to the red color on the palette
computer.mouse.single_click() # Click to select the red color

Visual-language misalignment

# Thank you!

Feel free to reach out to dz1158@nyu.edu

Hope to see you at our poster at ICML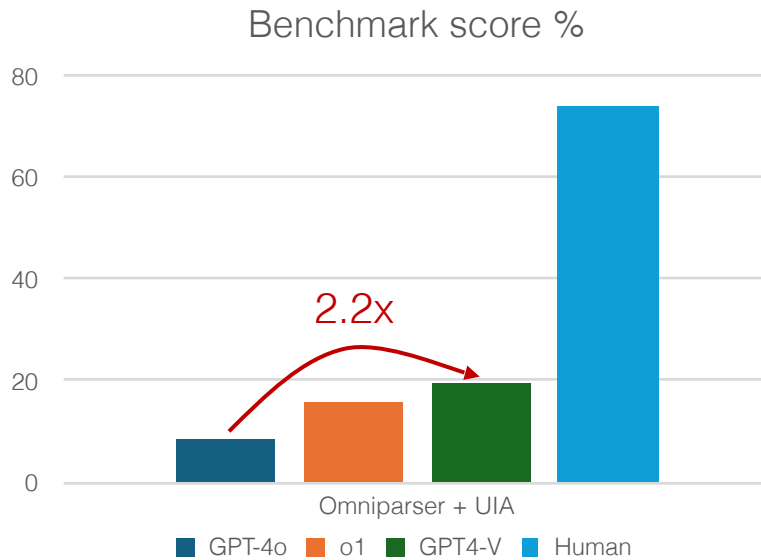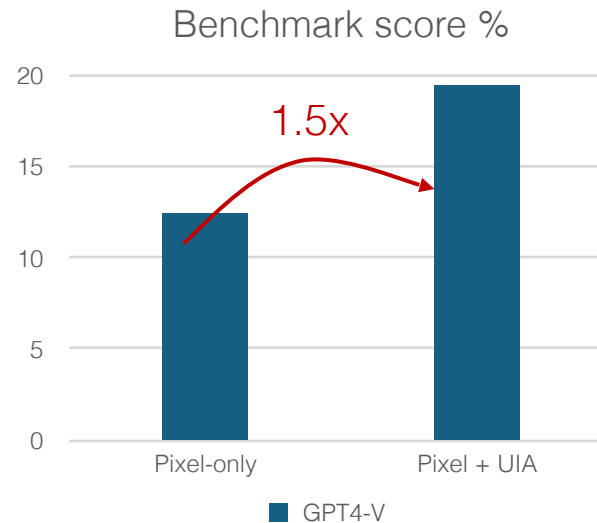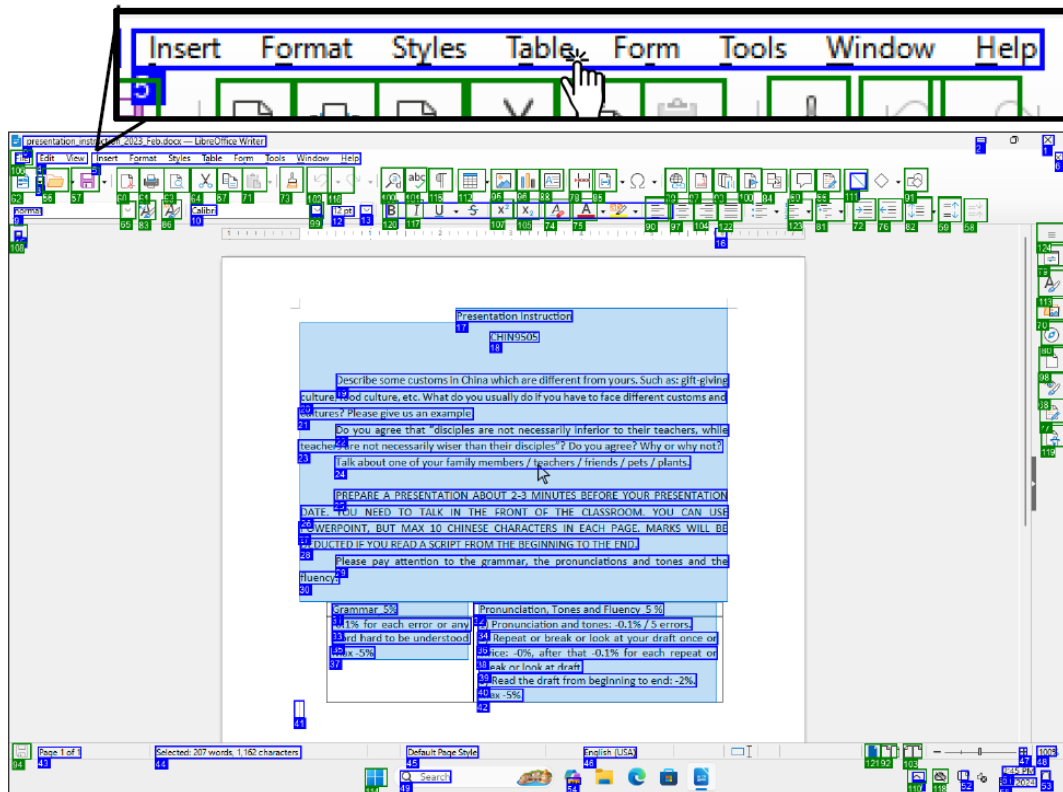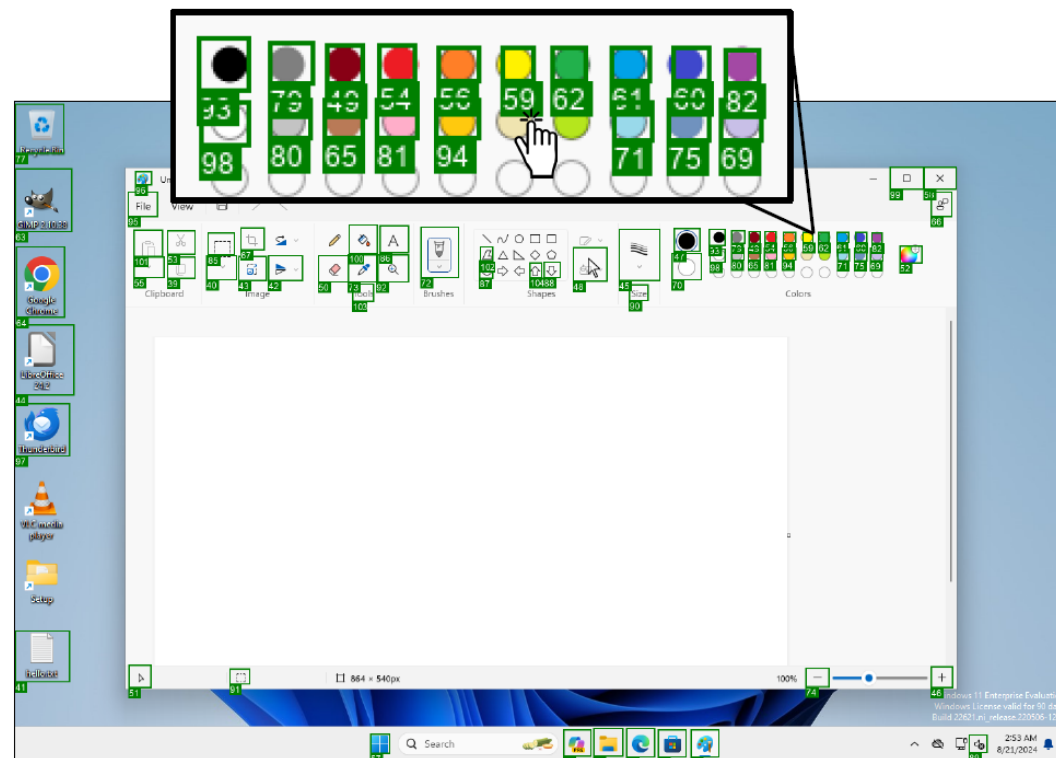